
Creating an Automated Classifier to Approximate the By-eye Selection of HAT Exoplanet Candidates

Will Coulton

Department of Physics
Princeton University
wcoulton@princeton.edu

Joshua Wallace

Department of Astrophysical Sciences
Princeton University
joshua.jw@princeton.edu

Abstract

A proposal to use machine learning techniques to aid in the binary classification of data from the Hungarian Automated Telescope (HAT) arrays. The HAT arrays are designed to detect periodic dimmings of stars, potentially caused by orbiting planets around those stars. There are many signals, both physical and spurious, that can approximate planet signals, so it is important to weed out as many false positives as possible. After an automated pipeline removes probable false planets, a final manual by-eye selection is made to create the final set of planet candidates. We propose to train classifier models to approximate this by-eye selection with the hope of finding a sufficiently well-performing model that can stand as a proxy for or even entirely replace the manual portion of the selection pipeline.

1 Introduction

Since the discovery of the first extra-solar planet (exoplanet) in the early 90's [6], nearly 3000 exoplanets have been discovered [3], with potentially thousands more exoplanets remaining unconfirmed or as yet undiscovered in existent data. This explosion in exoplanet discovery has been fueled by the funding and construction of many exoplanet detection surveys. The majority of exoplanets have been found by the space-based *Kepler* telescope and its primary and "K2" surveys [2], [5]. Ground-based surveys, while not as sensitive to smaller planets as *Kepler* (due to atmospheric distortions and other problems associated with observing from the ground), are sufficiently sensitive to discover Jupiter-sized planets on short-period orbits. (Such planets are called "hot Jupiters" because they are the same size as Jupiter but are sufficiently close in to the stars they orbit (much closer than Mercury is to our sun) that the star heats them up to several thousands of degrees Celsius.) Since ground-based surveys are much cheaper to run per area of sky monitored than spaced-based surveys, large-scale ground-based surveys have discovered more hot Jupiters than space-based surveys.

Of the few hundred hot Jupiters that have been found to date, the plurality (~ 100) have been found by Princeton's own Hungarian Automated Telescope (HAT) collaboration, headed by Prof. Gaspar Bakos of the Department of Astrophysical Sciences. This collaboration runs two surveys, both of which use telescopes that are no bigger than large camera lenses. The longest-running of the two surveys is HAT-Net, a collection of five telescopes in Arizona and two telescopes in Hawaii [4]. The other survey is HAT-South, which consists three locations with eight telescopes each: Chile, Namibia, and Australia [1].

Planets are far too faint compared to their host stars to be able to detect directly. Instead, what the HAT arrays do is monitor the brightness of stars as a function of time. When a planet crosses in front the star it orbits, the amount of light detected from that star is decreased. When a star is found to have periodic dips in brightness, it is labelled as a potentially planet-hosting star. Various cuts are made to help ensure that the signal is neither spurious nor caused by another astrophysical source other than a planet. Most of these cuts are automated, but the final step in the HAT pipeline is a

054 manual, by-eye examination of the signal. Planet candidates that pass this step then are sent for
055 follow up (and hopefully confirmation!) at bigger telescopes.

056
057 Other than this manual step, the vetting of possible planet candidates is completely automatic. This
058 single manual step is quite labor intensive, and prevents a fully automatic characterization of planet
059 detection efficiency, which is necessary for the calculation of statistics related to the occurrence rate
060 of these planets. Thus, it would be nice to have an automated approximation to the by-eye work that
061 has occurred. This is the purpose of this project.

062 **2 Related Work**

063 **3 Data**

064
065 The data consists of ~ 1000 examples of potential planet candidates that were decided, by eye, to
066 be included as planet candidates, and $\sim 30,000$ examples of potential planet candidates that were
067 decided to not be regarded as planet candidates. There are ~ 90 features, all continuous, for each
068 potential planet candidate. These features include parameters and statistics related to the fitting of
069 the brightness dips, parameters related to the measured planet size and orbital period, and properties
070 of the host star (how bright it is in different colors, etc.). The data were extracted from the proprietary
071 HAT database using a script provided by Joel Hartman, Research Scientist at Princeton University.
072 The database itself has not reference currently.

073 **4 Methods**

074
075 Our first approach to this problem will be to investigate the simple binary classifiers that we have
076 covered in class so far. These include Naive Bayes's, random forests, support vector machine (SVM),
077 and logistic regression. With logistic regression will investigate two types of regularization using l_1
078 and l_2 penalties. We will also explore a set of feature selection methods such as mutual information.

079
080 Whilst our main goal is to predict which objects would be manually selected, we are also interested
081 in running a set of unsupervised learning algorithms to view any hidden structure in the data set.
082 These methods will include K-means and gaussian mixture models. This will be useful for poten-
083 tially identifying new approaches to automating this process as well as potentially probing some
084 new relationships.

085 **5 Results**

086 **6 Discussion and Conclusion**

087 **Acknowledgments**

088
089 We are grateful to Joel Hartman for providing a script to extract the necessary data from the HAT
090 database. We are also grateful to the entire HAT team (PI: Gaspar Bakos) for their many dedicated
091 years of constant observations and the huge pile of astronomical data they've collected.

092 **References**

- 093
094
095
096
097
098
099
100
101 [1] G Bakos, C Afonso, T Henning, A Jordán, M Holman, RW Noyes, PD Sackett, D Sasselov,
102 Gábor Kovács, Z Csubry, et al. Hat-south: a global network of southern hemisphere automated
103 telescopes to detect transiting exoplanets. *Proceedings of the International Astronomical Union*,
104 4(S253):354–357, 2008.
105 [2] William J Borucki, David Koch, Gibor Basri, Natalie Batalha, Timothy Brown, Douglas Cald-
106 well, John Caldwell, Jørgen Christensen-Dalsgaard, William D Cochran, Edna DeVore, et al.
107 Kepler planet-detection mission: introduction and first results. *Science*, 327(5968):977–980,
2010.

108 [3] exoplanets.org. Exoplanets.org main page. <http://exoplanets.org/>. Accessed: 2017-
109 04-21.

110 [4] JD Hartman, G Bakos, KZ Stanek, and RW Noyes. Hatnet variability survey in the high stel-
111 lar density kepler field with millimagitude image subtraction photometry. *The Astronomical*
112 *Journal*, 128(4):1761, 2004.

113 [5] Steve B Howell, Charlie Sobeck, Michael Haas, Martin Still, Thomas Barclay, Fergal Mul-
114 lally, John Troeltzsch, Suzanne Aigrain, Stephen T Bryson, Doug Caldwell, et al. The k2 mis-
115 sion: Characterization and early results. *Publications of the Astronomical Society of the Pacific*,
116 126(938):398, 2014.

117 [6] A. Wolszczan and D. A. Frail. A planetary system around the millisecond pulsar PSR1257 +
118 12. *Nature*, 355:145–147, January 1992.

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161