

Exoplanet Classification with Data Mining

Israel dos Santos^{a,b,1}, Adriana Valio^{a,c}, Nizam Omar^{a,b} and A. H. F. Guimarães^c

^aElectrical Engineering and Computation Graduate Program (PPGEEC),
Mackenzie Presbyterian University, São Paulo, SP, Brazil

^bFaculty of Computation and Informatics (FCI),
Mackenzie Presbyterian University, São Paulo, SP, Brazil

^cCenter for Radio Astronomy and Astrophysics Mackenzie, Engineering School,
Mackenzie Presbyterian University, São Paulo, SP, Brazil

Received on November 21, 2016 / accepted on December 30, 2016

Abstract

This article presents a case study on Data Mining in Astronomy to classify exoplanets and help the identification of terrestrial planets within the habitable zone of their host stars. To achieve this objective, Data Mining techniques were applied on the planets listed on the Exoplanet Orbit Database (EOD). These planets can be grouped into at least 6 types according to their mass and radius. Two classification algorithms were used: Decision Tree and k-Nearest Neighbor. The latter method, or clustering technique, classifies the multidimensional array of attributes based on resource creation using the mass and the radius of the planet. This in turn creates the special attribute for the classification of Planet Type. In this analysis, both the preliminary validation and the results obtained by the decision tree indicated that there is no correlation between stars and planetary parameters. The classification model returned the results for two dimensions with the following classes: Mercurian, Subterran, Terran, Superterran, Neptunian and Jovian, however did not present results for Mercurian. *As a conclusion, both kNN and Decision Tree algorithms can be used in an ambient with several processing cores, however the kNN algorithm may be clustered. If a high correlation index and low entropy of the data is desired, the kNN may not converge, whereas the Decision Tree presents an instant result based on conditional parameters.*

Keywords: astronomy, classification, data mining, exoplanets.

1. Introduction

To date, about 1,500 exoplanets have been identified, orbiting other stars, and 3,700 planet candidates from the Kepler satellite [1]. Initially,

¹E-mail Corresponding Author: israel.santos@mackenzie.br

39 the vast majority of discovered planets were Hot Jupiters, that is, gaseous
 40 giants very close to their host star. However, the huge increase in the num-
 41 ber of exoplanets in the last years occurred mainly with the Kepler mission,
 42 a space probe able to identify even smaller bodies, aimed at finding Earth-
 43 like planets within the habitable zone (HZ) of their host star. The habitable
 44 zone of a stars is loosely defined as the distance from the star where the
 45 temperature on the surface of a solid planet is such that liquid water can
 46 exist.

47 Usually, exoplanets are divided into classes according to their radius and
 48 mass. Presently, they are separated into the following types: Asteroid, Mer-
 49 curian, Subterran, Terran, Superterran, Neptunian and Jovian, as presented
 50 in Table 1.

Table 1: Mass classification for solar and extrasolar planets

Planet Type	Mass (Earth Units)	Radius (Earth Units)
Asteroids	0 - 0.00001	0 - 0.03
Mercurians	0.00001 - 0.1	0.03 - 0.7
SubEarth	0.1 - 0.5	0.5 - 1.2
Earth	0.5 - 2	0.8 - 1.9
SuperEarth	2 - 10	1.3 - 3.3
Neptunians	10 - 50	2.1 - 5.7
Jovians	50 - 5000	3.5 - 27

Source: Planetary Habitability Laboratory (PHL)
 University of Puerto Rico at Arecibo

51 Asteroids are small irregular bodies that are not able to maintain a stable
 52 environment. Mercurians are able to maintain a significant atmosphere only
 53 in cold distant regions, as found on Saturn's largest natural satellite, Titan,
 54 the second largest in the entire Solar System and nearly one-and-a-half times
 55 the size of the moon. It is the only natural satellite known to have a dense
 56 atmosphere, being even denser than Earth's, and the only object to have
 57 clear evidence of liquid bodies on its surface, beyond our planet.

58 Subterran planets are able to maintain a significant atmosphere beyond
 59 the outer edges of the habitable zone, being identical to Mars. Terrans can
 60 retain a significant atmosphere with liquid water within the habitable zone,
 61 identical to Earth. Superterrans are able to sustain dense environments with
 62 liquid water within the habitable zone. Netunians, on the other hand, may
 63 have dense atmospheres in the hot zone, and finally, Jovians may have super
 64 dense environments in the hot zone.

65 In this context, this study aims to classify all the identified exoplanets,

66 according to the definition presented by the Laboratory of Habitability Plan-
67 ets [2] of the University of Puerto Rico at Arecibo (Table 1). This may help
68 research on the identification of terrestrial planets that are located within
69 the HZ of their star where life may originate and develop [3].

70

71 **2. The quest for exoplanets**

72 CoRoT (Convection, Rotation and Planetary Transits) [4] was the first
73 space mission dedicated to exoplanetary research. Launched in December
74 2006, the mission had a nominal life of 2.5 years, later extended until March
75 31, 2013, a period that studied the interior of the stars and searched for exo-
76 planets. CoRoT was designed by the CNES (National Center for Space Stud-
77 ies) in partnership with ESA (European Space Agency) and AEB (Brazilian
78 Space Agency).

79 The Kepler [5] space mission is a NASA-designed (National Aeronautics
80 and Space Administration) observatory to search for extrasolar planets. To
81 this end, the probe observed one hundred thousand stars for a period of
82 four years. The main goal of the Kepler was to search for telluric exoplanets
83 similar to Earth, that is, those between half to twice the size of the Earth.
84 Especially those terrestrial planets within the habitable zone of their host
85 stars, where liquid water may exist in their surface.

86 The method for planet identification used by both the CoRoT and the
87 Kepler satellites is the detection of small decreases in the light of a star
88 caused by the eclipse or transit of a planet. Planetary transit is the passage
89 of a planet in front of the stellar disk during its orbit, when a dimming of a
90 fraction of a percent in the light of the star can be measured.

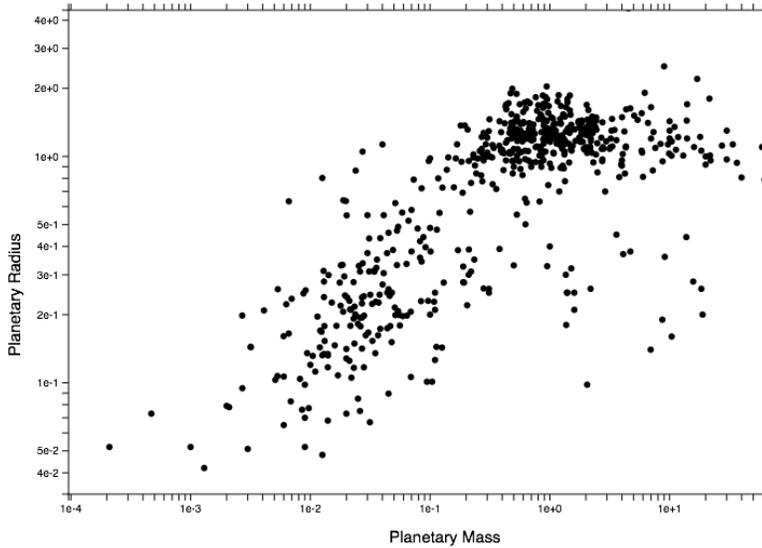
91 The data used in this study was extracted from the *Exoplanet Orbit*
92 *Database* (EOD) [1]. The content of the database is a compilation of data
93 collected using radial velocity and transit methods for planet detection, in-
94 cluding the planets discovered by the CoRoT [4] and Kepler [5] missions. The
95 masses and radius of stars are obtained from precise spectroscopic observa-
96 tions. Moreover, accurate masses of planets are also determined from radial
97 velocity measurements from spectroscopic observation of stars, whereas the
98 radius of the planet can only be measured with transit observations. The
99 amount of dimming in the light of a star during a transit is proportional to
100 the ratio of the planet to stellar area (or radius).

Table 2: Quantity of records on the database

Classification	Quantity
Confirmed Planets (EDO)	1491
Kepler Planet Candidates	3704
Total of Planets	5195

Source:Exoplanet Orbit Database [1]

101 The main classification attributes, that is mass and radius of the planets,
 102 are plotted on Figure 1. The data represent planets of different densities,
 103 thus explaining the scatter in the figure. Table 2 lists the number of planets
 104 in the Exoplanet Orbit Database, whereas

Figure 1: Planet radius (R_{Jup}) X Planet mass (M_{Jup})

105 Table 3 displays the metadata that characterizes a Star. The text type
 106 is defined for the *STAR* attribute and the numeric type is defined for the
 107 attributes *MSTAR*, *RSTAR*, *TEFF*, *RHOSTAR*, *LOGG*, *KP*. The planet
 108 metadata is shown in Table 4. The numeric type is defined for the attributes
 109 *ID*, *PER*, *ECC*, *K*, *T0*, *MSINI*, *A*, *MASS*, *R*. In this work we are manly
 110 interested in the *MASS* and *R* parameters.

Table 3: Star Metadata

ID	Name	Unit	Description
STAR	Star Name		Standard name for star
MSTAR	Mass of Star	Solar Mass	Estimated mass of the star, usually based on the association of effective temperature of the star and its luminosity compared to stellar models
RSTAR	Radius of Star	Solar Ray	Estimated Star Ray
TEFF	Teff	Kelvin	Effective Star Temperature
RHOSTA	Density of star	Grass / Centimeter	Density of star as measured from transit photometry and radial velocity information
LOGG	$\log_{10}(g)$		\log_{10} of gravity on the surface of the star (unit cgs)
KP	KP		magnitude of the Kepler band

Source: *Exoplanet Orbit Database*[1]

Table 4: Planet and Orbit metadata

ID	Name	Unit	Description
PER	Orbital Period	Days	The Orbital Period of the Planet
ECC	Orbital Eccentricity		The eccentricity of the orbit, on the usual scale from 0 to 1, where 0 is circular and 1 is extremely flattened
K	Velocity Semi-amplitude	m/s	The half-amplitude of Doppler variation (half of the variation of the peak-to-peak radial velocity).
T0	Time of Periastron	JD	The time of a planet periastron passage in JD
MSINI	Msin (i)	Jupiter Mass	Minimum Planet Mass. True masses are usually higher by about 15 % due to the geometric inclination effects
A	Semi-Major Axis	AU	Semi-major axis
MASS	Planet Mass	Jupiter Mass	Planet mass
R	Planet Radius	Jupiter Radius	Planet radius

Source: *Exoplanet Orbit Database* [1]

111

112

3. Data Mining Concepts

113

114

According to Tan, Steinbach and Kumar [6] the input data for the classification are a set of records, each record being known as an instance and

115 characterized by a double (x, y) , where x is the set of attributes and y the
 116 special attribute. Classification is the task of learning a target function f
 117 that maps each set of attributes x to one of the y labels of pre-determined
 118 classes [6]. Also, according these authors a classification technique (or clas-
 119 sifier) is a systematic approach to constructing classification models from a
 120 set of input data. Among different algorithms, the Hunt algorithm is the
 121 basis for many induction algorithms for existing decision trees, including
 122 ID3, C4.5 and CART.

123 According to the Hunt algorithm [6], a decision tree grows recursively by
 124 partitioning the training records into successive more pure subsets. Since
 125 D_t is the set of training records that are associated with the node t and
 126 $y = y_1, y_2, \dots, y_c$ are the labels of the classes, we have:

- 127 • If all records in D_t belong to the same class y'_t then t is a leaf node
 128 labeled y_t .
- 129 • If D_t contains records that belong to more than one class, an attribute
 130 test condition is selected to partition the records into smaller sub-
 131 sets. A child node is created for each test condition result, and the
 132 D_t records are distributed to the children based on the results. The
 133 algorithm is then applied recursively to each child node.

134 There are different approaches to determine the best way to divide the
 135 records, the metrics being based on the degree of impurity of the child
 136 nodes. As examples are the Entropy (Equation 1), Gini Rate (Equation 2),
 137 and Classification Error (Equation 3), among others:

$$\text{Entropy}(s) = - \sum_{i=0}^{c-1} p(i|s) \log_2 p(i|s), \quad (1)$$

$$\text{Gini}(s) = 1 - \sum_{i=0}^{c-1} [p(i|s)]^2, \quad (2)$$

$$\text{Classification Error}(s) = \max_i [p(i|s)]. \quad (3)$$

138 The *k-Nearest Neighbors* (kNN) algorithm [7] stores the training data and
 139 when a new object is submitted for classification, the algorithm searches for
 140 the closest k records (Euclidean distance) of this new record. Thus, the new
 141 record is ranked in the most common class among all the k closest records.

142 The algorithm applies the concept of centroids. Therefore, for a given a
 143 data set, the algorithm randomly selects k records, each representing a clus-
 144 ter. For each remaining record the similarity between the analyzed record

145 and the center of each grouping is calculated. The object is inserted in the
146 cluster with the shortest Euclidean distance, that is, greater similarity and,
147 with each new element inserted, the centroid is recalculated. Thus, one has:

- 148 1. Provide values for centroids
- 149 2. Generate a distance matrix between each point and the centroids
- 150 3. Place each point in the classes according to its distance from the class
151 centroid
- 152 4. Calculate the new centroids for each class
- 153 5. Return to step 2 and repeat until convergence

154 It should be emphasized that different variations implement optimiza-
155 tions for the choice of the k value, measures of dissimilarity and strategies
156 for the calculation of the cluster center. One variation of this is the k-Means
157 [8] algorithm, that uses the mode to calculate the centroids.

158
159

4. Planet classification

160 The clustering technique was used to generate a multidimensional array
161 of *NAME* attributes, the name of the candidate or confirmed planet, *R*, the
162 radius of the candidate or confirmed planet, *MASS*, the mass of candidate
163 or planet, *MSTAR*, the mass of the star, and finally the *RSTAR*, radius of
164 the star. Unfortunately, the Exoplanet Orbit Database [1] has missing data
165 for many records. Thus, for data preprocessing, the data cleaning technique
166 favors the option of excluding such records to avoid inconsistencies in the
167 elaboration of the construction process for later analysis. The technique of
168 resource creation used the mass and the radius of the planet to generate
169 the special attribute for the classification of Planet Type, according to the
170 characteristics presented in Table 1.

171 The Rapid-Miner[®] tool 5.3.000 was used to process the data, which was
172 selected from the *Exoplanet Orbit Database* [1], converted to a Microsoft
173 Excel worksheet[®], downloaded to the local machine and then submitted
174 to pre-processing. When loading the worksheet to RapidMiner[®], heav-
175 ily typed attributes are used. In this step, the data types of the *NAME*
176 attributes were defined for *ID*, the *MASS*, *R*, *MSTAR* and *RSTAR* for *NU-*
177 *MERIC*, and finally the *Class* attribute, which was generated by the resource
178 creation technique, has been defined as *LABEL*.

179 The parametrization criteria of the algorithms were empirically tested
180 for the software, algorithm and specificities that make up the Data Mining

181 technique. Finally, the preliminary validation of the data was performed
 182 based on the metric presented in Table 1.

183 We applied the two algorithms described in Section 3 to classify the plan-
 184 ets into Planet Types. First, the Decision Tree algorithm was used. Figure 2
 185 presents the result of sorting by the decision tree according to *MASS*, using
 186 the Gini index. Sorting by the decision tree using the information gain yields
 187 the same results. Therefore, no change in the structure of the trees with the
 188 index exchange occurs, basically because the initial parameters were not al-
 189 tered. The decision tree algorithm applied to the planet radius attribute *R*
 190 yields the same results, and thus are not shown here.

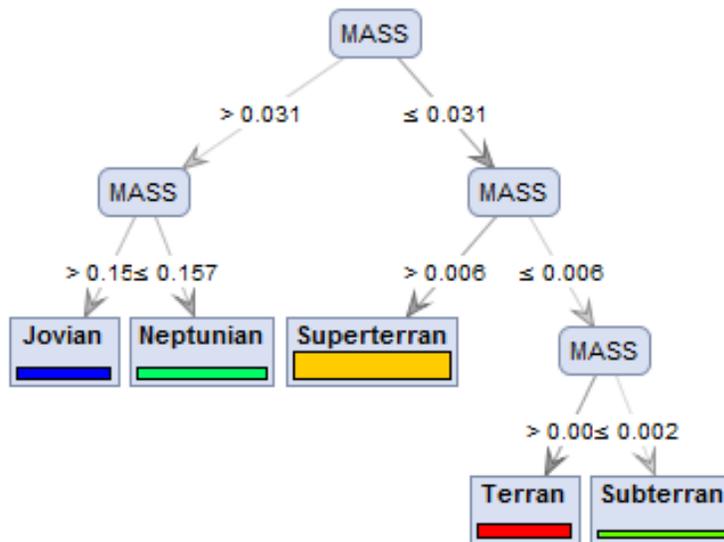


Figure 2: Decision tree for *MASS* attribute.

191 Next, the second algorithm, kNN, was applied to the planet database
 192 using both the planet mass and radius attributes. Table 5 presents the
 193 first clustering attempt using the *k-Means* algorithm without the special
 194 attribute, where it can be seen that the *k* variation generated groups in
 195 border areas.

Table 5: k-Means with 6 Clusters

Cluster	items
0	544
1	15
2	767
3	207
4	13
5	1035
6	181
Total number of items	2762

196 Considering the initial planet types given in Table 1, the algorithm
 197 matches the clustering data of Table 5 into the predicted planet classes, that
 198 yields the accuracy listed on Table 6. The classification model returned the
 199 result for two dimensions with the following classes: Mercurian, Subterran,
 200 Terran, Superterran, Neptunian and Jovian, as summarised in Table 6. As
 201 can be seen from the values presented on the Table, the kNN algorithm show
 202 accuracies above 80% in most of the cases.

Table 6: K-Means Performance for 6 Clusters

Prediction Class	Accuracy						Precision%
	Merc.	SubT.	Terran	SuperT.	Nept.	Jovian	
Mer.	0	0	0	0	0	0	0.00%
SubTer.	2	44	15	0	0	0	72.13%
Ter	0	31	384	47	0	0	83.12%
SupTer.	0	1	60	1434	44	0	93.18%
Nep.	0	0	0	23	323	7	91.50%
Jov.	0	0	0	0	4	343	98.85%
%	0.00	57.89	83.66	95.35	87.06	98.00	

203 For internal validation of the k-NN algorithm, the number of initial clus-
 204 ters was increased to 12. Doubling the number of clusters to 12 does not
 205 alter significantly the results, thus validating the method.

206

207

5. Conclusions

208

209

210

211

212

213

Two classification algorithms were used to classify almost three thousand planets in the Exoplanet Orbit Database (EOD), according the their mass and radius (Table 1). The initial classification performed in this work, will help identify terrestrial planets within the HZ of their host star, since life as we know it needs liquid water on the surface of a terrestrial planet to develop.

214 Both algorithms, Decision Tree and kNN, have contributed in some way
215 to the validation or clustering of data. However, the performance between
216 the algorithms tends to become computationally costly depending on the
217 initial configurations. The decision tree algorithm required a preprocessing
218 of the data to ensure that there was no absence of values, in the definition
219 of the choice of a special attribute for the classification of the samples,
220 whereas the kNN proved to be an algorithm that has supervised learning
221 based on the instances present in the data mass. Thus, it is not possible
222 to choose a special class attribute in this algorithm because the comparison
223 is performed in relation to all other known copies in the knowledge base,
224 and this comparison can consume a lot of computing time. In particular,
225 the Rapid-Miner[®] tool used in the present study allowed a more accurate
226 analysis of the data through different algorithms.

227 The current paradigm of science consists in allying the use of compu-
228 tational techniques with the other areas of knowledge, such as Astronomy.
229 More specifically, the results obtained from the planet classification scheme
230 presented here constrained the number of terrestrial planets candidates, thus
231 making it possible to focus on the set of life bearing planets.

232
233 **Acknowledgments.** Israel dos Santos acknowledges financial support
234 from MackPesquisa, A.V. thanks FAPESP for partial financial support, and
235 A.H.F.G. acknowledges a fellowship from CAPES.

236 References

- 237 [1] Exoplanet Data Explorer. Available: <http://www.exoplanets.org>
- 238 [2] Wright, J. T., Fakhouri, O., Marcy, G. W., Han, E., Feng, Y., Johnson,
239 J. A., ... & Piskunov, N. (2011). The exoplanet orbit database. Publi-
240 cations of the Astronomical Society of the Pacific, 123(902), 412-422.
241 DOI: 10.1086/659427
- 242 [3] Holman, M. J., & Wiegert, P. A. (1999). Long-term stability of planets
243 in binary systems. The Astronomical Journal, 117(1), 621.
244 DOI: 10.1086/300695
- 245 [4] Astronomy mission. From stars to habitable planets. Available:
246 <http://smc.cnes.fr/COROT/>
- 247 [5] Kepler. A search for Habitable Planets. Available:
248 <http://kepler.nasa.gov/>

- 249 [6] Pang-Ning, T., Steinbach, M., & Kumar, V. (2006). Introduction to
250 data mining. In Library of congress (Vol. 74).
- 251 [7] McNames, J. (2001). A fast nearest-neighbor algorithm based on a prin-
252 cipal axis search tree. *IEEE Transactions on Pattern Analysis and Ma-*
253 *chine Intelligence*, 23(9), 964-976.
254 DOI:10.1109/34.955110
- 255 [8] Hartigan, J. A., & Wong, M. A. (1979). Algorithm AS 136: A k-means
256 clustering algorithm. *Journal of the Royal Statistical Society. Series C*
257 *(Applied Statistics)*, 28(1), 100-108.
258 DOI: 10.2307/2346830