# Identifying deep moonquake nests using machine learning model on single lunar station on the far side of the Moon

Josipa Majstorović<sup>1</sup>, Philippe Lognonné<sup>2</sup>, Taichi Kawamura<sup>3</sup>, and Mark Paul Panning<sup>4</sup>

<sup>1</sup>Institut de Physique du Globe de Paris

<sup>2</sup>Université Paris Cité, Institute de physique de globe de Paris, CNRS <sup>3</sup>Université Paris Cité, Institut de physique du globe de Paris, CNRS <sup>4</sup>Jet Propulsion Laboratory, California Institute of Technology

December 8, 2023

## Abstract

One of the future NASA space program includes the Farside Seismic Suite (FSS) payload, a single station with two seismometers, on the far side of the Moon. During FSS operations, the processing of the data will provide us with new insight into the Moon's seismic activity. One of Apollo mission finding is the existence of deep moonquakes (DMQ), and the nature of their temporal occurrence patterns as well as the spatially clustering. It has been shown that DMQs reside in about 300 source regions. In this paper we tackle how we can associate new events with these source regions using the single station data. We propose a machine learning model that is trained to differentiate between DMQ nests using only the lunar orbital parameters related to DMQ time occurrences. We show that ML models perform well (with an accuracy >70%) when they are trained to classify less than 4 nests.

# Identifying deep moonquake nests using machine learning model on single lunar station on the far side of the Moon

# Josipa Majstorović<sup>1</sup>, Philippe Lognonné<sup>1</sup>, Taichi Kawamura<sup>1</sup>, Mark P. Panning<sup>2</sup>

 $^1$ Université Paris Cité, Institut de physique du globe de Paris, CNRS, Paris, France $^2$ Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA

# Key Points:

1

2

3

4

5

6

8

9	• As a part of the future space mission NASA will deploy a new seismic station to
10	Schrödinger Basin on the far side of the Moon.
11	• We propose a machine learning model trained to classify deep moonquakes using
12	the lunar orbital parameter.
13	• The models perform with accuracy greater than 70% when trained to classify com
14	binations of four or fewer nests.

 $Corresponding \ author: \ Josipa \ Majstorović, \ \texttt{josipa.majstorovic} @\texttt{protonmail.com}$ 

#### 15 Abstract

One of the future NASA space program includes the Farside Seismic Suite (FSS) pay-16 load, a single station with two seismometers, on the far side of the Moon. During FSS 17 operations, the processing of the data will provide us with new insight into the Moon's 18 seismic activity. One of Apollo mission finding is the existence of deep moonquakes (DMQ), 19 and the nature of their temporal occurrence patterns as well as the spatially clustering. 20 It has been shown that DMQs reside in about 300 source regions. In this paper we tackle 21 how we can associate new events with these source regions using the single station data. 22 We propose a machine learning model that is trained to differentiate between DMQ nests 23 using only the lunar orbital parameters related to DMQ time occurrences. We show that 24 ML models perform well (with an accuracy > 70%) when they are trained to classify 25 less than 4 nests. 26

# 27 Plain Language Summary

The future space missions will provide us with various new lunar data, one of which 28 will be ground vibration measurement. The studies of these measurements from the Apollo 29 era in 70s, showed that Moon can host various events. The most intriguing ones are deep 30 moonquakes (DMQs), which are events associated with the displacement deep within the 31 lunar interior. It has been shown that DMQs occur in the specific locations, which are 32 called nests, and that their temporal occurrence is related to the monthly motion of the 33 Moon around the Earth. In this paper we tackle how we can associate new events from 34 only one station located on the far side of the Moon with these known locations of DMQs. 35 We propose a machine learning model that is trained to classify DMQ nests, only using 36 the information about their temporal occurrences, e.g. time of the event, described in 37 terms of different lunar events. We report that models are performing well (with an ac-38 (uracy > 70%) when they are trained to classify 4 or fewer nests. This gives us a good 39 first approximation about the nest identification. 40

#### 41 **1** Introduction

We are at the beginning of a golden age of lunar exploration as many nations, to-42 gether with private companies, are establishing numerous efforts to obtain new scien-43 tific measurement of the Moon (Weber et al., 2021; Pickrell, 2022; Kawamura et al., 2022). 44 In light of this, NASA established the Artemis program which should land a crewed mis-45 sion at the lunar south pole (Witze, 2022). This would be the first attempt of a crewed 46 landing after the successful Apollo missions in 1970's. Before Artemis missions land on 47 the Moon, NASA has also established the Commercial Lunar Payload Services (CLPS) 48 program to land scientific payloads on the Moon using commercial landers. The Farside 49 Seismic Suite (FSS) is one of the selected payloads, and it will deliver two seismometers 50 to Schrödinger Basin on the far side of the Moon (Panning et al., 2021; Standley et al., 51 2023; Cutler et al., 2023): one vertical Very BroadBand seismometer, and Short Period 52 sensor, both spare or derived from the SEIS experiment sensors (Lognonné et al., 2019, 53 2020) from the InSight mission to Mars (Banerdt et al., 2020). 54

The Apollo missions showed the importance of deploying sensors on the surface of 55 the Moon, since a great deal of our knowledge about the Moon comes from the analy-56 sis of data acquired during the Apollo era (Lognonné & Johnson, 2015). Thus, analyz-57 ing ground motion measurements provided the community with the first constraints on 58 the lunar interior and the activity at the surface (Nakamura et al., 1982a, 1982b; Khan 59 et al., 2000; Khan & Mosegaard, 2002; Khan et al., 2014; Lognonné et al., 2003; Gagnepain-60 Beyneix et al., 2006; Weber et al., 2010; Garcia et al., 2011; Kawamura et al., 2017; Gar-61 cia et al., 2019; Nunn et al., 2020). It has also revealed that the Moon can host events 62 of various origins, such as shallow and deep moonquakes, meteoroid and artificial impacts 63 (Toksöz et al., 1974; Dainty et al., 1975; Lammlein, 1977a; Nakamura, 1983, 2003, 2005). 64

Today, we have more than twelve thousands events, out of which the deep moonquakes (DMQs) form the most numerous group (Nakamura et al., 1981; Nakamura, 2005).

DMQs are a distinctive group of seismic events that originate from depths between 67 700 and 1200 km, at high pressure and temperature conditions, where little brittle de-68 formation is expected. Due to very high waveform similarity between quakes, the DMQs 69 have been clustered into about 300 source regions or nests (Nakamura, 2003). This has 70 interpreted to be a consequence of DMQs occurring repeatedly at the fixed nests, which 71 are located mostly on the near side of the Moon. It has been shown that time occurrence 72 73 of the DMQs is correlated with the monthly motion of the Moon around the Earth. Thus, DMQ occurrences exhibit tidal periodicities and furthermore, the associated high strain 74 rates might explain brittle processes (Kawamura et al., 2017). However, the real causes 75 of their origins are yet to be discovered. There are two puzzling fact about their origin: 76 a) cyclic tidal stresses, caused by the monthly motion of Moon around Earth, are far less 77 than the estimated confining pressures where DMQs occurs (Cheng & Toksöz, 1978; Min-78 shull & Goulty, 1988a); b) do we need both, tectonic and ambient tidal stresses, to ex-79 plain their mechanical origin (Frohlich & Nakamura, 2009). 80

To better constrain the lunar interior and unravel the cause of DMQs, it is impor-81 tant to locate new events and associate them with the known nest locations from Apollo 82 with future lunar missions like FSS. These new observations will add, for each new nest, 83 a new differential  $t_s - t_p$  measurement constraining the deep interior with a different 84 epicentral distance. However, due to the mission requirements, it is extremely likely that, 85 at the beginning, we might have only one lunar station at the disposal. Therefore, in this 86 paper we study the problem of DMQ nest identification without using waveform infor-87 mation. This is due to the new location of the recording station, which will not match 88 existing Apollo-era waveform templates due to different propagation paths. We propose 89 a machine learning (ML) model that is trained to identify nests within the set of nests 90 of similar differential travel times. The main features used for the model training are re-91 lated to the fact that different nests respond differently to lunar cycle. 92

Very early studies have shown correlation between lunar transient events and po-93 sition of the Moon related to the Earth (Middlehurst, 1967; Cameron & Gilheany, 1967; 94 Moore, 1968). This further encouraged observations that some moonquakes occur with 95 periods that reflect Earth-Moon-Sun relationship (Ewing et al., 1971). Later, it has been 96 shown that the occurrence of DMQs are related to tidal stress cycles, and correlations 97 between DMQs occurrence times and lunar monthly tidal cycles have been indicated (Lammlein 98 et al., 1974; Toksöz et al., 1977; Lammlein, 1977b; Cheng & Toksöz, 1978; Minshull & 99 Goulty, 1988b). The lunar cycle can be explained with three lunar months: synodic, dra-100 conic, anomalistic. Synodic month is the period of lunar phases such as New Moon, First 101 Quarter, Full Moon, Last Quarter. Draconic month is the period between two nodes, as-102 cending or descending, where the nodes are points at which the Moon's orbit plane crosses 103 the ecliptic plane towards which it is inclined of about 5.14°. Anomalistic month is the 104 period between two extreme points, perigee and apogee, since the Moon's orbit approx-105 imates an ellipse rather than a circle. Earlier studies counted the number of events per 106 day as a function of time and found 0.5 and 1 month signals in the occurrence times re-107 lated to anomalistic and draconic period of 27 days (Lammlein et al., 1974; Lammlein, 108 1977b). The same studies also indicated 206-day and 6-year periods, related to the Sun's 109 perturbation on the lunar orbit and the relative precession of the perigee of the Moon's 110 orbit. Subsequently, many recent papers studied and confirmed tidal periodicities of DMQs 111 and more (Bulow et al., 2005, 2007; Bills et al., 2008; Frohlich & Nakamura, 2009; We-112 ber et al., 2009, 2010; Turner et al., 2022). 113

Based on the previous papers, it is clear that DMQ nests exhibit some clear temporal patterns in their occurrences, and that these are correlated with Moon-Earth system. Therefore, the open question is whether we can design features which would reflect these temporal patterns and to further use those features to study the nest identification with one lunar station. In this paper we tackle the question of defining optimal features and the machine learning model. The paper is organised as follows: first, we discuss data used in the analysis, the existing catalog of DMQ events. Second, we discuss the feature design. Third, we introduce a machine learning model. Fourth, we discuss successes and pitfalls of the machine learning model for nest identification applied to different combination of nests. We conclude how this study can offer some first estimates of the nest location in the future lunar missions.

#### 125 2 Data

We start with the existing catalog of lunar events (Nakamura et al., 1981), which 126 was updated in 2008 and modified in 2018 (Nunn et al., 2020). Catalog contains a list 127 of events (shallow and deep moonquakes, meteoroid and artificial impacts) with attributes 128 such as date and time of the event occurrence, signal envelope amplitude as measured 129 in mm on a standard plot, data availability per station, and the nest (source) classifi-130 cation for DMQs. It is important to note that the source classification is not an exact 131 location defined by latitude and longitude, rather a result from the waveform cross-correlation 132 and single-link cluster analysis (Nakamura, 2003). This analysis positively clustered around 133 7k DMQs into 77 nests, where the largest nest is associated with label A1. This nest con-134 tains 443 guakes and it is placed on the near side of the Moon. 135

136

# 2.1 Catalog processing: nest sets based on travel time information

Earlier studies published lunar interior models and location of DMQ nests in terms 137 of latitude and longitude by picking P and S travel times on the quake waveforms (Garcia 138 et al., 2019, review of these picks). Using lunar interior models and nest locations we can 139 define nests that are close by in distance if we consider only the  $t_{sp}$  travel time measure-140 ments. To do so we assume that a) our single lunar station is located on the far side in 141 Schrödinger Basin (FSS landing site at 71.378°S, 138.248°E), b) nests' latitudes and lon-142 gitudes from (Lognonné et al., 2003), c) calculated P- and S- wave travel times ( $t_p$  and 143  $t_s$ , respectively) using lunar velocity model between landing site and locations of DMQ 144 nests. By having  $t_s$  and  $t_p$  we can calculate  $t_{sp} = t_s - t_p$  for all nests and models shown 145 in Figure 1 (see Text S1 and Figure S1 for further explanation). Next, we count for each 146 nest how many there are with the similar  $t_{sp}$  travel time measurement assuming a pick-147 ing error of 5 seconds as shown in Figure 1A, consistent with the average picking error 148 in Lognonné et al. (2003). This count provides us with the different sets  $S_i$ , shown in Figure 1A, that contain nests  $N_j$  of similar travel times. In other words, if we are able 150 to measure  $t_s$  and  $t_p$  of the new lunar event with accuracy within 5 seconds, we are not 151 able to differentiate between nests that belong to different sets  $S_i$ . Therefore, to further 152 tackle the nest identification problem we proceed to associate each event with a com-153 bination of lunar orbital measurements. 154

155

#### 2.2 Feature selection based on the lunar orbital information

It has been shown that the DMQ temporal patterns in time occurrences are related 156 to different lunar cycles and that these patterns differ from nest to nest. Three lunar cy-157 cles are synodic, draconic, anomalistic, and they all have similar periods, but are marked 158 by different motions, either as the motion between two Full Moons phases, or two nodes, 159 or two apsis, respectively. One can list all the events when Moon is in the Full Moon (New 160 Moon) phase, passing through ascending (descending) node or perigee (apogee) by sim-161 ply looking at the Moon's ephemeris (Meeus, 1991). To make sure that we take into ac-162 count the temporal patterns, we design the main three features as a time difference be-163 tween the time of the quake in the nest and the time of the Moon's Full Moon, ascend-164 ing node and perigee, denoting it as  $\Delta t_{FullMoon}, \Delta t_{AscendingNode}, \Delta t_{Perigee}$ , respectively. 165 We can achieve the same effect by taking the other three time axis as referent one (New 166



Figure 1. Location study of the DMQ nests from the perspective of  $t_{sp} = t_s - t_p$  travel time measurements if we place station in the Schrödinger Basin and consider four different lunar models. A) Upper panel:  $t_{sp}$  travel time measurements for four lunar models from Garcia et al. (2011) (G11), Garcia et al. (2019) (issi2, issi3), (Khan et al., 2014) (K14) with nest labels; A) Lower panel: Sets  $S_i$  which represent nests with <u>similar</u> travel times if we consider a travel time error of 5 seconds. B) Lunar map with the nests locations where the color indicate the median  $t_{sp}$  for four lunar model.

Moon, descending node, and apogee). The next feature is related to the Moon position 167 within its orbit as in Frohlich and Nakamura (2009). The angle between the direction 168 of perigee and the current position of the body, as seen from the main focus of the el-169 lipse, is called the true anomaly, denoted further as  $\gamma$ . Further, as one of the feature we 170 also use the interval time between two quakes in the nest, noted as  $e_{i+1}-e_i$ , as in Weber 171 et al. (2010). And the last two features are related to the position of the Moon with re-172 spect to the Earth, and these are the distance, d, itself and the rate of the distance change, 173 d, as in Bills et al. (2008). 174

175 The selected features all have different ranges and we refer to them as raw data. To train a model that is able to generalise well for a given problem sometimes it is nec-176 essary to transform raw data to a form that is more suitable for training (Langer et al., 177 2019). By applying transformation on the raw data we may obtain a mapping which bet-178 ter reveals patterns in our data. Therefore, we chose to apply trigonometric transforma-179 tion of the true anomaly angle  $\gamma$ , to properly address the jump discontinuities in the fea-180 ture when angle goes from  $2\pi$  to 0, due to it's cyclic nature. This is addressed by trans-181 forming true anomaly angle  $\gamma$  to pair of  $[\cos \gamma, \sin \gamma]$ . An example of all eight features 182 are shown in Figure S2 for nest A1. 183

### <sup>184</sup> 3 Methodology

When new lunar data arrives, we shall be able to differentiate events in groups based 185 on the waveform similarity measurement. And we shall be able to measure their P and 186 S travel times, and thus form set of nests from Section 2.1. Final step would be to as-187 sociate these new events with the existing Apollo nests if possible. This nest identifica-188 tion from a single lunar station is a supervised classification problem. The model is trained 189 in a predictive way by taking into account nest locations as labels and nest lunar orbital 190 parameters as input data. Since we want to predict a class (nest), but we do not have 191 statistically large data set (as previously mentioned A1 has 443 quakes), we choose to 192 train a Random Forest (RF) Classifier, since RF can perform well with any size of datasets 193 and tend not to overfit (Ho, 1995; Breiman, 2001). 194

Random Forest (RF) is a machine learning technique that is based on decision trees 195 (Breiman et al., 1984; Quinlan, 1986) and bootstrap aggregating (Breiman, 1996), where 196 the main output is reached by majority votes among an ensemble of randomised deci-197 sion trees. A main building unit, a decision tree, is a tree-like learning algorithm where 198 each internal node tests on attribute, each branch corresponds to attribute value and each 199 leaf node represents the final prediction. Usually, during the training phase thresholds, 200 order and number of inequality operations within internal nodes are learned. The hy-201 perparameters that define a RF structure, such as the number of trees, and measure which 202 maximises diversity between classes, are determined beforehand (see Text S2 and Fig-203 ure S3). 204

RF also provides an assessment of the feature or input variable importance which might give us an insight of how the model reached its prediction. To assess the feature importance, the RF removes one of the features while it keeps the rest constant, and it measures, among others, the accuracy decrease (Breiman, 2001). RF models are able to solve regression and classification problems, as well as two- and multi-class problems. It has been show that RF can perform with high accuracy even though there are only a few parameters to tune.

In our case, during the training phase, the RF model has access to the extracted features of the individual quakes and the nest labels. The training is performed on a subset of the data, while the model performance is evaluated on the test subset, which the model has never seen. Evaluation is accomplished by comparing the model's predicted class (nest) with the ground truth one. The statistical performance of the model is presented with confusion matrix and Receiver Operating Characteristic (ROC) curve. We
expect that in the case of the ideal RF Classifer the diagonal of the confusion matrix is
equal to 1 (and off-diagonal elements are zero), while ROC curve is passing through the
left upper corner.

# 4 **Results and discussion**

4.1 Training and testing on two largest nests

We first test the hypothesis whether it is possible to differentiate two DMQ nests using the lunar orbital parameters (features). For this, we select the two largest nests, A1 and A8, with a total size of 768 events and ratio A1:A8=0.57 : 0.43 (see feature distributions in Figure S4).

Training and testing our base RF model (see Text S2) with the normalised and not 227 normalised input data, we end up selecting to work with the normalised input data since 228 this model performed better (see Figure S5 and S6). The base model trained with the 229 normalised input data performed with an accuracy of 89%, while precision, recall and 230 f1-score for the A1 nest is 88%, 94%, 91%, respectively, and for A8 is 90%, 80%, 85%, 231 respectively (see Figure S6B) with only the occurrence time knowledge. The ROC curve 232 is above the random classifier curve, meaning that the base model is not randomly clas-233 sifying A1 and A8 nests (see Figure S6C). Out of eight features, the first five most im-234 portant are  $cos(\gamma)$ ,  $\Delta t_{Periqee}$ , d,  $\Delta t_{AscendingNode}$ ,  $e_{i+1} - e_i$  (see Figure S6D). We no-235 tice that  $cos(\gamma)$  is the feature with the most important contribution to the model learn-236 ing. This might be because A1 and A8 have reversed distributions for  $cos(\gamma)$  feature (see 237 Figure S 4D). 238

We proceed into testing learning robustness of our base model in a series of exper-239 iments (see Text S2 and Figures S7-S12), all of which indicate that the model is stable. 240 This implies that the base model generalizes well, and not over fit the results. Further, 241 if we examine why the base model sometimes mislabels the nests (Figure S13), we no-242 tice that the 2D manifold (see Text S3) of feature space spanned by the input data, cal-243 culated by t-sne method (van der Maaten & Hinton, 2008), is not perfectly separated. 244 It seems this segregation might be dominated by a single feature, and that is  $\Delta t_{AscendingNode}$ 245 (see Figure S14 and S15A). 246

#### 247

### 4.2 Training and testing on three and more nests

In this section we study how the performance of our base RF model from Section 4.1 changes by adding more nests. We carry out three tests for the next combinations and their ratios: A1-A8-A18 (45%-33%-22%), A1-A8-A18-A6 (38%-28%-18%-15%), A1-A8-A18-A6-A14 (33%-25%-16%-13%-12%), where the three added nests are the three largest nests besides A1 and A8.

The analysis shows that by adding more nests, the performance of our base model 253 deteriorates since the accuracy drops from 88% to 59% (see Figures S16-S19). By adding 254 a 3rd nest, and we notice that A1 and A8 recalls deteriorate slightly, and 50% is A18 events 255 are classified either as A1 or A8 (see Figure S16). Yet, the precision of A18 is the high-256 est. Features,  $e_{i+1}-e_i$  and d, gain importance. Yet, the importance of all features be-257 come more equalized. By adding a 4th nest, A6, the recall of A1 nest improves, recall 258 of A8 nest deteriorates even more than before, recall of A18 improves notably, and the 259 new added nest A6 has a recall of 46%, by having most of its events misclassified only 260 as A1 nest, and not a single event as A18 (see Figure S17). This might implies that A18 261 and A6 nests have completely different source mechanisms. Less notably than before, 262 the importance of all features is becoming more equalized. Lastly, by adding a 5th nest, 263 A14, the recall of A1 and A8 become the highest, and three other nests perform with 264

recall less than 50%, and their most mislabelled data points are associated with A1 nest (see Figure S18). The importance between features is almost equalized, yet the interval time  $e_{i+1} - e_i$  is the only feature that stands out.

These results might imply that by adding more nests, we add more complexity into the problem, since we might be adding nests that have similar source mechanisms. Having similar source mechanisms means that sources are triggered by tides is the same way, so their lunar orbital features have similar characteristics, and we cannot differentiate between nests without having more data. Furthermore, it seems that the only significantly important feature is the interval time, the only feature that does not reflect the lunar orbital information.

Checking the two dimensional representation of the feature space constructed by
the feature combination of nests A1-A8-A18-A6-A14, we might conclude that for this
particular set it is to some degree impossible to completely differentiate between nests
due to the lack of data (see Figure S20).

279

# 4.3 Training and testing on nest sets

Using the same base RF model from Section 4.1, we proceed to train and test how well we can differentiate nests that belong to the same set shown in Figure 1. We analyze them in three separate groups by the frequency of the nest they contain: A) S1, S2, S3, S4, S12, S13, S14, S15; B) S5, S6, S7, S8; C) S9, S10, S11. The results are shown in Figure 2A, B, and C, respectively.

We observe high value of recall for most of the nests, as well as high accuracy for most of the sets (see Figure 2). Sets that have <= 4 nests perform better than those with more nests, as in sets from group A shown in Figure 2B. When the nest's recall is very low or zero (A11, A30, A41, A42, A50), it signifies a nest with very few events (see ratio of nests in all sets in Figure S18). If we take an example of nest A20, we notice that it has constant recall in many sets (see Figure 2B and C), even though it is not the biggest nest in the set (see Figure S21). Thus, not only the size but probably also the uniqueness of the features determine the success of identifying the nest.

The importance of different features is shown in Figure 3 for all three groups. On 293 one hand, removing just one nest could change the feature importance, as in the case 294 of S2 (where we remove A16) versus S1. On the other hand, we notice that the feature 295 importance does not drastically change when comparing results for sets S3 and S4, where 296 we add nest A44, even though the nest itself is large in size (see Figure S21). For the sets 297 in group B, the feature importance is stable with respect to adding or removing nests. 298 It is quite similar for group C, where only one set S8 has different feature importance. We notice that sets which contain  $\leq 4$  nests (as in group A), there is usually one or 300 two important features, while for sets with > 4 nests there is equalisation of the feature 301 importance (as in groups B and C). This might imply that a single lunar orbital param-302 eter is enough to explain the occurrence of the nests, which are unique in nature. Mix-303 ing more nests suggests that we might be mixing nests with similar temporal patterns, 304 thus learning how to differentiate them is more challenging. Moreover, the feature im-305 portance changes for sets that have unique combinations of nests, which may hint that 306 these nests have different source mechanisms. 307

If we consider a 2D manifold spanned by the sets from groups A, B, C (see Figures S22, S23, S24, respectively), we notice that unique segregation in this space correlates with the RF model accuracy. Nests that form closely spaced homogenized clusters in the 2D manifold tend to be correlated with models that scored high recall for these nests.



**Figure 2.** Performance of RF models designed to classify nests within different sets. A) Travel times  $t_{sp}$  for four lunar model from Garcia et al. (2011) (G11), Garcia et al. (2019) (issi2, issi3), (Khan et al., 2014) (K14) with nest labels. B) Recall for individual nests within each set with respect to their travel times labeled with sets to which they belong and the scored accuracy of this set. C) and D) same as B) just for different group of sets.

5° -	0.13	0.04	0.13	0.18	0.11	0.15	0.15	0.12		
Sr -	0.15	0.04	0.14	0.17	0.07	0.18	0.16	0.08		- 0.225
දිා -	0.24	0.04	0.12	0.17	0.1	0.11	0.13	0.08		- 0 200
5 <sup>A</sup> -	0.24	0.05	0.12	0.13	0.11	0.13	0.12	0.09		0.200
52-	0.12	0.05	0.14	0.13	0.16	0.15	0.12	0.13		- 0.175
53-	0.08	0.07	0.09	0.09	0.24	0.15	0.09	0.19		
51A -	0.18	0.07	0.09	0.15	0.14	0.11	0.12	0.12		- 0.150
Sets -	0.12	0.03	0.14	0.17	0.13	0.12	0.17	0.12		
్ హ -	0.16	0.07	0.11	0.12	0.13	0.17	0.12	0.12		- 0.125
- OS	0.15	0.06	0.09	0.12	0.14	0.21	0.11	0.12		0 1 0 0
51 -	0.16	0.05	0.09	0.12	0.15	0.19	0.11	0.13		- 0.100
- °ې	0.14	0.04	0.1	0.14	0.17	0.14	0.12	0.14		- 0.075
දුව -	0.16	0.07	0.11	0.12	0.12	0.17	0.12	0.12		
570 -	0.16	0.07	0.11	0.12	0.13	0.19	0.12	0.11		- 0.050
51-	0.16	0.07	0.12	0.12	0.12	0.18	0.12	0.11		
	$\Delta t_A$	$\Delta t_F$	$\Delta t_P$	$\cos(\gamma)$	$sin(\gamma)$	Δe	ď	ģ		

Figure 3. Feature importance for Random Forest models associated with different travel time sets shown in Figure 1.

# **5** Conclusion

In this paper we propose how to tackle DMQ nest identification during future lu-314 nar missions that will likely host only one station on the far side of the Moon. We pro-315 pose constraining their location by using differential time travel measurement  $t_{sp}$  and 316 parameters related to the temporal patterns of the DMQ occurrence. First, in our anal-317 ysis we assume that we cannot differentiate between nests whose differences in travel time 318 are less than 5 seconds. Thus, we form set of nests that have similar travel times. Sec-319 ond, for each event within the nests we calculate features that are used to build a Ran-320 321 dom Forest model. This model is trained to differentiate between distinct nests. The features used for training are build by associating each event in all nests with the time dif-322 ference between events' origin time and time of lunar ascending node, Full Moon phase, 323 perigee, then position of the Moon in its orbit expressed by true anomaly angle, distance 324 of the Moon from the Earth, rate change of this distance, and the time between two suc-325 cessive quakes. We show that by training Random Forest models to differentiate between 326 distinct nests within sets, we can obtain models with high accuracy (more than half of 327 the models score above 70% accuracy). Yet, the performances of these models depend 328 on the number of nests within the set. More nests implies that the problem is more dif-329 ficult to solve, probably because a) nests might have similar source mechanisms, b) the 330 number of events within nests is unbalanced, and c) we don't have enough data. Since 331 RF models also arrange features by their importance to make a final classification de-332 cision, we observe that the importance of the features change with different sets. This 333 complements the findings of previous papers, since it signifies that nests do correspond 334 to different lunar events, which eventually might be connected to the distribution of tidal 335 stresses during these events. Finally, our model provides a good first approximation of 336 the nest identification. And as the catalog of new events grows, it will be straightforward 337 to retrain RF model with the new enlarged dataset. 338

# 339 Open Research Section

The deep moonquake catalog used in this study is published in Nakamura et al. (1981), and revisited in Nunn et al. (2020). Python package Skyfield used to calculate Moon's orbital parameters based on JPL ephemeris can be found on the website https:// rhodesmill.org/skyfield/ (Rhodes, 2019, Software). For our implementation of the Random Forest algorithm we use Scikit-learn machine learning Python library (Pedregosa et al., 2011).

# 346 Acknowledgments

French co-authors thanks the French Space Agency, CNES, for supporting this research in the frame of the French contribution to FSS as well as IDEX Paris Cité (ANR-18-IDEX-0001). MPP was supported by funds from the Jet Propulsion Laboratory, under contract with the National Aeronautics and Space Administration (NASA).

# 351 References

- Banerdt, W. B., Smrekar, S. E., Banfield, D., Giardini, D., Golombek, M., Johnson,
   C. L., ... others (2020). Initial results from the insight mission on mars.
   *Nature Geoscience*, 13(3), 183–189.
- Bills, B. G., Bulow, R., & Johnson, C. L. (2008). Influence of earth-moon orbit geometry on deep moonquake occurrence times. *Lunar Planet. Sci.*, XXXIX.
- Breiman, L. (1996). Bagging predictors. *Machine learning*, 24, 123–140.
- Breiman, L. (2001). Random forests. *Machine learning*, 45, 5–32.
- Breiman, L., Friedman, J., Stone, C. J., & Olshen, R. A. (1984). Classification and regression trees. CRC press.

361	Bulow, R., Johnson, C., Bills, B., & Shearer, P. (2007). Temporal and spatial prop-
362	erties of some deep moonquake clusters. Journal of Geophysical Research: Blue et = 110 (EQ)
363	Planets, 112 (E9).
364	bulow, R., Johnson, C., & Snearer, P. (2005). New events discovered in the apolio lunar seismic data. <i>Lournal of Ceonhusical Research: Planets</i> 110(E10)
305	Compared W S & Cilhoppy I I (1067) Operation mean blink and report of ab
366 367	cameron, W. S., & Gineany, J. J. (1967). Operation moon blink and report of ob- servations of lunar transient phenomena. <i>Icarus</i> , $\gamma(1-3)$ , 29–41.
368	Cheng, C. H., & Toksöz, M. N. (1978). Tidal stresses in the moon. Journal of
369	Geophysical Research: Solid Earth, 83(B2), 845-853. Retrieved from https://
370	agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JB083iB02p00845
371	doi: https://doi.org/10.1029/JB083iB02p00845
372	Cutler, J. W., Nguyen, T. A., Kano, T., Kumar, Y. R. L., Panning, M., April, S., &
373	Haque, S. (2023). Overview of the avionics design for the farside seismic suite.
374	In Aiaa scitech 2023 forum. Retrieved from https://arc.aiaa.org/doi/abs/
375	10.2514/6.2023-1880 doi: 10.2514/6.2023-1880
376	Dainty, A., Goins, N., & Toksoz, M. (1975). Natural lunar seismic events and the
377	structure of the moon. In In: Lunar science conference, 6th, houston, tex.,
378	march 17-21, 1975, proceedings. volume 3.(a78-46741 21-91) new york, perga-
379	mon press, inc., 1975, p. 2887-2897. (Vol. 6, pp. 2887–2897).
380	Ewing, M., Latham, G., Press, F., Sutton, G., Dorman, J., Nakamura, Y., Ko-
381	vach, R. (1971). Seismology of the moon and implications on internal struc-
382	ture, origin and evolution. <i>Highlights of astronomy</i> , 2, 155–172.
383	Frohlich, C., & Nakamura, Y. (2009). The physical mechanisms of deep moon-
384	quakes and intermediate-depth earthquakes: How similar and how differ-
385	ent? Physics of the Earth and Planetary Interiors, 173(3), 365-374. Re-
386	trieved from https://www.sciencedirect.com/science/article/pii/
387	S0031920109000338 doi: https://doi.org/10.1016/j.pepi.2009.02.004
388	Gagnepain-Beyneix, J., Lognonné, P., Chenet, H., Lombardi, D., & Spohn, T.
389	(2006). A seismic model of the lunar mantle and constraints on tempera-
390	ture and mineralogy. Physics of the Earth and Planetary Interiors, 159(3-4),
391	140-100. Carrie P. F. Carnensin Permeire I. Charnet S. & Lornenné P. (2011). Voru pro-
392	liminary reference moon model Physics of the Earth and Planetary Interiors
393	188(1-2), 96-113.
395	Garcia, R. F., Khan, A., Drilleau, M., Margerin, L., Kawamura, T., Sun, D., oth-
396	ers (2019). Lunar seismology: An update on interior structure models. Space
397	Science Reviews, 215, 1–47.
398	Ho, T. K. (1995). Random decision forests. In Proceedings of 3rd international con-
399	ference on document analysis and recognition (Vol. 1, pp. 278–282).
400	Kawamura, T., Grott, M., Garcia, R., Wieczorek, M., de Raucourt, S., Lognonné,
401	P., others (2022). An autonomous lunar geophysical experiment package
402	(algep) for future space missions: In response to call for white papers for the
403	voyage 2050 long-term plan in the esa science program. Experimental Astron-
404	omy, 54 (2-3), 617-640.
405	Kawamura, T., Lognonné, P., Nishikawa, Y., & Tanaka, S. (2017). Evaluation of
406	deep moonquake source parameters: Implication for fault characteristics and
407	thermal state. Journal of Geophysical Research: Planets, 122(7), 1487–1504.
408	Khan, A., Connolly, J. A. D., Pommier, A., & Noir, J. (2014). Geophysical evi-
409	dence for melt in the deep lunar interior and implications for lunar evolution.
410	Journal of Geophysical Research: Planets, 119(10), 2197-2221. Retrieved
411	IFOIII nttps://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2014 IE004661 doi: https://doi.org/10.1002/2014 IE004661
412	Z014JE004001 uoi: https://uoi.org/10.1002/2014JE004001 Khan A fr Moreorened K (2002) An inquiry into the lunger interior. A gravity
413	ear inversion of the apollo lunar seismic data. <i>Lowrad of Coonductional Descented</i> .
414	Planets $107(\text{F6})$ 3–1
410	1 0010000, 107(110), 0 1.

- Khan, A., Mosegaard, K., & Rasmussen, K. L. (2000). A new seismic velocity model 416 for the moon from a monte carlo inversion of the apollo lunar seismic data. 417 Geophysical Research Letters, 27(11), 1591–1594. 418 Lammlein, D. R. (1977a). Lunar seismicity and tectonics. *Physics of the* 419 Earth and Planetary Interiors, 14(3), 224-273. Retrieved from https:// 420 www.sciencedirect.com/science/article/pii/0031920177901753 doi: 421 https://doi.org/10.1016/0031-9201(77)90175-3 422 Lammlein, D. R. (1977b). Lunar seismicity and tectonics. *Physics of the Earth and* 423 *Planetary Interiors*, *14*(3), 224–273. 424 Lammlein, D. R., Latham, G. V., Dorman, J., Nakamura, Y., & Ewing, M. (1974).425 Lunar seismicity, structure, and tectonics. Reviews of Geophysics, 12(1), 1-21. 426 Langer, H., Falsaperla, S., & Hammer, C. (2019). Advantages and pitfalls of pattern 427 recognition: selected cases in geophysics. Elsevier. 428 Lognonné, P., Banerdt, W. B., Giardini, D., Pike, W. T., Christensen, U., Laudet, 429 P., ... others (2019). Seis: Insight's seismic experiment for internal structure 430 of mars. Space Science Reviews, 215, 1-170. 431 Lognonné, P., Banerdt, W. B., Pike, W. T., Giardini, D., Christensen, U., Garcia, 432 R. F., ... others (2020). Constraints on the shallow elastic and anelastic struc-433 ture of mars from insight seismic data. Nature Geoscience, 13(3), 213–220. 434 Lognonné, P., Gagnepain-Beyneix, J., & Chenet, H. (2003). A new seismic model of 435 the moon: implications for structure, thermal evolution and formation of the 436 moon. Earth and Planetary Science Letters, 211(1-2), 27–44. 437 Lognonné, P., & Johnson, C. (2015). 10.03—planetary seismology. Treatise on geo-438 physics, 2, 65–120. 439 Meeus, J. (1991). Astronomical algorithms. *Richmond*. 440 Middlehurst, B. M. (1967). An analysis of lunar events. Reviews of Geophysics, 441 5(2), 173-189.442 Minshull, T., & Goulty, N. (1988a). The influence of tidal stresses on deep moon-443 quake activity. Physics of the Earth and Planetary Interiors, 52(1), 41-55. 444 Retrieved from https://www.sciencedirect.com/science/article/pii/ 445 0031920188900568 doi: https://doi.org/10.1016/0031-9201(88)90056-8 446 Minshull, T., & Goulty, N. (1988b). The influence of tidal stresses on deep moon-447 quake activity. Physics of the earth and planetary interiors, 52(1-2), 41–55. 448 Moore, P. (1968). Transient lunar phenomena: A review, 1967. J. Br. Astron. As-449 soc., 78, 138-144. 450 Nakamura, Y. (1983). Seismic velocity structure of the lunar mantle. Journal of 451 Geophysical Research: Solid Earth, 88(B1), 677-686. Retrieved from https:// 452 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JB088iB01p00677 453 doi: https://doi.org/10.1029/JB088iB01p00677 454 Nakamura, Y. (2003).New identification of deep moonquakes in the apollo lunar 455 seismic data. Physics of the Earth and Planetary Interiors, 139(3), 197-205. 456 Retrieved from https://www.sciencedirect.com/science/article/pii/ 457 S0031920103002103 doi: https://doi.org/10.1016/j.pepi.2003.07.017 458 Farside deep moonquakes and deep interior of the moon. Nakamura, Y. (2005).459 Journal of Geophysical Research: Planets, 110(E1). Retrieved from https:// 460 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2004JE002332 doi: 461 https://doi.org/10.1029/2004JE002332 462 Nakamura, Y., Latham, G., Dorman, H., & Harris, J. (1981). Passive seismic ex-463 Galveston Geophysics Laboratory Contriperiment long-period event catalog. 464 bution, 491, 314. 465 Nakamura, Y., Latham, G. V., & Dorman, H. J. (1982a). Apollo lunar seismic 466 experiment—final summary. Journal of Geophysical Research: Solid Earth, 467 87(S01), A117–A123. 468 Nakamura, Y., Latham, G. V., & Dorman, H. J. (1982b). Apollo lunar seismic 469
- 470 experiment—final summary. Journal of Geophysical Research: Solid Earth,

471	87(S01), A117-A123. Retrieved from https://agupubs.onlinelibrary
472	.wiley.com/doi/abs/10.1029/JB087iS01p0A117 doi: https://doi.org/
473	10.1029/JB087iS01p0A117
474	Nunn, C., Garcia, R. F., Nakamura, Y., Marusiak, A. G., Kawamura, T., Sun, D.,
475	others (2020). Lunar seismology: A data and instrumentation review.
476	Space Science Reviews, 216(5), 89.
477	Panning, M., Kedar, S., Bowles, N., Calcutt, S., Cutler, J., Elliott, J., others
478	(2021). Farside seismic suite (fss): First seismic data from the farside of the
479	moon delivered by a commercial lander. In Agu fall meeting abstracts (Vol.
480	2021, pp. P54C–01).
481	Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,
482	others (2011). Scikit-learn: Machine learning in python. the Journal of
483	machine Learning research, 12, 2825–2830.
484	Pickrell, J. (2022, May). These six countries are about to go to the moon — here's
485	why. Nature, 605, 208–211. doi: 10.1038/d41586-022-01252-7
486	Quinlan, J. R. (1986). Induction of decision trees. Machine learning, 1, 81–106.
487	Rhodes, B. (2019). Skyfield: Generate high precision research-grade positions for
488	stars, planets, moons, and earth satellites. Astrophysics Source Code Library,
489	ascl-1907.
490	Standley, I. M., Pike, W. T., Calcutt, S., & Hoffman, J. P. (2023). Short period seis-
491	mometer for the lunar farside seismic suite mission. In 2023 ieee aerospace con-
492	ference (p. 1-9). doi: 10.1109/AERO55745.2023.10115559
493	Toksoz, M. N., Dainty, A. M., Solomon, S. C., & Anderson, K. R. (1974). Structure
494	of the moon. <i>Reviews of Geophysics</i> , 12(4), 539-567. Retrieved from https://
495	agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/RG0121004p00539
496	doi: https://doi.org/10.1029/RG0121004p00539
497	lotion to tidal strongen. <i>Science</i> 106(4202) 070 081
498	Turner A B Hawtherne I C & Caddes M (2022) Stresses in the lunar
499	interior: Insights from slip directions in the 201 doop moonquelo nest
500	nd of Geonbusical Research: Planets 197(12) a2022 IE007364 Betrieved
501	from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/
502	2022.IE007364 doi: https://doi.org/10.1029/2022.IE007364
504	van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-sne. Journal
505	of Machine Learning Research, 9(86), 2579–2605. Retrieved from http://imlr
506	.org/papers/v9/vandermaaten08a.html
507	Weber, R., Bills, B., & Johnson, C. (2009). Constraints on deep moonquake focal
508	mechanisms through analyses of tidal stress. <i>Journal of Geophysical Research</i> :
509	<i>Planets</i> , 114 (E5).
510	Weber, R., Bills, B., & Johnson, C. (2010). A simple physical model for deep moon-
511	quake occurrence times. Physics of the Earth and Planetary Interiors, 182(3-
512	4), 152–160.
513	Weber, R., Neal, C., Grimm, R., Grott, M., Schmerr, N., Wieczorek, M., oth-
514	ers (2021). The scientific rationale for deployment of a long-lived geophysical
515	network on the moon. Bulletin of the AAS, $53(4)$ .
516	Witze, A. (2022, May). The \$93-billion plan to put astronauts back on the moon.

<sup>517</sup> Nature, 605, 212–216. doi: 10.1038/d41586-022-01253-6

-14-

# Identifying deep moonquake nests using machine learning model on single lunar station on the far side of the Moon

# Josipa Majstorović<sup>1</sup>, Philippe Lognonné<sup>1</sup>, Taichi Kawamura<sup>1</sup>, Mark P. Panning<sup>2</sup>

 $^1$ Université Paris Cité, Institut de physique du globe de Paris, CNRS, Paris, France $^2$ Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA

# Key Points:

1

2

3

4

5

6

8

9	• As a part of the future space mission NASA will deploy a new seismic station to
10	Schrödinger Basin on the far side of the Moon.
11	• We propose a machine learning model trained to classify deep moonquakes using
12	the lunar orbital parameter.
13	• The models perform with accuracy greater than 70% when trained to classify com
14	binations of four or fewer nests.

 $Corresponding \ author: \ Josipa \ Majstorović, \ \texttt{josipa.majstorovic} @\texttt{protonmail.com}$ 

#### 15 Abstract

One of the future NASA space program includes the Farside Seismic Suite (FSS) pay-16 load, a single station with two seismometers, on the far side of the Moon. During FSS 17 operations, the processing of the data will provide us with new insight into the Moon's 18 seismic activity. One of Apollo mission finding is the existence of deep moonquakes (DMQ), 19 and the nature of their temporal occurrence patterns as well as the spatially clustering. 20 It has been shown that DMQs reside in about 300 source regions. In this paper we tackle 21 how we can associate new events with these source regions using the single station data. 22 We propose a machine learning model that is trained to differentiate between DMQ nests 23 using only the lunar orbital parameters related to DMQ time occurrences. We show that 24 ML models perform well (with an accuracy > 70%) when they are trained to classify 25 less than 4 nests. 26

# 27 Plain Language Summary

The future space missions will provide us with various new lunar data, one of which 28 will be ground vibration measurement. The studies of these measurements from the Apollo 29 era in 70s, showed that Moon can host various events. The most intriguing ones are deep 30 moonquakes (DMQs), which are events associated with the displacement deep within the 31 lunar interior. It has been shown that DMQs occur in the specific locations, which are 32 called nests, and that their temporal occurrence is related to the monthly motion of the 33 Moon around the Earth. In this paper we tackle how we can associate new events from 34 only one station located on the far side of the Moon with these known locations of DMQs. 35 We propose a machine learning model that is trained to classify DMQ nests, only using 36 the information about their temporal occurrences, e.g. time of the event, described in 37 terms of different lunar events. We report that models are performing well (with an ac-38 (uracy > 70%) when they are trained to classify 4 or fewer nests. This gives us a good 39 first approximation about the nest identification. 40

#### 41 **1** Introduction

We are at the beginning of a golden age of lunar exploration as many nations, to-42 gether with private companies, are establishing numerous efforts to obtain new scien-43 tific measurement of the Moon (Weber et al., 2021; Pickrell, 2022; Kawamura et al., 2022). 44 In light of this, NASA established the Artemis program which should land a crewed mis-45 sion at the lunar south pole (Witze, 2022). This would be the first attempt of a crewed 46 landing after the successful Apollo missions in 1970's. Before Artemis missions land on 47 the Moon, NASA has also established the Commercial Lunar Payload Services (CLPS) 48 program to land scientific payloads on the Moon using commercial landers. The Farside 49 Seismic Suite (FSS) is one of the selected payloads, and it will deliver two seismometers 50 to Schrödinger Basin on the far side of the Moon (Panning et al., 2021; Standley et al., 51 2023; Cutler et al., 2023): one vertical Very BroadBand seismometer, and Short Period 52 sensor, both spare or derived from the SEIS experiment sensors (Lognonné et al., 2019, 53 2020) from the InSight mission to Mars (Banerdt et al., 2020). 54

The Apollo missions showed the importance of deploying sensors on the surface of 55 the Moon, since a great deal of our knowledge about the Moon comes from the analy-56 sis of data acquired during the Apollo era (Lognonné & Johnson, 2015). Thus, analyz-57 ing ground motion measurements provided the community with the first constraints on 58 the lunar interior and the activity at the surface (Nakamura et al., 1982a, 1982b; Khan 59 et al., 2000; Khan & Mosegaard, 2002; Khan et al., 2014; Lognonné et al., 2003; Gagnepain-60 Beyneix et al., 2006; Weber et al., 2010; Garcia et al., 2011; Kawamura et al., 2017; Gar-61 cia et al., 2019; Nunn et al., 2020). It has also revealed that the Moon can host events 62 of various origins, such as shallow and deep moonquakes, meteoroid and artificial impacts 63 (Toksöz et al., 1974; Dainty et al., 1975; Lammlein, 1977a; Nakamura, 1983, 2003, 2005). 64

Today, we have more than twelve thousands events, out of which the deep moonquakes (DMQs) form the most numerous group (Nakamura et al., 1981; Nakamura, 2005).

DMQs are a distinctive group of seismic events that originate from depths between 67 700 and 1200 km, at high pressure and temperature conditions, where little brittle de-68 formation is expected. Due to very high waveform similarity between quakes, the DMQs 69 have been clustered into about 300 source regions or nests (Nakamura, 2003). This has 70 interpreted to be a consequence of DMQs occurring repeatedly at the fixed nests, which 71 are located mostly on the near side of the Moon. It has been shown that time occurrence 72 73 of the DMQs is correlated with the monthly motion of the Moon around the Earth. Thus, DMQ occurrences exhibit tidal periodicities and furthermore, the associated high strain 74 rates might explain brittle processes (Kawamura et al., 2017). However, the real causes 75 of their origins are yet to be discovered. There are two puzzling fact about their origin: 76 a) cyclic tidal stresses, caused by the monthly motion of Moon around Earth, are far less 77 than the estimated confining pressures where DMQs occurs (Cheng & Toksöz, 1978; Min-78 shull & Goulty, 1988a); b) do we need both, tectonic and ambient tidal stresses, to ex-79 plain their mechanical origin (Frohlich & Nakamura, 2009). 80

To better constrain the lunar interior and unravel the cause of DMQs, it is impor-81 tant to locate new events and associate them with the known nest locations from Apollo 82 with future lunar missions like FSS. These new observations will add, for each new nest, 83 a new differential  $t_s - t_p$  measurement constraining the deep interior with a different 84 epicentral distance. However, due to the mission requirements, it is extremely likely that, 85 at the beginning, we might have only one lunar station at the disposal. Therefore, in this 86 paper we study the problem of DMQ nest identification without using waveform infor-87 mation. This is due to the new location of the recording station, which will not match 88 existing Apollo-era waveform templates due to different propagation paths. We propose 89 a machine learning (ML) model that is trained to identify nests within the set of nests 90 of similar differential travel times. The main features used for the model training are re-91 lated to the fact that different nests respond differently to lunar cycle. 92

Very early studies have shown correlation between lunar transient events and po-93 sition of the Moon related to the Earth (Middlehurst, 1967; Cameron & Gilheany, 1967; 94 Moore, 1968). This further encouraged observations that some moonquakes occur with 95 periods that reflect Earth-Moon-Sun relationship (Ewing et al., 1971). Later, it has been 96 shown that the occurrence of DMQs are related to tidal stress cycles, and correlations 97 between DMQs occurrence times and lunar monthly tidal cycles have been indicated (Lammlein 98 et al., 1974; Toksöz et al., 1977; Lammlein, 1977b; Cheng & Toksöz, 1978; Minshull & 99 Goulty, 1988b). The lunar cycle can be explained with three lunar months: synodic, dra-100 conic, anomalistic. Synodic month is the period of lunar phases such as New Moon, First 101 Quarter, Full Moon, Last Quarter. Draconic month is the period between two nodes, as-102 cending or descending, where the nodes are points at which the Moon's orbit plane crosses 103 the ecliptic plane towards which it is inclined of about 5.14°. Anomalistic month is the 104 period between two extreme points, perigee and apogee, since the Moon's orbit approx-105 imates an ellipse rather than a circle. Earlier studies counted the number of events per 106 day as a function of time and found 0.5 and 1 month signals in the occurrence times re-107 lated to anomalistic and draconic period of 27 days (Lammlein et al., 1974; Lammlein, 108 1977b). The same studies also indicated 206-day and 6-year periods, related to the Sun's 109 perturbation on the lunar orbit and the relative precession of the perigee of the Moon's 110 orbit. Subsequently, many recent papers studied and confirmed tidal periodicities of DMQs 111 and more (Bulow et al., 2005, 2007; Bills et al., 2008; Frohlich & Nakamura, 2009; We-112 ber et al., 2009, 2010; Turner et al., 2022). 113

Based on the previous papers, it is clear that DMQ nests exhibit some clear temporal patterns in their occurrences, and that these are correlated with Moon-Earth system. Therefore, the open question is whether we can design features which would reflect these temporal patterns and to further use those features to study the nest identification with one lunar station. In this paper we tackle the question of defining optimal features and the machine learning model. The paper is organised as follows: first, we discuss data used in the analysis, the existing catalog of DMQ events. Second, we discuss the feature design. Third, we introduce a machine learning model. Fourth, we discuss successes and pitfalls of the machine learning model for nest identification applied to different combination of nests. We conclude how this study can offer some first estimates of the nest location in the future lunar missions.

#### 125 2 Data

We start with the existing catalog of lunar events (Nakamura et al., 1981), which 126 was updated in 2008 and modified in 2018 (Nunn et al., 2020). Catalog contains a list 127 of events (shallow and deep moonquakes, meteoroid and artificial impacts) with attributes 128 such as date and time of the event occurrence, signal envelope amplitude as measured 129 in mm on a standard plot, data availability per station, and the nest (source) classifi-130 cation for DMQs. It is important to note that the source classification is not an exact 131 location defined by latitude and longitude, rather a result from the waveform cross-correlation 132 and single-link cluster analysis (Nakamura, 2003). This analysis positively clustered around 133 7k DMQs into 77 nests, where the largest nest is associated with label A1. This nest con-134 tains 443 guakes and it is placed on the near side of the Moon. 135

136

# 2.1 Catalog processing: nest sets based on travel time information

Earlier studies published lunar interior models and location of DMQ nests in terms 137 of latitude and longitude by picking P and S travel times on the quake waveforms (Garcia 138 et al., 2019, review of these picks). Using lunar interior models and nest locations we can 139 define nests that are close by in distance if we consider only the  $t_{sp}$  travel time measure-140 ments. To do so we assume that a) our single lunar station is located on the far side in 141 Schrödinger Basin (FSS landing site at 71.378°S, 138.248°E), b) nests' latitudes and lon-142 gitudes from (Lognonné et al., 2003), c) calculated P- and S- wave travel times ( $t_p$  and 143  $t_s$ , respectively) using lunar velocity model between landing site and locations of DMQ 144 nests. By having  $t_s$  and  $t_p$  we can calculate  $t_{sp} = t_s - t_p$  for all nests and models shown 145 in Figure 1 (see Text S1 and Figure S1 for further explanation). Next, we count for each 146 nest how many there are with the similar  $t_{sp}$  travel time measurement assuming a pick-147 ing error of 5 seconds as shown in Figure 1A, consistent with the average picking error 148 in Lognonné et al. (2003). This count provides us with the different sets  $S_i$ , shown in Figure 1A, that contain nests  $N_j$  of similar travel times. In other words, if we are able 150 to measure  $t_s$  and  $t_p$  of the new lunar event with accuracy within 5 seconds, we are not 151 able to differentiate between nests that belong to different sets  $S_i$ . Therefore, to further 152 tackle the nest identification problem we proceed to associate each event with a com-153 bination of lunar orbital measurements. 154

155

#### 2.2 Feature selection based on the lunar orbital information

It has been shown that the DMQ temporal patterns in time occurrences are related 156 to different lunar cycles and that these patterns differ from nest to nest. Three lunar cy-157 cles are synodic, draconic, anomalistic, and they all have similar periods, but are marked 158 by different motions, either as the motion between two Full Moons phases, or two nodes, 159 or two apsis, respectively. One can list all the events when Moon is in the Full Moon (New 160 Moon) phase, passing through ascending (descending) node or perigee (apogee) by sim-161 ply looking at the Moon's ephemeris (Meeus, 1991). To make sure that we take into ac-162 count the temporal patterns, we design the main three features as a time difference be-163 tween the time of the quake in the nest and the time of the Moon's Full Moon, ascend-164 ing node and perigee, denoting it as  $\Delta t_{FullMoon}, \Delta t_{AscendingNode}, \Delta t_{Perigee}$ , respectively. 165 We can achieve the same effect by taking the other three time axis as referent one (New 166



Figure 1. Location study of the DMQ nests from the perspective of  $t_{sp} = t_s - t_p$  travel time measurements if we place station in the Schrödinger Basin and consider four different lunar models. A) Upper panel:  $t_{sp}$  travel time measurements for four lunar models from Garcia et al. (2011) (G11), Garcia et al. (2019) (issi2, issi3), (Khan et al., 2014) (K14) with nest labels; A) Lower panel: Sets  $S_i$  which represent nests with <u>similar</u> travel times if we consider a travel time error of 5 seconds. B) Lunar map with the nests locations where the color indicate the median  $t_{sp}$  for four lunar model.

Moon, descending node, and apogee). The next feature is related to the Moon position 167 within its orbit as in Frohlich and Nakamura (2009). The angle between the direction 168 of perigee and the current position of the body, as seen from the main focus of the el-169 lipse, is called the true anomaly, denoted further as  $\gamma$ . Further, as one of the feature we 170 also use the interval time between two quakes in the nest, noted as  $e_{i+1}-e_i$ , as in Weber 171 et al. (2010). And the last two features are related to the position of the Moon with re-172 spect to the Earth, and these are the distance, d, itself and the rate of the distance change, 173 d, as in Bills et al. (2008). 174

175 The selected features all have different ranges and we refer to them as raw data. To train a model that is able to generalise well for a given problem sometimes it is nec-176 essary to transform raw data to a form that is more suitable for training (Langer et al., 177 2019). By applying transformation on the raw data we may obtain a mapping which bet-178 ter reveals patterns in our data. Therefore, we chose to apply trigonometric transforma-179 tion of the true anomaly angle  $\gamma$ , to properly address the jump discontinuities in the fea-180 ture when angle goes from  $2\pi$  to 0, due to it's cyclic nature. This is addressed by trans-181 forming true anomaly angle  $\gamma$  to pair of  $[\cos \gamma, \sin \gamma]$ . An example of all eight features 182 are shown in Figure S2 for nest A1. 183

### <sup>184</sup> 3 Methodology

When new lunar data arrives, we shall be able to differentiate events in groups based 185 on the waveform similarity measurement. And we shall be able to measure their P and 186 S travel times, and thus form set of nests from Section 2.1. Final step would be to as-187 sociate these new events with the existing Apollo nests if possible. This nest identifica-188 tion from a single lunar station is a supervised classification problem. The model is trained 189 in a predictive way by taking into account nest locations as labels and nest lunar orbital 190 parameters as input data. Since we want to predict a class (nest), but we do not have 191 statistically large data set (as previously mentioned A1 has 443 quakes), we choose to 192 train a Random Forest (RF) Classifier, since RF can perform well with any size of datasets 193 and tend not to overfit (Ho, 1995; Breiman, 2001). 194

Random Forest (RF) is a machine learning technique that is based on decision trees 195 (Breiman et al., 1984; Quinlan, 1986) and bootstrap aggregating (Breiman, 1996), where 196 the main output is reached by majority votes among an ensemble of randomised deci-197 sion trees. A main building unit, a decision tree, is a tree-like learning algorithm where 198 each internal node tests on attribute, each branch corresponds to attribute value and each 199 leaf node represents the final prediction. Usually, during the training phase thresholds, 200 order and number of inequality operations within internal nodes are learned. The hy-201 perparameters that define a RF structure, such as the number of trees, and measure which 202 maximises diversity between classes, are determined beforehand (see Text S2 and Fig-203 ure S3). 204

RF also provides an assessment of the feature or input variable importance which might give us an insight of how the model reached its prediction. To assess the feature importance, the RF removes one of the features while it keeps the rest constant, and it measures, among others, the accuracy decrease (Breiman, 2001). RF models are able to solve regression and classification problems, as well as two- and multi-class problems. It has been show that RF can perform with high accuracy even though there are only a few parameters to tune.

In our case, during the training phase, the RF model has access to the extracted features of the individual quakes and the nest labels. The training is performed on a subset of the data, while the model performance is evaluated on the test subset, which the model has never seen. Evaluation is accomplished by comparing the model's predicted class (nest) with the ground truth one. The statistical performance of the model is presented with confusion matrix and Receiver Operating Characteristic (ROC) curve. We
expect that in the case of the ideal RF Classifer the diagonal of the confusion matrix is
equal to 1 (and off-diagonal elements are zero), while ROC curve is passing through the
left upper corner.

# 4 **Results and discussion**

4.1 Training and testing on two largest nests

We first test the hypothesis whether it is possible to differentiate two DMQ nests using the lunar orbital parameters (features). For this, we select the two largest nests, A1 and A8, with a total size of 768 events and ratio A1:A8=0.57 : 0.43 (see feature distributions in Figure S4).

Training and testing our base RF model (see Text S2) with the normalised and not 227 normalised input data, we end up selecting to work with the normalised input data since 228 this model performed better (see Figure S5 and S6). The base model trained with the 229 normalised input data performed with an accuracy of 89%, while precision, recall and 230 f1-score for the A1 nest is 88%, 94%, 91%, respectively, and for A8 is 90%, 80%, 85%, 231 respectively (see Figure S6B) with only the occurrence time knowledge. The ROC curve 232 is above the random classifier curve, meaning that the base model is not randomly clas-233 sifying A1 and A8 nests (see Figure S6C). Out of eight features, the first five most im-234 portant are  $cos(\gamma)$ ,  $\Delta t_{Periqee}$ , d,  $\Delta t_{AscendingNode}$ ,  $e_{i+1} - e_i$  (see Figure S6D). We no-235 tice that  $cos(\gamma)$  is the feature with the most important contribution to the model learn-236 ing. This might be because A1 and A8 have reversed distributions for  $cos(\gamma)$  feature (see 237 Figure S 4D). 238

We proceed into testing learning robustness of our base model in a series of exper-239 iments (see Text S2 and Figures S7-S12), all of which indicate that the model is stable. 240 This implies that the base model generalizes well, and not over fit the results. Further, 241 if we examine why the base model sometimes mislabels the nests (Figure S13), we no-242 tice that the 2D manifold (see Text S3) of feature space spanned by the input data, cal-243 culated by t-sne method (van der Maaten & Hinton, 2008), is not perfectly separated. 244 It seems this segregation might be dominated by a single feature, and that is  $\Delta t_{AscendingNode}$ 245 (see Figure S14 and S15A). 246

#### 247

### 4.2 Training and testing on three and more nests

In this section we study how the performance of our base RF model from Section 4.1 changes by adding more nests. We carry out three tests for the next combinations and their ratios: A1-A8-A18 (45%-33%-22%), A1-A8-A18-A6 (38%-28%-18%-15%), A1-A8-A18-A6-A14 (33%-25%-16%-13%-12%), where the three added nests are the three largest nests besides A1 and A8.

The analysis shows that by adding more nests, the performance of our base model 253 deteriorates since the accuracy drops from 88% to 59% (see Figures S16-S19). By adding 254 a 3rd nest, and we notice that A1 and A8 recalls deteriorate slightly, and 50% is A18 events 255 are classified either as A1 or A8 (see Figure S16). Yet, the precision of A18 is the high-256 est. Features,  $e_{i+1}-e_i$  and d, gain importance. Yet, the importance of all features be-257 come more equalized. By adding a 4th nest, A6, the recall of A1 nest improves, recall 258 of A8 nest deteriorates even more than before, recall of A18 improves notably, and the 259 new added nest A6 has a recall of 46%, by having most of its events misclassified only 260 as A1 nest, and not a single event as A18 (see Figure S17). This might implies that A18 261 and A6 nests have completely different source mechanisms. Less notably than before, 262 the importance of all features is becoming more equalized. Lastly, by adding a 5th nest, 263 A14, the recall of A1 and A8 become the highest, and three other nests perform with 264

recall less than 50%, and their most mislabelled data points are associated with A1 nest (see Figure S18). The importance between features is almost equalized, yet the interval time  $e_{i+1} - e_i$  is the only feature that stands out.

These results might imply that by adding more nests, we add more complexity into the problem, since we might be adding nests that have similar source mechanisms. Having similar source mechanisms means that sources are triggered by tides is the same way, so their lunar orbital features have similar characteristics, and we cannot differentiate between nests without having more data. Furthermore, it seems that the only significantly important feature is the interval time, the only feature that does not reflect the lunar orbital information.

Checking the two dimensional representation of the feature space constructed by
the feature combination of nests A1-A8-A18-A6-A14, we might conclude that for this
particular set it is to some degree impossible to completely differentiate between nests
due to the lack of data (see Figure S20).

279

# 4.3 Training and testing on nest sets

Using the same base RF model from Section 4.1, we proceed to train and test how well we can differentiate nests that belong to the same set shown in Figure 1. We analyze them in three separate groups by the frequency of the nest they contain: A) S1, S2, S3, S4, S12, S13, S14, S15; B) S5, S6, S7, S8; C) S9, S10, S11. The results are shown in Figure 2A, B, and C, respectively.

We observe high value of recall for most of the nests, as well as high accuracy for most of the sets (see Figure 2). Sets that have <= 4 nests perform better than those with more nests, as in sets from group A shown in Figure 2B. When the nest's recall is very low or zero (A11, A30, A41, A42, A50), it signifies a nest with very few events (see ratio of nests in all sets in Figure S18). If we take an example of nest A20, we notice that it has constant recall in many sets (see Figure 2B and C), even though it is not the biggest nest in the set (see Figure S21). Thus, not only the size but probably also the uniqueness of the features determine the success of identifying the nest.

The importance of different features is shown in Figure 3 for all three groups. On 293 one hand, removing just one nest could change the feature importance, as in the case 294 of S2 (where we remove A16) versus S1. On the other hand, we notice that the feature 295 importance does not drastically change when comparing results for sets S3 and S4, where 296 we add nest A44, even though the nest itself is large in size (see Figure S21). For the sets 297 in group B, the feature importance is stable with respect to adding or removing nests. 298 It is quite similar for group C, where only one set S8 has different feature importance. We notice that sets which contain  $\leq 4$  nests (as in group A), there is usually one or 300 two important features, while for sets with > 4 nests there is equalisation of the feature 301 importance (as in groups B and C). This might imply that a single lunar orbital param-302 eter is enough to explain the occurrence of the nests, which are unique in nature. Mix-303 ing more nests suggests that we might be mixing nests with similar temporal patterns, 304 thus learning how to differentiate them is more challenging. Moreover, the feature im-305 portance changes for sets that have unique combinations of nests, which may hint that 306 these nests have different source mechanisms. 307

If we consider a 2D manifold spanned by the sets from groups A, B, C (see Figures S22, S23, S24, respectively), we notice that unique segregation in this space correlates with the RF model accuracy. Nests that form closely spaced homogenized clusters in the 2D manifold tend to be correlated with models that scored high recall for these nests.



**Figure 2.** Performance of RF models designed to classify nests within different sets. A) Travel times  $t_{sp}$  for four lunar model from Garcia et al. (2011) (G11), Garcia et al. (2019) (issi2, issi3), (Khan et al., 2014) (K14) with nest labels. B) Recall for individual nests within each set with respect to their travel times labeled with sets to which they belong and the scored accuracy of this set. C) and D) same as B) just for different group of sets.

\$ <sup>3</sup> -	0.13	0.04	0.13	0.18	0.11	0.15	0.15	0.12		
Sr -	0.15	0.04	0.14	0.17	0.07	0.18	0.16	0.08		- 0.225
දිා -	0.24	0.04	0.12	0.17	0.1	0.11	0.13	0.08		- 0 200
5 <sup>A</sup> -	0.24	0.05	0.12	0.13	0.11	0.13	0.12	0.09		0.200
52-	0.12	0.05	0.14	0.13	0.16	0.15	0.12	0.13		- 0.175
53-	0.08	0.07	0.09	0.09	0.24	0.15	0.09	0.19		
51A -	0.18	0.07	0.09	0.15	0.14	0.11	0.12	0.12		- 0.150
Sets -	0.12	0.03	0.14	0.17	0.13	0.12	0.17	0.12		
్ హ -	0.16	0.07	0.11	0.12	0.13	0.17	0.12	0.12		- 0.125
- OS	0.15	0.06	0.09	0.12	0.14	0.21	0.11	0.12		0 1 0 0
51 -	0.16	0.05	0.09	0.12	0.15	0.19	0.11	0.13		- 0.100
- °ې	0.14	0.04	0.1	0.14	0.17	0.14	0.12	0.14		- 0.075
දුව -	0.16	0.07	0.11	0.12	0.12	0.17	0.12	0.12		
570 -	0.16	0.07	0.11	0.12	0.13	0.19	0.12	0.11		- 0.050
51-	0.16	0.07	0.12	0.12	0.12	0.18	0.12	0.11		
	$\Delta t_A$	$\Delta t_F$	$\Delta t_P$	$\cos(\gamma)$	$sin(\gamma)$	$\Delta e$	ď	ģ		

Figure 3. Feature importance for Random Forest models associated with different travel time sets shown in Figure 1.

# **5** Conclusion

In this paper we propose how to tackle DMQ nest identification during future lu-314 nar missions that will likely host only one station on the far side of the Moon. We pro-315 pose constraining their location by using differential time travel measurement  $t_{sp}$  and 316 parameters related to the temporal patterns of the DMQ occurrence. First, in our anal-317 ysis we assume that we cannot differentiate between nests whose differences in travel time 318 are less than 5 seconds. Thus, we form set of nests that have similar travel times. Sec-319 ond, for each event within the nests we calculate features that are used to build a Ran-320 321 dom Forest model. This model is trained to differentiate between distinct nests. The features used for training are build by associating each event in all nests with the time dif-322 ference between events' origin time and time of lunar ascending node, Full Moon phase, 323 perigee, then position of the Moon in its orbit expressed by true anomaly angle, distance 324 of the Moon from the Earth, rate change of this distance, and the time between two suc-325 cessive quakes. We show that by training Random Forest models to differentiate between 326 distinct nests within sets, we can obtain models with high accuracy (more than half of 327 the models score above 70% accuracy). Yet, the performances of these models depend 328 on the number of nests within the set. More nests implies that the problem is more dif-329 ficult to solve, probably because a) nests might have similar source mechanisms, b) the 330 number of events within nests is unbalanced, and c) we don't have enough data. Since 331 RF models also arrange features by their importance to make a final classification de-332 cision, we observe that the importance of the features change with different sets. This 333 complements the findings of previous papers, since it signifies that nests do correspond 334 to different lunar events, which eventually might be connected to the distribution of tidal 335 stresses during these events. Finally, our model provides a good first approximation of 336 the nest identification. And as the catalog of new events grows, it will be straightforward 337 to retrain RF model with the new enlarged dataset. 338

# 339 Open Research Section

The deep moonquake catalog used in this study is published in Nakamura et al. (1981), and revisited in Nunn et al. (2020). Python package Skyfield used to calculate Moon's orbital parameters based on JPL ephemeris can be found on the website https:// rhodesmill.org/skyfield/ (Rhodes, 2019, Software). For our implementation of the Random Forest algorithm we use Scikit-learn machine learning Python library (Pedregosa et al., 2011).

# 346 Acknowledgments

French co-authors thanks the French Space Agency, CNES, for supporting this research in the frame of the French contribution to FSS as well as IDEX Paris Cité (ANR-18-IDEX-0001). MPP was supported by funds from the Jet Propulsion Laboratory, under contract with the National Aeronautics and Space Administration (NASA).

# 351 References

- Banerdt, W. B., Smrekar, S. E., Banfield, D., Giardini, D., Golombek, M., Johnson,
   C. L., ... others (2020). Initial results from the insight mission on mars.
   *Nature Geoscience*, 13(3), 183–189.
- Bills, B. G., Bulow, R., & Johnson, C. L. (2008). Influence of earth-moon orbit geometry on deep moonquake occurrence times. *Lunar Planet. Sci.*, XXXIX.
- Breiman, L. (1996). Bagging predictors. *Machine learning*, 24, 123–140.
- Breiman, L. (2001). Random forests. *Machine learning*, 45, 5–32.
- Breiman, L., Friedman, J., Stone, C. J., & Olshen, R. A. (1984). Classification and regression trees. CRC press.

361	Bulow, R., Johnson, C., Bills, B., & Shearer, P. (2007). Temporal and spatial prop-
362	erties of some deep moonquake clusters. Journal of Geophysical Research: Blue et = 110 (EQ)
363	Planets, 112 (E9).
364	bulow, R., Johnson, C., & Snearer, P. (2005). New events discovered in the apolio lunar seismic data. <i>Lournal of Ceonhusical Research: Planets</i> 110(E10)
305	Compared W S & Cilhoppy I I (1067) Operation mean blink and report of ab
366 367	cameron, W. S., & Gineany, J. J. (1967). Operation moon blink and report of ob- servations of lunar transient phenomena. <i>Icarus</i> , $\gamma(1-3)$ , 29–41.
368	Cheng, C. H., & Toksöz, M. N. (1978). Tidal stresses in the moon. Journal of
369	Geophysical Research: Solid Earth, 83(B2), 845-853. Retrieved from https://
370	agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JB083iB02p00845
371	doi: https://doi.org/10.1029/JB083iB02p00845
372	Cutler, J. W., Nguyen, T. A., Kano, T., Kumar, Y. R. L., Panning, M., April, S., &
373	Haque, S. (2023). Overview of the avionics design for the farside seismic suite.
374	In Aiaa scitech 2023 forum. Retrieved from https://arc.aiaa.org/doi/abs/
375	10.2514/6.2023-1880 doi: 10.2514/6.2023-1880
376	Dainty, A., Goins, N., & Toksoz, M. (1975). Natural lunar seismic events and the
377	structure of the moon. In In: Lunar science conference, 6th, houston, tex.,
378	march 17-21, 1975, proceedings. volume 3.(a78-46741 21-91) new york, perga-
379	mon press, inc., 1975, p. 2887-2897. (Vol. 6, pp. 2887–2897).
380	Ewing, M., Latham, G., Press, F., Sutton, G., Dorman, J., Nakamura, Y., Ko-
381	vach, R. (1971). Seismology of the moon and implications on internal struc-
382	ture, origin and evolution. <i>Highlights of astronomy</i> , 2, 155–172.
383	Frohlich, C., & Nakamura, Y. (2009). The physical mechanisms of deep moon-
384	quakes and intermediate-depth earthquakes: How similar and how differ-
385	ent? Physics of the Earth and Planetary Interiors, 173(3), 365-374. Re-
386	trieved from https://www.sciencedirect.com/science/article/pii/
387	S0031920109000338 doi: https://doi.org/10.1016/j.pepi.2009.02.004
388	Gagnepain-Beyneix, J., Lognonné, P., Chenet, H., Lombardi, D., & Spohn, T.
389	(2006). A seismic model of the lunar mantle and constraints on tempera-
390	ture and mineralogy. Physics of the Earth and Planetary Interiors, 159(3-4),
391	140-100. Carrie P. F. Carnensin Permeire I. Charnet S. & Lornenné P. (2011). Voru pro-
392	liminary reference moon model Physics of the Earth and Planetary Interiors
393	188(1-2), 96-113.
395	Garcia, R. F., Khan, A., Drilleau, M., Margerin, L., Kawamura, T., Sun, D., oth-
396	ers (2019). Lunar seismology: An update on interior structure models. Space
397	Science Reviews, 215, 1–47.
398	Ho, T. K. (1995). Random decision forests. In Proceedings of 3rd international con-
399	ference on document analysis and recognition (Vol. 1, pp. 278–282).
400	Kawamura, T., Grott, M., Garcia, R., Wieczorek, M., de Raucourt, S., Lognonné,
401	P., others (2022). An autonomous lunar geophysical experiment package
402	(algep) for future space missions: In response to call for white papers for the
403	voyage 2050 long-term plan in the esa science program. Experimental Astron-
404	omy, 54 (2-3), 617-640.
405	Kawamura, T., Lognonné, P., Nishikawa, Y., & Tanaka, S. (2017). Evaluation of
406	deep moonquake source parameters: Implication for fault characteristics and
407	thermal state. Journal of Geophysical Research: Planets, 122(7), 1487–1504.
408	Khan, A., Connolly, J. A. D., Pommier, A., & Noir, J. (2014). Geophysical evi-
409	dence for melt in the deep lunar interior and implications for lunar evolution.
410	Journal of Geophysical Research: Planets, 119(10), 2197-2221. Retrieved
411	IFOIII nttps://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2014 IE004661 doi: https://doi.org/10.1002/2014 IE004661
412	Z014JE004001 uoi: https://uoi.org/10.1002/2014JE004001 Khan A fr Moreorened K (2002) An inquiry into the lunger interior. A gravity
413	ear inversion of the apollo lunar seismic data. <i>Lowrad of Coonductional Descented</i> .
414	Planets $107(\text{F6})$ 3–1
410	1 0010000, 107(110), 0 1.

- Khan, A., Mosegaard, K., & Rasmussen, K. L. (2000). A new seismic velocity model 416 for the moon from a monte carlo inversion of the apollo lunar seismic data. 417 Geophysical Research Letters, 27(11), 1591–1594. 418 Lammlein, D. R. (1977a). Lunar seismicity and tectonics. *Physics of the* 419 Earth and Planetary Interiors, 14(3), 224-273. Retrieved from https:// 420 www.sciencedirect.com/science/article/pii/0031920177901753 doi: 421 https://doi.org/10.1016/0031-9201(77)90175-3 422 Lammlein, D. R. (1977b). Lunar seismicity and tectonics. *Physics of the Earth and* 423 *Planetary Interiors*, *14*(3), 224–273. 424 Lammlein, D. R., Latham, G. V., Dorman, J., Nakamura, Y., & Ewing, M. (1974).425 Lunar seismicity, structure, and tectonics. Reviews of Geophysics, 12(1), 1-21. 426 Langer, H., Falsaperla, S., & Hammer, C. (2019). Advantages and pitfalls of pattern 427 recognition: selected cases in geophysics. Elsevier. 428 Lognonné, P., Banerdt, W. B., Giardini, D., Pike, W. T., Christensen, U., Laudet, 429 P., ... others (2019). Seis: Insight's seismic experiment for internal structure 430 of mars. Space Science Reviews, 215, 1-170. 431 Lognonné, P., Banerdt, W. B., Pike, W. T., Giardini, D., Christensen, U., Garcia, 432 R. F., ... others (2020). Constraints on the shallow elastic and anelastic struc-433 ture of mars from insight seismic data. Nature Geoscience, 13(3), 213–220. 434 Lognonné, P., Gagnepain-Beyneix, J., & Chenet, H. (2003). A new seismic model of 435 the moon: implications for structure, thermal evolution and formation of the 436 moon. Earth and Planetary Science Letters, 211(1-2), 27–44. 437 Lognonné, P., & Johnson, C. (2015). 10.03—planetary seismology. Treatise on geo-438 physics, 2, 65–120. 439 Meeus, J. (1991). Astronomical algorithms. *Richmond*. 440 Middlehurst, B. M. (1967). An analysis of lunar events. Reviews of Geophysics, 441 5(2), 173-189.442 Minshull, T., & Goulty, N. (1988a). The influence of tidal stresses on deep moon-443 quake activity. Physics of the Earth and Planetary Interiors, 52(1), 41-55. 444 Retrieved from https://www.sciencedirect.com/science/article/pii/ 445 0031920188900568 doi: https://doi.org/10.1016/0031-9201(88)90056-8 446 Minshull, T., & Goulty, N. (1988b). The influence of tidal stresses on deep moon-447 quake activity. Physics of the earth and planetary interiors, 52(1-2), 41–55. 448 Moore, P. (1968). Transient lunar phenomena: A review, 1967. J. Br. Astron. As-449 soc., 78, 138-144. 450 Nakamura, Y. (1983). Seismic velocity structure of the lunar mantle. Journal of 451 Geophysical Research: Solid Earth, 88(B1), 677-686. Retrieved from https:// 452 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JB088iB01p00677 453 doi: https://doi.org/10.1029/JB088iB01p00677 454 Nakamura, Y. (2003).New identification of deep moonquakes in the apollo lunar 455 seismic data. Physics of the Earth and Planetary Interiors, 139(3), 197-205. 456 Retrieved from https://www.sciencedirect.com/science/article/pii/ 457 S0031920103002103 doi: https://doi.org/10.1016/j.pepi.2003.07.017 458 Farside deep moonquakes and deep interior of the moon. Nakamura, Y. (2005).459 Journal of Geophysical Research: Planets, 110(E1). Retrieved from https:// 460 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2004JE002332 doi: 461 https://doi.org/10.1029/2004JE002332 462 Nakamura, Y., Latham, G., Dorman, H., & Harris, J. (1981). Passive seismic ex-463 Galveston Geophysics Laboratory Contriperiment long-period event catalog. 464 bution, 491, 314. 465 Nakamura, Y., Latham, G. V., & Dorman, H. J. (1982a). Apollo lunar seismic 466 experiment—final summary. Journal of Geophysical Research: Solid Earth, 467 87(S01), A117–A123. 468 Nakamura, Y., Latham, G. V., & Dorman, H. J. (1982b). Apollo lunar seismic 469
- 470 experiment—final summary. Journal of Geophysical Research: Solid Earth,

471	87(S01), A117-A123. Retrieved from https://agupubs.onlinelibrary
472	.wiley.com/doi/abs/10.1029/JB087iS01p0A117 doi: https://doi.org/
473	10.1029/JB087iS01p0A117
474	Nunn, C., Garcia, R. F., Nakamura, Y., Marusiak, A. G., Kawamura, T., Sun, D.,
475	others (2020). Lunar seismology: A data and instrumentation review.
476	Space Science Reviews, 216(5), 89.
477	Panning, M., Kedar, S., Bowles, N., Calcutt, S., Cutler, J., Elliott, J., others
478	(2021). Farside seismic suite (fss): First seismic data from the farside of the
479	moon delivered by a commercial lander. In Agu fall meeting abstracts (Vol.
480	2021, pp. P54C–01).
481	Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,
482	others (2011). Scikit-learn: Machine learning in python. the Journal of
483	machine Learning research, 12, 2825–2830.
484	Pickrell, J. (2022, May). These six countries are about to go to the moon — here's
485	why. Nature, 605, 208–211. doi: 10.1038/d41586-022-01252-7
486	Quinlan, J. R. (1986). Induction of decision trees. Machine learning, 1, 81–106.
487	Rhodes, B. (2019). Skyfield: Generate high precision research-grade positions for
488	stars, planets, moons, and earth satellites. Astrophysics Source Code Library,
489	ascl-1907.
490	Standley, I. M., Pike, W. T., Calcutt, S., & Hoffman, J. P. (2023). Short period seis-
491	mometer for the lunar farside seismic suite mission. In 2023 ieee aerospace con-
492	ference (p. 1-9). doi: 10.1109/AERO55745.2023.10115559
493	Toksoz, M. N., Dainty, A. M., Solomon, S. C., & Anderson, K. R. (1974). Structure
494	of the moon. <i>Reviews of Geophysics</i> , 12(4), 539-567. Retrieved from https://
495	agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/RG0121004p00539
496	doi: https://doi.org/10.1029/RG0121004p00539
497	lotion to tidal strongen. <i>Science</i> 106(4202) 070 081
498	Turner A B Hawtherne I C & Caddes M (2022) Stresses in the lunar
499	interior: Insights from slip directions in the 201 doop moonquelo nest
500	nd of Geonbusical Research: Planets 197(12) a2022 IE007364 Betrieved
501	from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/
502	2022.IE007364 doi: https://doi.org/10.1029/2022.IE007364
504	van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-sne. Journal
505	of Machine Learning Research, 9(86), 2579–2605. Retrieved from http://imlr
506	.org/papers/v9/vandermaaten08a.html
507	Weber, R., Bills, B., & Johnson, C. (2009). Constraints on deep moonquake focal
508	mechanisms through analyses of tidal stress. <i>Journal of Geophysical Research</i> :
509	<i>Planets</i> , 114 (E5).
510	Weber, R., Bills, B., & Johnson, C. (2010). A simple physical model for deep moon-
511	quake occurrence times. Physics of the Earth and Planetary Interiors, 182(3-
512	4), 152–160.
513	Weber, R., Neal, C., Grimm, R., Grott, M., Schmerr, N., Wieczorek, M., oth-
514	ers (2021). The scientific rationale for deployment of a long-lived geophysical
515	network on the moon. Bulletin of the AAS, $53(4)$ .
516	Witze, A. (2022, May). The \$93-billion plan to put astronauts back on the moon.

<sup>517</sup> Nature, 605, 212–216. doi: 10.1038/d41586-022-01253-6

-14-

# Supporting Information for "Identifying deep moonquake nests using machine learning model on single lunar station on the far side of the Moon"

Josipa Majstorović<sup>1</sup>, Philippe Lognonné<sup>1</sup>, Taichi Kawamura<sup>1</sup>, Mark Panning<sup>2</sup>

<sup>1</sup>Université Paris Cité, Institut de physique du globe de Paris, CNRS, Paris, France

<sup>2</sup>Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA

# Contents of this file

- 1. Text S1 to S3  $\,$
- 2. Figures S1 to S24

This file contains supplement text, and figures for the manuscript "Identifying deep moonquake nests using machine learning model on single lunar station on the far side of the Moon". Supplemental Text S1 provides a detailed description of P- and S- travel times calculating for deep moonquake nests using several lunar models. Supplemental Text S2 details about the Random Forest algorithm and the particularity of its training and testing. Supplemental Text S3 provides additional information on the t-distributed stochastic neighbor embedding (t-sne) method. Supplemental Figure S1 show travel time  $t_{sp} = t_s - t_p$  for around 20 nests and seven existing lunar models. Supplemental Figure

Corresponding author: J. Majstorović, Université Paris Cité, Institut de physique du globe de Paris, CNRS, Paris, France (josipa.majstorovic@protonmail.com)

X - 2 MAJSTOROVIĆ, J. ET. AL.: IDENTIFYING DEEP MOONQUAKES WITH ML

S4 illustrates distributions of features used to train machine learning model for two nests, A1 and A8. Supplemental Figure S2 shows features used as an input data for machine learning model. Supplemental Figures S5 to S13 show statistics of the Random Forest models on the test dataset that are trained to dissociate between A1 and A8 nest. Supplemental Figure S13 shows prediction values and negatively labelled data points for the base Random Forest model trained to dissociate between A1 and A8 nests. Supplemental Figure S14 and S15 illustrate 2-D graphic manifold of the space spanned by the 8-D feature space of A1 and A8 nests. Supplemental Figures S16 to S19 show statistic of the Random Forest model trained and tested on combination of nests A1-A8-A18, A1-A8-A18-A6, A1-A8-A18-A6-A14, A8-A18-A6-A14, respectively. Supplemental Figure S20 illustrates 2D graphic manifold representation of the 8D feature space spanned by five nests A1-A8-A18-A6-A14. Supplemental Figure S21 show nests ratios for different sets. Supplemental Figures S22 to S24 illustrates 2-D manifolds of the feature spanned by different sets.

**Text S1.** Calculation of travel times of seismic waves between sources and station can be obtained using two programming packages: Python package Obspy and its module 'taup' and TauP Java package, both based on the paper Crotwell, Owens, Ritsema, et al. (1999). In our experiment we placed single station on the far side of the Moon in the Schrödinger Crater, while our sources are located nests from the paper Lognonné, Gagnepain-Beyneix, and Chenet (2003). Also, we utilize the existing lunar interior models from papers Garcia et al. (2019) (ISSI 1, ISSI 2, ISSI 3), Garcia, Gagnepain-Beyneix, Chevrot, and Lognonné (2011), Khan, Connolly, Pommier, and Noir (2014), Matsumoto et al. (2015), Weber, Lin, Garnero, Williams, and Lognonné (2011). We first calculate epicentral distances between

nests and station, then travel time of P- and S- seismic waves for the seven models using Python and Java packages. This leaves us with:  $t'_{p;i,j}, t''_{s;i,j}, t''_{p;i,j}, t''_{s;i,j}$  where *i* indicates nest, *j* indicates lunar model, *t'* and *t''* travel times calculated using Python and Java, respectively. Next, we calculate the average over P- and S- wave travel times for two programming packages, leaving us with  $t_{p;i,j} = (t'_{p;i,j} + t''_{p;i,j})/2$ ,  $t_{s;i,j} = (t'_{s;i,j} + t''_{s;i,j})/2$ . Further, we calculate the travel time difference,  $t_{sp;i,j} = t_{s;i,j} - t_{p;i,j}$  values that we plot in Figure S1, with *i* running over the y-axis for all nests and *j* running over the x-axis over all lunar models. We can notice that some combination of nests and lunar models don't have travel time  $t_{sp}$  estimation, and some underestimate or overestimate it, when compared to the average value per nest. Due to these discrepancies we decide to further work with only four models: ISSI 2, ISSI 3 Garcia et al. (2019), Garcia PEPI 2011 (Garcia et al., 2011), Khan JGR 2014 (Khan et al., 2014).

Text S2. As discussed in the main manuscript, Random Forest is a machine learning algorithm that consist of ensemble of randomised decision trees. A decision tree consists of decision (internal) nodes, followed by inequality branches, and leaf nodes that hold the final prediction of individual trees shown in Figure S3. Thus, within each tree the beginning is at root node that doesn't have incoming branches. Next in line are internal nodes where based on the available features/attributes and inequality operations, the decision whether the feature is smaller or larger than some threshold is made. These translate to leaf nodes, which represent all possible outcomes. The hyperparamters that define a RF structure and need to be defined before a training process are: the number of decision trees, the maximum depth of trees, the measure that maximises diversity between classes, the minimum samples in the internal node, and the minimum number of samples in leaf

node for it to be considered, the maximum number of features when looking for the best split, the maximum number of leaf nodes, the maximum samples to be draw from the main training dataset when training each decision tree. We proceed to test the learning robustness of our base model that is trained with normalised features (shown in Figure S6) by carry out several experiments: a) changing the randomness of the bootstrapping initialization of the samples that are used when building decision trees, the randomness of the feature sampling when considering for the best split at each internal node, as well as the randomness of the training and test dataset split (see Figure S7); b) changing the optimal number of decision trees (see Figure S8); c) equalizing the size of nests within the dataset by randomly downsampling the largest nest A1 to be the same size as A8 (see Figure S9); d) reducing the number of input feature data to five most important from the base model  $(cos(\gamma), \Delta t_{Perigee}, d, \Delta t_{AscendingNode}, e_{i+1} - e_i)$  (see Figure S10). In all test beside those in experiment a), we keep the random state fixed. Finally, the results do not vary between different tests, indicating that in all above configurations models are able to learn how to classify two nest with the similar performances. Further, we notice that in experiment c) the statistic for A8 nest improved compared to the base model statistic, indicating that having a balanced classes while training a ML is important. Moreover, we observe that the model trained with the fewer features statistically perform worst than the base model. This might be because all eight features are uncorrelated, thus equally important for model learning. Next, we cross-validate our base model. A cross-validation is a technique to assess how the model will generalize to an independent data set by using the resampling technique. A resampling technique uses different portions of the training data to train and validate model during several iteration. Usually, the training

dataset is divided into k equally sized folds, and then k iterations is performed. In each iteration (k-1) folds are used for training, and one fold is used for validation. During the cross-validation, the test set is kept aside. Eventually, the full dataset is divided into three sets: training (55%), validation (20%), test (25%), where training and validation set change in each iteration. We calculate cross-validation with k = 5, while keeping the base model parameters. The choice of k = 5, has been proven to be a good practice (Witten & James, 2013). This test produces 5 models with the same performance as indicated with ROC curves (see Figure S11). This implies that the base model generalizes well, and to not over fit the results. Moreover, we also perform the grid search over several other RF parameters, besides the number of decision trees. Grid search represents a set of many models, where each model is build with unique set of parameters, and each is trained and tested with the same datasets. The tested parameters are: the maximum features ('auto', 'sqrt'), the maximum depth of the decision trees (10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110, None), the minimum samples in the internal nodes (2, 5, 10), the minimum number of samples in leaf node (1, 2, 4), the bootstraping technique on or off. However, it seems that the model that performed the best during the grid search (the minimum samples in internal node equal to 10, the minimum samples in leaf node equal to 2, the maximum features equal to 'auto', the maximum depth equal to 110, with bootstraping turned on) does not perform better than our base model (see Figure S12). Even though our base RF model as a final output predicts a class (A1 or A8), it also associates each class prediction with the class probability. This value ranges between 0 and 1, where 1 indicated that a model is absolutely certain that a given event belongs to a predicted class. From our base model the correctly predicated classes score 81% cases higher than

X - 6 MAJSTOROVIĆ, J. ET. AL.: IDENTIFYING DEEP MOONQUAKES WITH ML

0.80 for test dataset (see Figure S13A). This suggest that model is confident in its prediction. The mislabelled prediction values show uniform distribution of values between 0.5 to 1 (see Figure S13B). We can further examine these mislabelled events from the perspective of the feature input data. We choose three features,  $cos(\gamma)$ ,  $\Delta t_{Perigee}$ , d, with high importance value. We could argue that the mislabelled data points from both nests show characteristics which better suits the opposite class (see Figure S13C-E). Yet, the decision is not defined using only one feature. To gain an insight how all eight features contribute into separating two nests, we calculate the 2D manifold of their feature space using t-sne method (van der Maaten & Hinton, 2008). Visualisation of the feature space that is color-coded based on two nests, show us that two classes are well but not perfectly separated (see Figure S14). If we further color 2D manifold space with the values of the individual features, we notice that  $\Delta t_{AscendingNode}$  feature might be the most responsible for imperfect split between nests (see Figure S15).

**Text S3.** One of the statistical dimensionality reduction algorithm that helps to visual high-dimensional data is t-distributed Stochastic Neighbor Embedding (t-sne) algorithm (van der Maaten & Hinton, 2008). It is an unsupervised non-linear reduction technique, since it allows us to separate data that cannot be separated by any straight line. Once it is applied on the input data, first, it starts by calculating the probability distribution of neighbours around each points. The term neighbour stands for the set of points that are closest to a given point. In the input original high-dimensional space the probability is modeled as a Gaussian distribution. Second, the algorithm models the probability of neighbours around given points in the lower-dimensional space using a Student's t-distribution. Third, the algorithm minimizes the divergence, usually Kullback-Leibler

divergence, between two probabilities using gradient descent. The result is a lowerdimensional manifold of the data, that still preserves the pairwise similarities between original data points, optimized to a stable state. This optimisation process generates clusters and sub-clusters of similar data points that become visually better understand in the lower-dimensional space by keeping the relationship of the data from the higherdimensional space. There are several t-sne hyperparameters that need to be adjusted by the user, and the most important one is perplexity. The preplexity parameter defines the number of influential neighbours used to calculate the Gaussian probabilities around given point in the high-dimensional space. Its value range from 5 - 50 (Wattenberg et al., 2016), and can significantly impact the resulting mapping of the input data. For our implementation of the t-sne algorithm we use Scikit-learn machine learning Python library (Pedregosa et al., 2011). After trying some combinations we choose to work with the given set of parameters: n\_components=2, perplexity=30, n\_iter=5000, verbose=1, random\_state=133, while keeping the rest of the them as default by the package implementation.

# References

X - 8

- Crotwell, H. P., Owens, T. J., Ritsema, J., et al. (1999). The taup toolkit: Flexible seismic travel-time and ray-path utilities. *Seismological Research Letters*, 70, 154–160.
- Garcia, R. F., Gagnepain-Beyneix, J., Chevrot, S., & Lognonné, P. (2011). Very preliminary reference moon model. *Physics of the Earth and Planetary Interiors*, 188(1-2), 96–113.
- Garcia, R. F., Khan, A., Drilleau, M., Margerin, L., Kawamura, T., Sun, D., ... others (2019). Lunar seismology: An update on interior structure models. Space Science Reviews, 215, 1–47.
- Khan, A., Connolly, J. A. D., Pommier, A., & Noir, J. (2014). Geophysical evidence for melt in the deep lunar interior and implications for lunar evolution. *Journal of Geophysical Research: Planets*, 119(10), 2197-2221. Retrieved from https://agupubs.onlinelibrary.wiley .com/doi/abs/10.1002/2014JE004661 doi: https://doi.org/10.1002/2014JE004661
- Lognonné, P., Gagnepain-Beyneix, J., & Chenet, H. (2003). A new seismic model of the moon: implications for structure, thermal evolution and formation of the moon. *Earth and Planetary Science Letters*, 211(1-2), 27–44.
- Matsumoto, K., Yamada, R., Kikuchi, F., Kamata, S., Ishihara, Y., Iwata, T., ...
  Sasaki, S. (2015). Internal structure of the moon inferred from apollo seismic data and selenodetic data from grail and llr. *Geophysical Research Letters*, 42(18), 7351-7358. Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2015GL065335
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... others (2011). Scikit-learn: Machine learning in python. the Journal of machine Learning research, 12, 2825–2830.

- van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-sne. Journal of Machine Learning Research, 9(86), 2579-2605. Retrieved from http://jmlr.org/papers/v9/ vandermaaten08a.html
- Wattenberg, M., Viégas, F., & Johnson, I. (2016). How to use t-sne effectively. *Distill*. Retrieved from http://distill.pub/2016/misread-tsne doi: 10.23915/distill.00002
- Weber, R. C., Lin, P.-Y., Garnero, E. J., Williams, Q., & Lognonné, P. (2011). Seismic detection of the lunar core. *Science*, 331 (6015), 309-312. Retrieved from https://www.science.org/ doi/abs/10.1126/science.1199375 doi: 10.1126/science.1199375
- Witten, D., & James, G. (2013). An introduction to statistical learning with applications in r. springer publication.

	100.0	105.0	101 5	100.0	100 7	1070	200.2	
A01 -	186.8	185.3	191.5	186.9	190.7	187.9	200.3	
A06 -	244.0	229.6	241.3	233.3	240.7	542.6		- 500
A07 -	212.8	207.3	215.9	209.8	214.9	211.5	226.1	500
A08 -	172.0	171.2	225.4	172.6	175.7	173.2	184.2	
A09 -	159.4	158.6	163.1	159.9	162.5	160.1	170.0	
A11 -	379.2	187.0	195.8	190.1	196.8	253.6		
A14 -	172.0	172.5	226.9	173.8	176.8	174.5	185.3	- 400
A16 -	381.4	192.6	201.6	195.9	201.1	454.3	210.9	
A17 -	241.0	301.9	237.3	229.8	236.5	532.7	544.8	
A18 -	214.5	209.5	218.1	211.9	217.1	213.7	228.5	
A20 -	361.6	212.3	222.9	215.1	223.3	17.7		t t
អ្ហ៍ A24 -	161.1	160.0	164.7	161.3	164.0	161.7	171.7	- 300 I
≝ A25 -	217.8	214.6	224.4	217.7	223.4	283.5	235.1	= t
A27 -	371.9	204.2	214.5	207.9	214.1	483.6		sp
A30 -	216.8	213.9	224.2	217.2	223.1	219.3	234.6	-
A33 -	171.4	171.7	176.6	173.1	176.0	173.8	184.5	- 200
A34 -	210.3	206.7	215.8	209.3	214.6	211.0	225.7	200
A40 -	204.6	201.7	209.2	203.6	208.3	205.2	219.1	
A41 -	83.8	212.7	223.0	216.2	222.3	500.7	233.6	
A42 -	364.7	216.1	227.0	219.4	227.1	269.4		
A44 -	363.1	223.4	234.7	226.4	234.8	278.6		- 100
A50 -	223.8	218.3	227.5	221.0	226.6	223.1	238.5	
A51 -	209.8	205.6	213.8	207.8	212.8	209.5	223.9	
A97 -	189.8	186.5	193.6	188.4	192.6	189.6	202.3	
	1	2	ň		4	ں۔ '		 -
		≥_		101	101	101	102	
	NO	NO	NO					
	Q	Q	Q	EP	5	GRI	DC.	
	<u>_</u>	<u> </u>	<u>_</u> I	E E	Ē	。'	cie.	
	ISS	ISS	ISS	<u>.</u>	(ha	lot	Ň	
				Ga	×	ш	Jer	
				-		ats	Vek	
						Ĕ	5	
			L	unar mode	ls			

Figure S1. Travel time  $t_{sp} = t_s - t_p$  calculations, where  $t_s$  and  $t_p$  stands for S- and P- waves travel times, respectively, between the range of nests and the station placed on the far side of the Moon in the Schrödinger crater. Calculation are done for seven existing lunar model, from papers Garcia et al. (2019) (ISSI M1, ISSI M2, ISSI M3), Garcia et al. (2011), Khan et al. (2014), Matsumoto et al. (2015), Weber et al. (2011).



Figure S2. Time evolution of features during the Apollo mission and their histograms for A1 nest. First three features are time difference between quakes and A) the instance when Moon was passing through ascending node, B) Full Moon phase, C) instance when Moon was in perigee; next D)  $\cos \gamma$ , E)  $\sin \gamma$  where  $\gamma$  is the true anomaly angle, indicating the position of the Moon in the orbit; F) time difference between two quakes, G) distance between Moon and Earth at the quake occurrence, H) the rate of distance change from G.



**Figure S3.** Schematic representation of a Random Forest algorithm. In this example, the model is trained with 3 features, it consists of 3 decision trees with a maximum tree depth equal to 3. A tree consists of decision nodes (circles), followed by inequality branches (dashed lines), and leaf nodes (rectangles). The prediction is taking place in each tree by yes or no questions, while the final prediction is made upon majority voting considering individual tree predictions.



Figure S4. Distributions of eight features used for training Random Forest model for dissociating between nests A1 and A8.



**Figure S5.** Statistics performance of the base Random Forest model on the dataset trained to dissociate between nests A1 and A8 using raw feature data without normalisation: A) confusion matrix, B) precision, recall, f1-score per nests and accuracy of the model, C) receiver operating characteristic (ROC) curve, D) feature importance for the model to make decisions.



Figure S6. Same as Figure S5 but using feature data that are normalised between 0 and 1.



**Figure S7.** Statistics performance of the Random Forest models while changing the randomness of the data split and decision tree initialisation compared to the base model shown in Figure S6. The models are trained to dissociated between nests A1 and A8. The randomness is changed from 600 to 1400 from test 1 to test 10 by step of 100. Statistics are: A) confusion matrix, B) precision, recal, f1-score and accuracy of the model, C) the importance of the feature used by the models to make a correct classification.



**Figure S8.** Same as Figure S7 while keeping the randomness fixed, but changing the number of trees used to build Random Forest model.



Figure S9. Same as Figure S6, but using the balanced dataset, thus having the same number of A1 and A8 events.



Figure S10. Same as Figure S6, but keeping only five out of eight features:  $\Delta t_{AscendingNode}$ ,  $\Delta t_{Perigee}, \cos(\gamma), e_{i+1} - e_i, d.$ 



**Figure S11.** 5-fold cross-validation of the base RF model shown in Figure S6 with the mean and standard deviation.



Figure S12. Same as Figure S6, but for the best performing model from the grid search analysis.



Figure S13. Classification results: A) prediction values for positively classified events, B) prediction values for negatively classified events. Events used for training and events used for testing but got mislabeled form the perspective of features: C)  $\cos(\gamma)$ , D)  $\Delta t_{Perigee}$ , E) d.



Figure S14. 2-D manifold of the feature space spanned by nests A1 and A8.



Figure S15. 2-D manifold of the feature space spanned by nests A1 and A8 colored by the features: A)  $\Delta t_{AscendingNode}$ , B)  $\Delta t_{FullMoon}$ , C)  $\Delta t_{Perigee}$ , D)  $\cos(\gamma)$ , E)  $\sin(\gamma)$ , F)  $e_{i+1} - e_i$ , G) d, H)  $\dot{d}$ .



Figure S16. Same as Figure S6, but classifying three nests A1, A8 and A18.



Figure S17. Same as Figure S6, but classifying four nests A1, A8, A18, A6.

X - 26



Figure S18. Same as Figure S6, but classifying five nests A1, A8, A18, A6, A14.



Figure S19. Same as Figure S6, but classifying four nests A8, A18, A6, A14.



Figure S20. 2-D manifold of the feature space spanned by nests A1, A8, A18, A6, A14.



**Figure S21.** Pie charts for all sets shown in Figure 1 displaying the composition of the set and the contribution of each nest within the each set.



Figure S22. 2-D manifold of the feature space spanned by nests belonging to defined sets S<sub>i</sub>:
A) S1, B) S2, C) S3, D) S4, E) S12, F) S13, G) S14, H) S15.



Figure S23. 2-D manifold of the feature space spanned by nests belonging to defined sets S<sub>i</sub>:
A) S5, B) S6, C) S7, D) S8



Figure S24. 2-D manifold of the feature space spanned by nests belonging to defined sets S<sub>i</sub>:A) S9, B) S10, C) S11.