

# A Logit-Based Binary Classifier of Tsunamigenic Earthquakes for the Northern Pacific Ocean

Alexey Konovalov (✉ [a.konovalov@geophystech.ru](mailto:a.konovalov@geophystech.ru))

Far Eastern Geological Institute FEB RAS: Dal'nevostocnyj geologiceskij institut Dal'nevostocnogo otdelenia Rossijskoj akademii nauk <https://orcid.org/0000-0003-2997-1524>

Grigory Samsonov

Sakhalin State University

Andrey Stepnov

Far Eastern Geological Institute FEB RAS: Dal'nevostocnyj geologiceskij institut Dal'nevostocnogo otdelenia Rossijskoj akademii nauk

---

## Research Article

**Keywords:** submarine earthquake, tsunami, logit, predictive model, binary classifier, early warning

**Posted Date:** March 15th, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1421995/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **A Logit-Based Binary Classifier of Tsunamigenic Earthquakes for the Northern Pacific**  
2 **Ocean**

3 A.V. Konovalov<sup>a\*</sup>, G.A. Samsonov<sup>b</sup>, A.A. Stepnov<sup>a</sup>

4 <sup>a</sup>Far East Geological Institute, Far Eastern Branch, Russian Academy of Sciences, Vladivostok, Russia

5 <sup>b</sup>Sakhalin State University, Yuzhno-Sakhalinsk, Russia

6 \*e-mail: a.konovalov@geophys.tech.ru

7

8 In this study, logistic regression was used as a tool for deriving the binary classifier of tsunamigenic and  
9 nontsunamigenic earthquakes in the near-source zone. Earthquake source depth and moment magnitude were  
10 considered as predictors. The training dataset consisted of 767 M6.0+ submarine earthquakes, including 80  
11 tsunamigenic and 687 nontsunamigenic events that occurred in the northern part of the Pacific Ocean from 1960 to  
12 2020. The target area has already experienced significant and catastrophic tsunamis. The current analysis clearly  
13 showed that the data-driven logit model had a significantly lower false discovery rate relative to the threshold  
14 magnitude criteria that are widely used by tsunami warning agencies. At the same time, the balanced accuracy was  
15 about 71%, which suggests optimism for accurate tsunami forecasting based on the rapid interpretation of earthquake  
16 source parameters.

17

18 *Keywords:* submarine earthquake; tsunami; logit; predictive model; binary classifier; early warning

## 19 **1 Introduction**

20 Despite the fact that it is almost impossible to find the exact date of the first issued tsunami warning, the  
21 scientific principles for predicting this phenomenon were formulated in the 1920s (e.g., Finch 1924). Levin and Nosov  
22 (2016) report that probably the first successful tsunami forecast occurred on February 3, 1923, based on an  
23 interpretation of earthquake data.

24 S.L. Soloviev made a significant contribution to the development of the seismotectonic principle of tsunami  
25 forecasting. He proposed a magnitude–geographical criterion for the discrimination of tsunamigenic earthquakes  
26 (Soloviev and Shebalin 1959), which involved using ratios of the height (intensity) of the tsunami and the magnitude  
27 of the tsunamigenic earthquake. The analysis of the tsunamis produced by subduction earthquakes in the Far East  
28 region of Russia contributed to the assessment of the threshold intensity of the tsunami; intensities exceeding that  
29 threshold led to significant impacts in the target region. Based on the linear relationship between tsunami intensity  
30 and the magnitude of tsunamigenic earthquakes, a threshold magnitude corresponding to tsunami threshold intensity  
31 was selected.

32 Later, these empirical relations and threshold magnitudes were specified (e.g., Iida 1970; Soloviev 1972; Abe  
33 1979; Gusakov and Chubarov 1987; Boschetti and Ioualalen 2021). Current tsunami warning centers take into  
34 account not only the magnitude of a submarine earthquake but also the distance from the possible tsunami source to  
35 the coastal location (e.g., Amato et al. 2021). Such rules are represented by several hazard classes governing a certain  
36 action of emergency agencies.

37 In early studies (e.g., Iida 1970; Ivashchenko and Go 1973), the threshold magnitude as a function of the source  
38 depth was found by visual fitting of the line separating the two classes (tsunamigenic and nontsunamigenic) of  
39 submarine earthquakes on the magnitude–depth plot. These equations were never used in the tsunami early warning  
40 surveys. Of late, earthquake depth has been considered as a simplified indicator variable categorizing the seismic  
41 events into two groups: the first group includes events with a source depth of less than 100 km; the second, 100 km  
42 or more (e.g., Users Guide... 2017).

43 Nosov et al. (2018) developed a fully automatic service (Tsunami Observer) for tsunami forecasting.  
44 Earthquake source parameters such as the focal mechanism solution and source depth, reported by seismological  
45 agencies, were used to calculate the free surface deformation. The proposed method gave an estimation of the initial  
46 elevation of the water surface in the tsunami source region and the potential energy of the tsunami associated with it.  
47 Tsunami intensity was estimated based on the calibrated relation between the calculated tsunami energy and observed  
48 tsunami height. Then, the "threshold intensity" principle was applied to determine whether the submarine earthquake

49 would be tsunamigenic or not. This approach has a number of advantages over the current magnitude–geographical  
50 criteria. Its clarity lies in the use of additional parameters from numerical simulation. The disadvantage is the  
51 ambiguity of the nodal plane selection. The authors of the experimental service noted (Nosov et al. 2020) that tsunami  
52 energy estimates are highly sensitive to the choice of nodal plane for shallow earthquakes.

53 Cienfuegos et al. (2018) showed that the use of a simplified earthquake source model with uniform slip  
54 distribution led to conservative estimates of tsunami parameters. Epistemic uncertainties associated with spatial  
55 heterogeneity of earthquake slip play an important role in accurate tsunami prediction. Rational estimates of tsunami  
56 parameters in terms of probability can be obtained by stochastic simulation of the slip vector at the planar source and  
57 the position of the nodal plane in space, and subsequent numerical simulation of tsunami characteristics for each  
58 random realization of the finite fault model (Davies 2019). Probabilistic assessment of a tsunami hazard can also be  
59 realized using the Bayesian approach or the total probability theorem (Selva et al. 2021).

60 On the other hand, the focal mechanism is most commonly determined more slowly than the hypocenter  
61 parameters and the earthquake magnitude. For this reason, simulation-based approaches may lose their relevance in  
62 the near-source zone. Therefore, tsunami forecasting still remains an important practical issue of seismological  
63 observations based on rapid interpretation of earthquake source parameters.

64 In the current study, logistic regression (Cramer 2002) was used as a tool for the classification of tsunamigenic  
65 and nontsunamigenic earthquakes. We used a preprocessed dataset of submarine earthquakes that caused or did not  
66 cause tsunamis in the target region. Earthquake source depth and moment magnitude were considered as predictors.  
67 This approach was used here for the first time for developing the data-driven predictive model of tsunamigenic  
68 earthquakes.

## 69 **2 Materials and methods**

### 70 **2.1 Methods**

71 Logit analysis is widely used for binary data classification in geoscience (e.g., Zhu et al. 2017; Stepnova et al.  
72 2021). In the binary logistic model, the logistic response function is determined by a binary dependent variable  
73 (Harrell, 2015), where the two possible outcomes are labeled as 0 and 1. The logistic response function predicts not  
74 the results themselves but the probability of their onset:

$$75 \quad h_{\theta}(\mathbf{X}) = \frac{1}{1+e^{-\theta^T \mathbf{X}}} = P(Y = 1|\mathbf{X}, \theta), \quad (1)$$

76 where  $h_{\theta}(\mathbf{X})$  is the logistic response function,  $\mathbf{X} = \{1, X_1, X_2, \dots, X_n\}^T$  is the column vector of independent variables  
77 (predictors),  $\theta = \{\theta_0, \theta_1, \theta_2, \dots, \theta_n\}^T$  is the column vector of the parameters to be estimated,  $n$  is the number of

78 independent variables,  $Y$  is a binary dependent variable that has two values  $Y = 0$  and  $Y = 1$ , and  $P(Y = 1|\mathbf{X}, \boldsymbol{\theta})$  is  
 79 the conditional probability of event 1 with the given  $\mathbf{X}$  and  $\boldsymbol{\theta}$ , wherein

$$80 \quad P(Y = 0|\mathbf{X}, \boldsymbol{\theta}) = 1 - h_{\theta}(\mathbf{X}). \quad (2)$$

81 Each predictor can be set either by a continuous or by a categorical variable. Logistic regression measures the  
 82 relationship between a categorical dependent variable and one or more predictors by estimating the logarithm of the  
 83 odds:

$$84 \quad \text{Logit}(P) = \ln \frac{P(Y = 1|\mathbf{X}, \boldsymbol{\theta})}{1 - P(Y = 1|\mathbf{X}, \boldsymbol{\theta})} = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_n X_n, \quad (3)$$

85 where  $\text{Logit}(P)$  is the logit unit.

86 The given model follows the generalized Bernoulli scheme, in which a random variable takes only two values,  
 87 0 and 1, and their probabilities are determined by Eq. 2 and Eq. 1, respectively. Therefore, the probability distribution  
 88 function of the variable  $Y = y$  is given by

$$89 \quad P(y|\mathbf{X}, \boldsymbol{\theta}) = h_{\theta}(\mathbf{X})^y [1 - h_{\theta}(\mathbf{X})]^{1-y} \quad (4)$$

90 and the likelihood can be written as

$$91 \quad L(y|\mathbf{X}, \boldsymbol{\theta}) = \prod_i P(y_i|\mathbf{X}_i, \boldsymbol{\theta}) = \prod_i h_{\theta}(\mathbf{X}_i)^{y_i} [1 - h_{\theta}(\mathbf{X}_i)]^{1-y_i}, \quad (5)$$

92 where  $i$  runs from 1 to  $N$  ( $N$  is the sample size).

93 Maximization of the likelihood function is achieved by optimization methods. However, in practice, the log-  
 94 likelihood is most commonly used instead of the original likelihood function. Since the logarithm is a monotonically  
 95 increasing function, the maximum value of the logarithm of the probability occurs at the same point as the original  
 96 probability function. Taking the logarithm of Eq. 5 gives:

$$97 \quad L_{\log}(y|\mathbf{X}_i, \boldsymbol{\theta}) = N^{-1} \ln \prod_i P(y_i|\mathbf{X}_i, \boldsymbol{\theta}) = N^{-1} (\sum_i y_i \ln h_{\theta}(\mathbf{X}_i) + (1 - y_i) \ln [1 - h_{\theta}(\mathbf{X}_i)]). \quad (6)$$

98 If Eq. 6 is multiplied by  $-1$ , the given expression will explain logistic losses. Minimizing the logistic loss  
 99 function by the gradient descent method allows us to estimate the best parameters  $\boldsymbol{\theta}$ .

100 In our case, the moment magnitude  $M$  and the source depth  $D$  are predictors. The binary dependent variable is  
 101 the tsunamigenic class label of submarine earthquakes, which has two values:  $Y = \text{true}$  and  $Y = \text{false}$ . According to  
 102 Eq. 3, the logarithm of the odds for tsunamigenic class is a linear combination of independent variables  $M$  and  $D$ :

$$103 \quad \text{Logit}(P) = \ln \frac{P(Y = \text{true}|\mathbf{X}, \boldsymbol{\theta})}{1 - P(Y = \text{true}|\mathbf{X}, \boldsymbol{\theta})} = \theta_0 + \theta_1 * M + \theta_2 * D, \quad (7)$$

104 where  $\theta_0$ ,  $\theta_1$  and  $\theta_2$  are parameters of our model to be estimated.

105 The  $Logit(P)$  is used as a classifier of tsunamigenic/nontsunamigenic events. This is achieved by using a  
106 threshold probability  $P_{th}$ : if the right-hand side of Eq. 7 with the given earthquake parameters is equal to or greater  
107 than  $Logit(P_{th})$ , then the seismic event is categorized as tsunamigenic, with all other values belonging to the  
108 nontsunamigenic class.

109 For those classification problems that have a class imbalance, the threshold probability  $P_{th}$  may significantly  
110 differ from the default value 0.5. We used F-metrics as an indicator of the performance of a classifier that leads the  
111 tuning of the threshold level:

$$112 F_{\beta} = (1 + \beta^2) \frac{Precision \cdot Recall}{\beta^2 \cdot Precision + Recall} \quad (8)$$

113 *Precision* and *Recall* are defined as follows (see Eq. 9 and Eq. 10, respectively):

$$114 Precision = \frac{TP}{TP+FP}, \quad (9)$$

$$115 Recall = \frac{TP}{TP+FN}, \quad (10)$$

116 where  $TP$  (true positives) and  $TN$  (true negatives) are correctly predicted classes, and  $FP$  (false positives) and  $FN$   
117 (false negatives) are type I and type II errors, respectively.

118 Parameter  $\beta$  is selected so that *Recall* is considered  $\beta$  times more essential than *Precision*. When the number  
119 of true tsunamigenic events that may be incorrectly categorized as nontsunamigenic (type I error) is minimized, the  $\beta$   
120 parameter can be different from 1. However, in the current study, it is taken as equal to 1. This means that *Precision*  
121 and *Recall* have the same significance or, in the other words, that the F1 score is the mean harmonic value between  
122 *Precision* and *Recall*.

123 Generally, metrics in Eq. 8–10 are functions of  $P$ . Maximizing the F1 score with respect to the  $P$  allows us to  
124 select the optimal threshold probability  $P_{th}$  used in the binary classification.

## 125 2.2 Study area

126 The study area is limited by four angular points that form a polygon with the geographical coordinates 33–  
127 55°N and 138–166°E (see Fig. 1). This part of the Pacific Ocean encompasses possible tsunami sources that may  
128 cause significant impacts to the coastal locations of island arcs considered as near-source zones of tsunamis. The target  
129 area has already experienced significant and catastrophic tsunamis (e.g., Tanioka and Satake 1996; Mimura et al.

130 2011; Gusakov 2011): the 1737 Kamchatka Tsunami, the 1886 Great Meiji Sanriku Tsunami, the 1952 North Kuril  
131 Tsunami and the 2011 Great East Japan Tsunami.

## 132 **2.3 Data**

133 The U.S. Geological Survey's earthquake database (U.S. Geological Survey 2021) was used in the current  
134 study. It is available at the website: <https://earthquake.usgs.gov/earthquakes/search/>. The time window from 1960 to  
135 2020 was considered as a period of intensive increase in the number of instrumental sites. Initially, we used the M 4.5+  
136 dataset. After statistical tests of the completeness level of the seismic database, the cutoff magnitude was increased to  
137 Mc 6 (see Sect. 2.4.1). Selection parameters for the earthquake depth were not limited.

138 As the seismological data centers report different types of magnitude such as Mw, Ms and mb, all magnitude  
139 estimates were unified to the Mw scale based on global intermagnitude relations (Gusev 1991). The magnitude  
140 conversion from M to Mw was achieved by polynomial approximation of the cubic type, where M is either Ms or mb:

$$141 \quad M_w = a * M^3 + b * M^2 + c * M + d, \quad (11)$$

142 where  $a$ ,  $b$ ,  $c$  and  $d$  are the regression coefficients (see Table 1).

143 The 2021 tsunami database provided by the US National Oceanic and Atmospheric Administration is free and  
144 available at the National Geophysical Data Center's website:  
145 <https://www.ngdc.noaa.gov/hazel/view/hazards/tsunami/>. The search criteria used for the tsunami dataset were the  
146 same as those used for the seismic database selection. The analysis of the completeness level of the tsunami database  
147 (see Sect. 2.4.2) indicated that events with a maximum water height of less than 10 cm could be missed due to the  
148 sensitivity of the ocean-bottom network. Therefore, tsunamis with reported maximum water heights of less than 10  
149 cm or no reported height values were removed from the dataset.

## 150 **2.4 Completeness check**

### 151 **2.4.1 Completeness level of seismic database**

152 A statistical model developed by Ogata and Katsura (1993) was applied in the current study for estimating the  
153 magnitude of completeness of the earthquake catalogue. The  $\mu$  (50% detection rate),  $b$ -value and  $\sigma$  can be  
154 simultaneously estimated by numerically maximizing the likelihood function:

$$155 \quad L(M, \mathbf{p}) = \frac{W(M, \mathbf{p})q(M, \mathbf{p})}{\int_{-\infty}^{+\infty} W(M, \mathbf{p})q(M, \mathbf{p})dM}, \quad (12)$$

156 where  $\mathbf{p}$  is a parameter vector  $(\mu, b, \sigma)$  to be estimated,  $W(M) = e^{-\beta M}$  with  $\beta = b \ln(10)$  is the distribution based on  
157 the Gutenberg–Richter relation, and  $q(M, \mathbf{p}) = \int_{-\infty}^M \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} dx$  is the detection rate function.

158 Given the observed magnitudes  $M_1, \dots, M_N$ , model parameters were estimated by minimizing the  
159  $\ln L(M, \mathbf{p}) = \sum_{i=1}^N \ln L(M_i, \mathbf{p})$  using the trust region method (Conn et al. 2000). The code is available on GitHub:  
160 <https://gist.github.com/jamm1985/79b43a6dc1655a32e44ae64ee437d5db> (accessed 10 February 2022).

161 The magnitude of completeness is determined by the equation

$$162 \quad m_c(k) = \mu + k \sigma, \quad (13)$$

163 where  $k$  is the confidence level;  $k = 0, 1, 2, 3$  means that 50%, 68%, 95% and 99% of events are detected above  $m_c$ .

164 According to Fig. 2, the model parameters with the given seismic dataset were well constrained. The magnitude  
165 of completeness (95% confidence level) determined from Eq. 13 was equal to  $m_c$  6.

## 166 2.4.2 Completeness level of tsunami database

167 The cumulative counts depending on the log of the maximum water height (as a magnitude unit) were plotted  
168 for the analysis of the completeness level of the tsunami database (Fig. 3). The cumulative frequency–magnitude  
169 relation had a log-linear form. At the smallest magnitudes, the log-linear dependence was not fulfilled. It is usually  
170 assumed that the dataset is complete up to the threshold magnitude for which the log-linear dependence on the  
171 cumulative frequency–magnitude plot is fulfilled. According to the recurrence plot in Fig. 3, events with heights of  
172 less than 10 cm could occur more frequently, but were missed due to the sensitivity of the equipment. Various studies  
173 (e.g., Baba et al. 2004) indicate that the minimum tsunami height recorded by ocean-bottom equipment is about 5 cm.  
174 Therefore, the threshold value of the maximum water height reflects the detection bounds of the instrumentation  
175 network aimed at tsunami registration in the study area. We considered this value (10 cm) as the completeness level  
176 of the tsunami database.

## 177 2.5 Data preprocessing

178 We filtered the earthquake catalogue in order to remove nonsubmarine seismic events. In addition, at this stage,  
179 catalogs of submarine earthquakes and tsunamis were merged into one seismic dataset, in which an additional binary  
180 variable associated with the tsunamigenic class (true/false) was assigned to each seismic event. The following rules  
181 of association of events from the two datasets were adopted: date matching, matching of the epicenter locations with  
182 valid differences of up to 0.5 degrees, and matching of earthquake magnitudes with valid differences of up to 0.5. If  
183 events from the two subsets were associated, then the tsunamigenic variable was labeled as true; otherwise, it was

184 mapped as false. This means that the nonassociated events from the earthquake catalogue were automatically assigned  
185 a nontsunamigenic class.

186 During the merging of the seismic and tsunami databases into one training dataset, seven nonassociated events  
187 from the tsunami subset were identified. The earthquake source parameters of the three nonassociated tsunamis were  
188 outside the valid differences defined in the rules of association. The epicenters of the remaining four nonassociated  
189 tsunamigenic events according to the earthquake catalogue were located on land. These events were removed from  
190 the earthquake catalogue at the initial stage of data preprocessing. A detailed analysis of earthquake locations indicated  
191 that one of them, the catastrophic 1995 Neftegorsk earthquake (Mw 7.1), occurred in Sakhalin Island and could not  
192 have caused a tsunami in the ocean, while the other three events occurred in the immediate vicinity of the coastline  
193 and were accompanied by tsunami waves. This can be explained by effects related to errors in the epicenter location  
194 or the size of the earthquake source, when an area of the finite fault is located on land, while another area is underwater.  
195 We applied the corresponding adjustments to the training dataset. Maps of earthquake epicenters and tsunami sources,  
196 which were both removed and used in the current study after preprocessing, are shown in Fig. 1a and Fig. 1b,  
197 respectively.

198 Finally, after careful verification of the training dataset, we found that two events with  $M_w \geq 7.7$  and source  
199 depths of less than 50 km were categorized as nontsunamigenic. Independent tsunami data centers (e.g., Web  
200 Encyclopedia on Natural Hazards 2021) reported the same information. The considered events occurred on 11.03.2011  
201 in the aftershock area of the Great Japanese Earthquake (Mw 9.1), which occurred a few minutes earlier and caused  
202 catastrophic tsunami waves. Wave processes induced by the mainshock and by the largest aftershocks could be  
203 superimposed. Therefore, these aftershocks were excluded from the analysis since they could not be objectively  
204 attributed to one of the two classes.

205 The final database (see SI) consisted of 767 submarine earthquakes, including 80 tsunamigenic and 687  
206 nontsunamigenic events (Fig. 4).

### 207 **3 Results and discussion**

208 Fig. 5 shows the distribution of submarine earthquakes on the 2D magnitude–depth plot. The class category is  
209 indicated by a color. It is clear that tsunamigenic earthquakes tend to be above a certain imaginary line, the parameters  
210 of which should be estimated.

211 Applying the logit tool to the joint earthquake–tsunami database, it was found that the conditional probability  
212 that a submarine earthquake with a given moment magnitude and source depth would cause a tsunami is given by the  
213 following expression:

214  $h_{\theta}(\mathbf{X}) = \text{Sig}(0.993 \cdot M - 0.0349 \cdot D - 7.3224),$  (14)

215 where Sig is the sigmoid function determined in Eq. 1.

216 The probability in Eq. 14 increases with magnitude and decreases with depth.

217 As we are dealing with class–imbalanced data (there were many more records of nontsunamigenic earthquakes  
218 than tsunamigenic cases), it was necessary to select an optimal threshold probability according to Eq. 8. Fig. 6 shows  
219 that the F1 score reached the highest value at a probability of approximately 0.22, and in this way, the balance between  
220 *Precision* and *Recall* was achieved.

221 Substituting the obtained threshold probability (0.22) in Eq. 14, we obtain the threshold magnitude given by  
222 the functional dependence on the source depth of the submarine earthquake:

223  $M_{th} = 5.8 + 0.035 \cdot D.$  (15)

224 The predictive model of tsunamigenic earthquakes obtained in the current study (see Eq. 15), as well as  
225 proposed by K. Iida (1963), are shown in Fig. 5. Iida was the first to show that the threshold magnitude increases with  
226 the source depth according to a linear law. Eq. 15 reflects a well-known fact; however, the model coefficients were  
227 statistically justified for the first time and determined from a dataset over a long duration of instrumental observation.

228 It is well known that the source depth of shallow earthquakes is characterized by uncertainty, especially for  
229 distant events. In seismic bulletins, this parameter is most commonly fixed. Generally, the interplate type of seismicity  
230 predominates in the Japan–Kurile–Kamchatka subduction zone. The source depth of shallow interplate earthquakes  
231 is fixed in the range of 30–40 km. From Eq. 15, it follows that the threshold magnitude for the depth of 33 km  
232 corresponds to Mw 7.0. This estimate corresponds to the minimum threshold magnitude used in the tsunami warning  
233 survey of Russia and of many other centers.

234 The hypocenters of the crustal earthquakes are most commonly located at a depth of 10 km. In this case, the  
235 logit model predicted the threshold magnitude of Mw 6.2 (see Eq. 15). Some shallow earthquakes with  $M \sim 6$ ,  
236 occurring in the coastal zone of the northwest part of the Pacific seismic belt, are accompanied by significant tsunami  
237 waves (e.g., Satake and Kanamori 1991; Konovalov et al. 2015). The results obtained in the current study correspond  
238 to the considerations expressed recently on clarifying the threshold magnitude of tsunamigenic earthquakes (e.g.,  
239 Gusiakov 2011).

240 The performance of a binary classifier was tested using balanced accuracy (*BA*) and the false discovery rate  
241 (*FDR*). The false discovery rate indicates the counts of tsunamigenic earthquakes that may be incorrectly categorized

242 as nontsunamigenic. We did not split the training and validation subsets, due to the low number of positively  
243 categorized classes. The results from testing the model should be considered as the results obtained with the training  
244 dataset. The earthquake source depth of the validation subset was in the range of 0–250 km.

245 The given metrics are defined as follows:

$$246 \quad BA = 0.5 \cdot \left( \frac{TP}{TP+FN} + \frac{TN}{TN+FP} \right), \quad (16)$$

$$247 \quad FDR = \frac{FP}{TP+FP}. \quad (17)$$

248 Table 2 summarizes the performance tests considering different models: the data-driven logit model (see  
249 Eq. 15), the Iida-1963 model (see Eq. 18) and the threshold magnitude criterion that is currently used by the Russian  
250 tsunami warning survey (see Eq. 18). The Iida-1963 model is applicable to tsunamis with an intensity greater than 1  
251 (Iida 1963):

$$252 \quad M_{th} = 6.46 + 0.019 \cdot D. \quad (18)$$

253 The threshold magnitude criterion is given by a simplified expression:

$$254 \quad M_{th}=7. \quad (19)$$

255 As can be seen from Table 2, the Iida-1963 model had a higher balanced accuracy relative to the logit model.  
256 However, the Iida-1963 model had a relatively high false discovery rate; that is, the Iida-1963 model was mainly  
257 applicable to significant tsunamis. The threshold magnitude criterion is widely used by tsunami warning agencies  
258 around the world. However, in this case, the counts of missed warnings remained significant. The current analysis  
259 clearly showed that the data-driven logit model had a significantly lower false discovery rate at the 70% accuracy  
260 level.

261 For future studies, it would be advisable to consider additional categorical variables, such as the type of focal  
262 mechanism (strike-slip/reverse/normal/unknown), the tectonic type of earthquake based on the geographical location  
263 of the hypocenter (subduction/collision/rift/unknown) as well as continuous variables, for example, underwater depth  
264 in the earthquake epicenter.

265 Chicco and Jurman (2020) showed that the Matthews correlation coefficient produces more informative and  
266 truthful metrics for deriving the binary classifier, especially for imbalanced datasets. They suggest that it should be  
267 used instead of the F1 score. We suggest that further evaluation with the Matthews correlation coefficient using  
268 additional independent variables would help to significantly improve the performance of the binary classifier.

269 **4 Conclusions**

270 The logistic regression tool was used for deriving the data-driven binary classifier of tsunamigenic and  
271 nontsunamigenic submarine earthquakes. It was assumed that source depth and moment magnitude are the most  
272 important factors controlling tsunami generation in the near-source zone. The logit model coefficients are statistically  
273 justified and determined using preprocessed tsunami and earthquake catalogues that contain empirical data over a  
274 long duration of instrumental observation in the Northern Pacific Ocean. The advantages of the considered approach  
275 are not only the use of the data on the intensity of tsunamigenic earthquakes but also data of submarine earthquakes  
276 that did not cause tsunamis.

277 The threshold magnitude criteria are widely used by tsunami warning agencies. The current study clearly  
278 showed that the data-driven logit model had a significantly lower false discovery rate. At the same time, the balanced  
279 accuracy was about 71%, which suggests optimism for accurate tsunami forecasting based on the rapid interpretation  
280 of earthquake source parameters. The authors hope that the given model improves tsunami forecasting.

281 **Declarations**

282 **Conflict of interest** The authors declare that they have no conflict of interest.

283 **Data availability** The primary data used in the current study are freely available at the indicated sources. Data  
284 that result through preprocessing are available as electronic supplementary material.

285 **Code availability** Visual Studio Code Source Editor and Python programming language were used in the  
286 current study. The following Python libraries were employed: NumPy, Matplotlib, Scikit-learn and Global-land-mask.  
287 Code is available upon request.

288 **References**

- 289 Gusev AA (1991) Intermagnitude relationships and asperity statistics. *Pure Appl Geophys* 136:515–527.
- 290 Gusiakov VK (2011) Magnitude-geographical criterion for operational tsunami prognosis: Analysis of application in  
291 1958–2009. *Seism Instr* 47:203. <https://doi.org/10.3103/S0747923911030078>
- 292 Gusakov VK, Chubarov LB (1987) Numerical simulation of tsunami excitation and propagation in the coastal zone.  
293 *Izv AN SSSR. Physics of the Earth* 21(11):53–64.
- 294 Ivashchenko AI, Go CHN (1973) Tsunamigennost' i glubina ochaga zemletryaseniya // V kn.: Volny tsunami.  
295 Yuzhno-Sakhalinsk: SakhKNII DVNTS AN SSSR 32:152–155 (in Russian).
- 296 Konovalov AV, Nagornyykh TV, Safonov DA, Lomtev VL (2015) Nevelsk earthquakes of August 2, 2007 and seismic  
297 setting in the southeastern margin of Sakhalin Island. *Russ J of Pac Geol* 9:451–466.  
298 <https://doi.org/10.1134/S1819714015060056>
- 299 Nosov MA, Kolesov SV, Bolshakova AV, Nurislamova GN (2020) The Effect of the Choice of the Nodal Plane on  
300 Tsunami Energy Estimates. *Moscow Univ Phys* 75:501–506. <https://doi.org/10.3103/S0027134920050197>
- 301 Nosov MA, Kolesov SV, Bolshakova AV, Nurislamova GN, Sementsov KA, Karpov VA (2018) Automated system  
302 for estimation of tsunami hazard of an earthquake. *Uch Zap Fiz Fak Mosk Univ* 5:1850901 (in Russian).  
303 <http://uzmu.phys.msu.ru/abstract/2018/5/1850901>
- 304 Solov'ev SL (1972) Tsunami and Earthquake Recurrence in the Pacific Ocean. *Tr SakhKNII* 29:7–47 (in Russian).
- 305 Solov'ev SL, Shebalin NV (1959) Tsunami and Intensity of Kuril-Kamchatka Earthquakes. *Izv AN SSSR. Ser Geofiz*  
306 8:1195–1198.
- 307 Abe K (1979) Size of great earthquakes of 1873-1974 inferred from tsunami data. *J Geophys Res* 84:1561–1568.
- 308 Amato A, Avallone A, Basili R et al. (2021) From seismic monitoring to tsunami warning in the Mediterranean Sea.  
309 *Seismol Res Lett* 92:1796–1816. <https://doi.org/10.1785/0220200437>
- 310 Baba T, Hirata K, Kaneda Y (2004) Tsunami magnitude determined from ocean-bottom pressure gauge data around  
311 Japan. *Geophysical Research Letters* 31:L08303. doi: 10.1029/2003GL019397
- 312 Boschetti L, Ioualalen M (2021) Integrated tsunami intensity scale based on maxima of tsunami amplitude and induced  
313 current. *Nat Hazards* 105:815–839. <https://doi.org/10.1007/s11069-020-04338-5>

314 Chicco D, Jurman G (2020) The advantages of the Matthews correlation coefficient (MCC) over F1 score and  
315 accuracy in binary classification evaluation. *BMC Genomics* 21:6. [https://doi.org/10.1186/s12864-019-6413-](https://doi.org/10.1186/s12864-019-6413-7)  
316 7

317 Cienfuegos R, Catalán PA, Urrutia A, Benavente R, Aránguiz R, González G (2018) What can we do to forecast  
318 tsunami hazards in the near field given large epistemic uncertainty in rapid seismic source inversions.  
319 *Geophysical Research Letters* 45:1–12. <https://doi.org/10.1029/2018GL076998>

320 Conn AR, Gould NI, Toint, PL (2000) *Trust region methods*, Siam, pp 169–200

321 Cramer JS (2002) *The Origins of Logistic Regression*. Tinbergen Institute Working Paper. No 2002-119/4. P. 167–  
322 168.

323 Davies G (2019) Tsunami variability from uncalibrated stochastic earthquake models: tests against deep ocean  
324 observations 2006–2016. *Geophysical Journal International* 218(3):1939–1960. doi: 10.1093/gji/ggz260

325 Iida K (1963) Magnitude of tsunamigenic earthquake, aftershock area and area of tsunami origin. *Geophysical papers*  
326 dedicated to prof. Kenzo Sassa.

327 Iida K (1970) The generation of tsunamis and the focal mechanism of earthquakes. In: Adams, W.M. (ed.). *Tsunamis*  
328 in the Pacific Ocean. Honolulu (Hawaii), East-West Center Press, 3–18.

329 Finch RH (1924) On the Prediction of Tidal Waves. *Monthly Weather Review* 52(3):147–148.

330 Harrell FE (2015) *Binary Logistic Regression*. In: *Regression Modeling Strategies*. Springer Series in Statistics.  
331 Springer, Cham. [https://doi.org/10.1007/978-3-319-19425-7\\_10](https://doi.org/10.1007/978-3-319-19425-7_10)

332 Levin BV, Nosov MA (2016) *Physics of Tsunamis* (Springer-Verlag, New York).

333 Mimura N, Yasuhara K, Kawagoe S, Yokoki H, So K (2011) Damage from the Great East Japan Earthquake and  
334 Tsunami – A quick report. *Mitig Adapt Strateg Glob Change* 16:803–818. [https://doi.org/10.1007/s11027-011-](https://doi.org/10.1007/s11027-011-9297-7)  
335 9297-7

336 National Geophysical Data Center / World Data Service: NCEI/WDS Global Historical Tsunami Database. NOAA  
337 National Centers for Environmental Information. doi:10.7289/V5PN93H7 (accessed 30 June 2021).

338 Ogata Y, Katsura K (1993) Analysis of temporal and spatial heterogeneity of magnitude frequency distribution  
339 inferred from earthquake catalogues. *Geophysical Journal International* 113(3):727–738.  
340 <https://doi.org/10.1111/j.1365-246x.1993.tb04663.x>

- 341 Satake K, Kanamori H (1991) Abnormal tsunamis caused by the June 13, 1984, Torishima, Japan, earthquake. J  
342 Geophys Res 96(B12):19933–19939. doi: 10.1029/91JB01903
- 343 Selva J, Lorito S, Volpe M et al. (2021) Probabilistic tsunami forecasting for early warning. Nat Commun 12:5677.  
344 <https://doi.org/10.1038/s41467-021-25815-w>
- 345 Stepnova YA, Stepnov AA, Konovalov AV et al. (2021) Predictive Model of Rainfall-Induced Landslides in High-  
346 Density Urban Areas of the South Primorsky Region (Russia). Pure Appl Geophys.  
347 <https://doi.org/10.1007/s00024-021-02822-y>
- 348 Tanioka Y, Satake K (1996) Fault parameters of the 1896 Sanriku tsunami earthquake estimated from tsunami  
349 numerical modeling. Geophys Res Lett 23:1549–1552. doi: 10.1029/96GL01479
- 350 U.S. Geological Survey, 2021, Search Earthquake Catalog: <https://earthquake.usgs.gov/earthquakes/search/> (accessed  
351 30 June 2021).
- 352 Users Guide for the Pacific Tsunami Warning Center Enhanced Products for the Tsunami and other Coastal Hazards  
353 Warning System for the Caribbean and Adjacent Regions (CARIBE-EWS). IOC Technical Series, No 135.  
354 UNESCO/IOC, 2017.
- 355 Web Encyclopedia on Natural Hazards, Institute of Computational Mathematics and Mathematical Geophysics SB  
356 RAS, Tsunami Laboratory, Novosibirsk, Russia <http://tsun.sbcc.ru/nh/list.html> (accessed 30 June 2021).
- 357 Zhu J, Baise LG, Thompson EM (2017) An Updated Geospatial Liquefaction Model for Global Application. Bulletin  
358 of the Seismological Society of America 107(3):1365–1385. doi: 10.1785/0120160198

359 **Figure captions**

360 **Fig. 1** Maps of study area in the Northern Pacific Ocean: **a** location of earthquake epicenters,  $M \geq 6$ , 1960–2020; **b**  
361 location of tsunami sources, 1960–2020.

362 **Fig. 2** The observed ('+' symbol) and calculated (black line) probability density function of the magnitude for the  
363 earthquake catalogue,  $M \geq 4.6$ , 1960–2020. The best estimates of the model parameters are shown at the top of the  
364 panel.

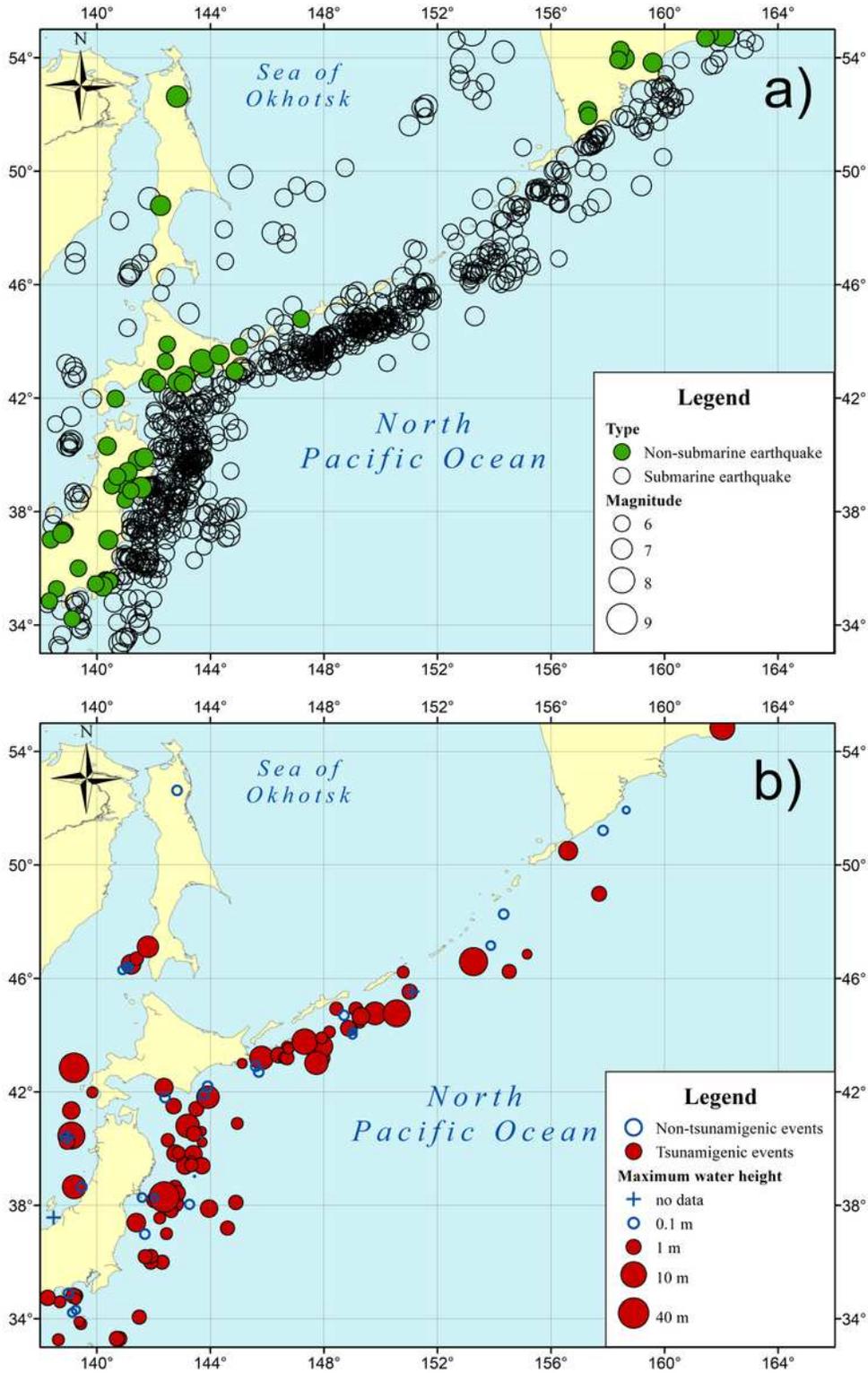
365 **Fig. 3** Cumulative frequency–magnitude distribution for the catalogue of tsunamis, 1960–2020. Tsunami magnitude  
366 is measured as  $m_t = \log_2 h_{max}$ , where  $h_{max}$  is the maximum water height.

367 **Fig. 4** Location of submarine earthquakes used in the logit model. The earthquake catalogue was merged with tsunami  
368 database and filtered.

369 **Fig. 5** Distribution of binary variable that categorizes the tsunamigenic class of submarine earthquake on the 2D  
370 magnitude–depth plot according to the training dataset (see Sect. 2.5). Predictive models of tsunamigenic earthquakes:  
371 data-driven logit model (black solid line) and Iida-1963 (black dotted line).

372 **Fig. 6** F1 score against the threshold probability of the logit model.

# Figures



**Figure 1**

Maps of study area in the Northern Pacific Ocean: **a** location of earthquake epicenters,  $M \geq 6$ , 1960–2020; **b** location of tsunami sources, 1960–2020.

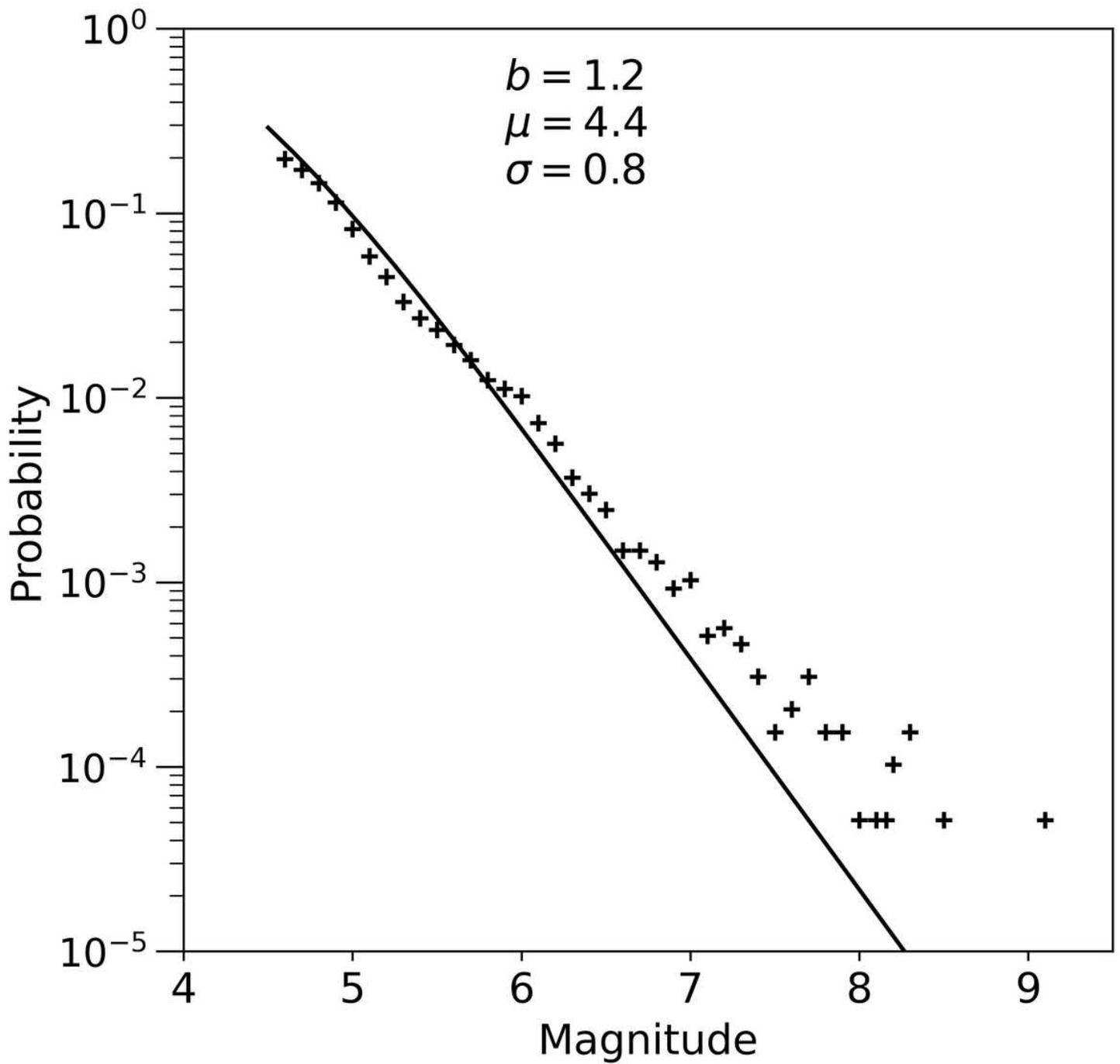
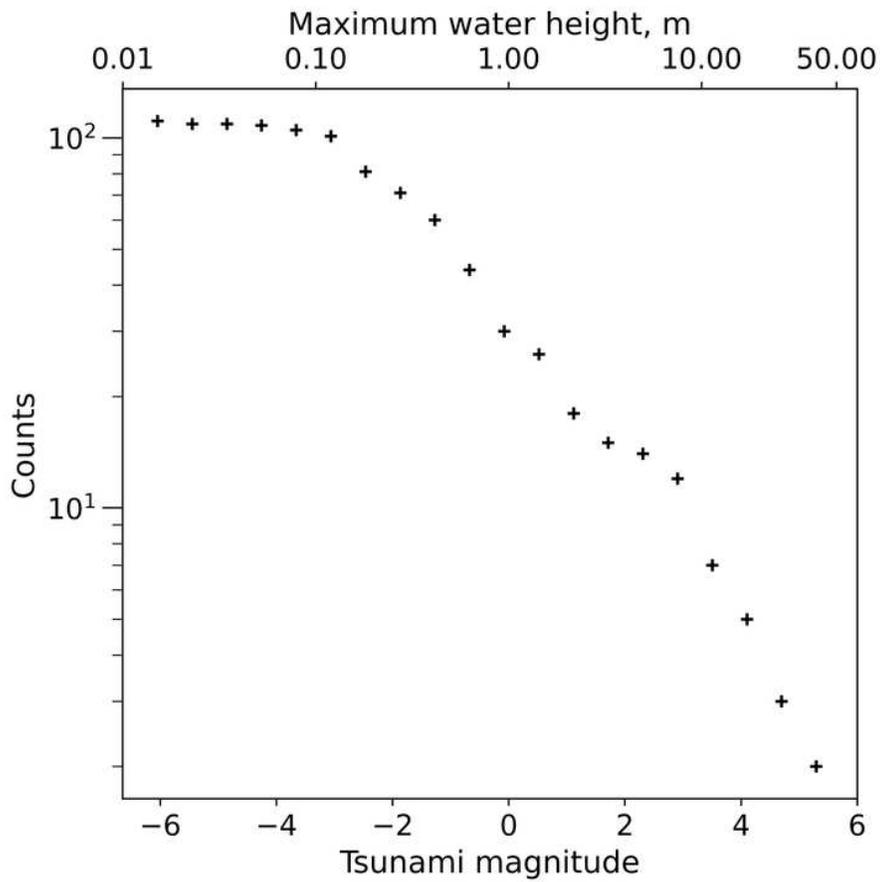


Figure 2

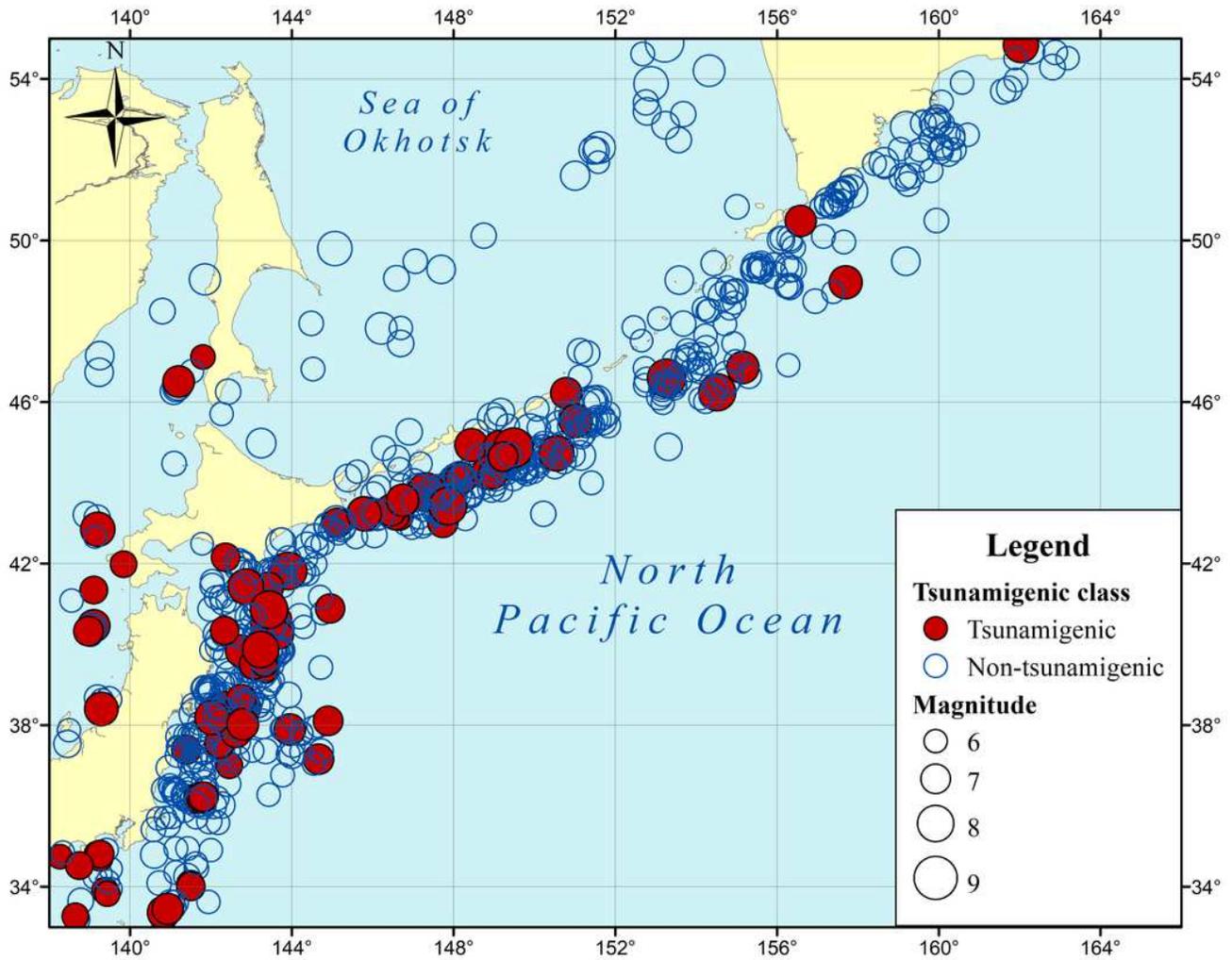
The observed ('+' symbol) and calculated (black line) probability density function of the magnitude for the earthquake catalogue,  $M \geq 4.6$ , 1960–2020. The best estimates of the model parameters are shown at the top of the panel.



**Fig. 3** Cumulative frequency–magnitude distribution for the catalogue of tsunamis, 1960–2020. Tsunami magnitude is measured as  $m_t = \log_2 h_{max}$ , where  $h_{max}$  is the maximum water height.

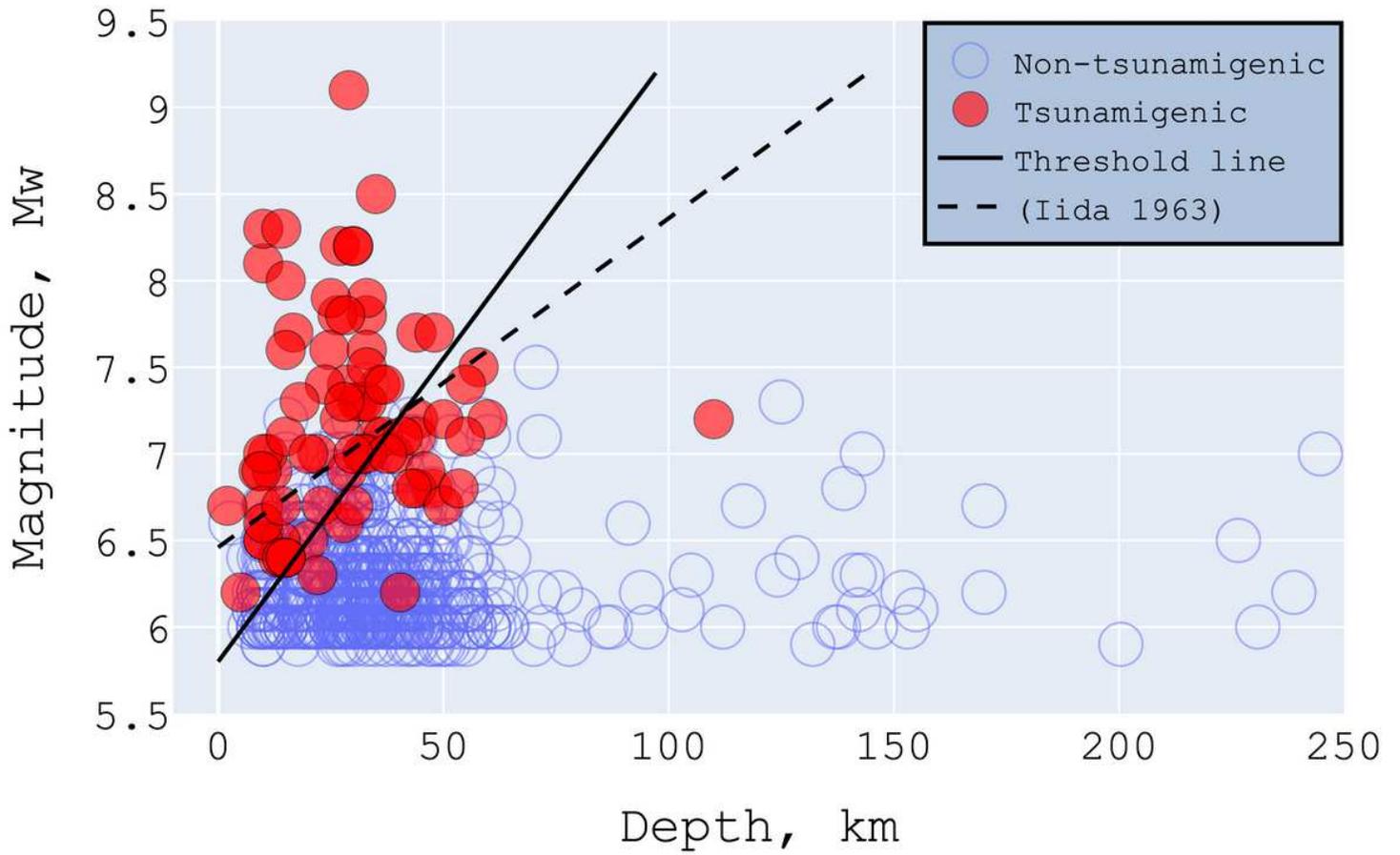
**Figure 3**

See figure for legend.



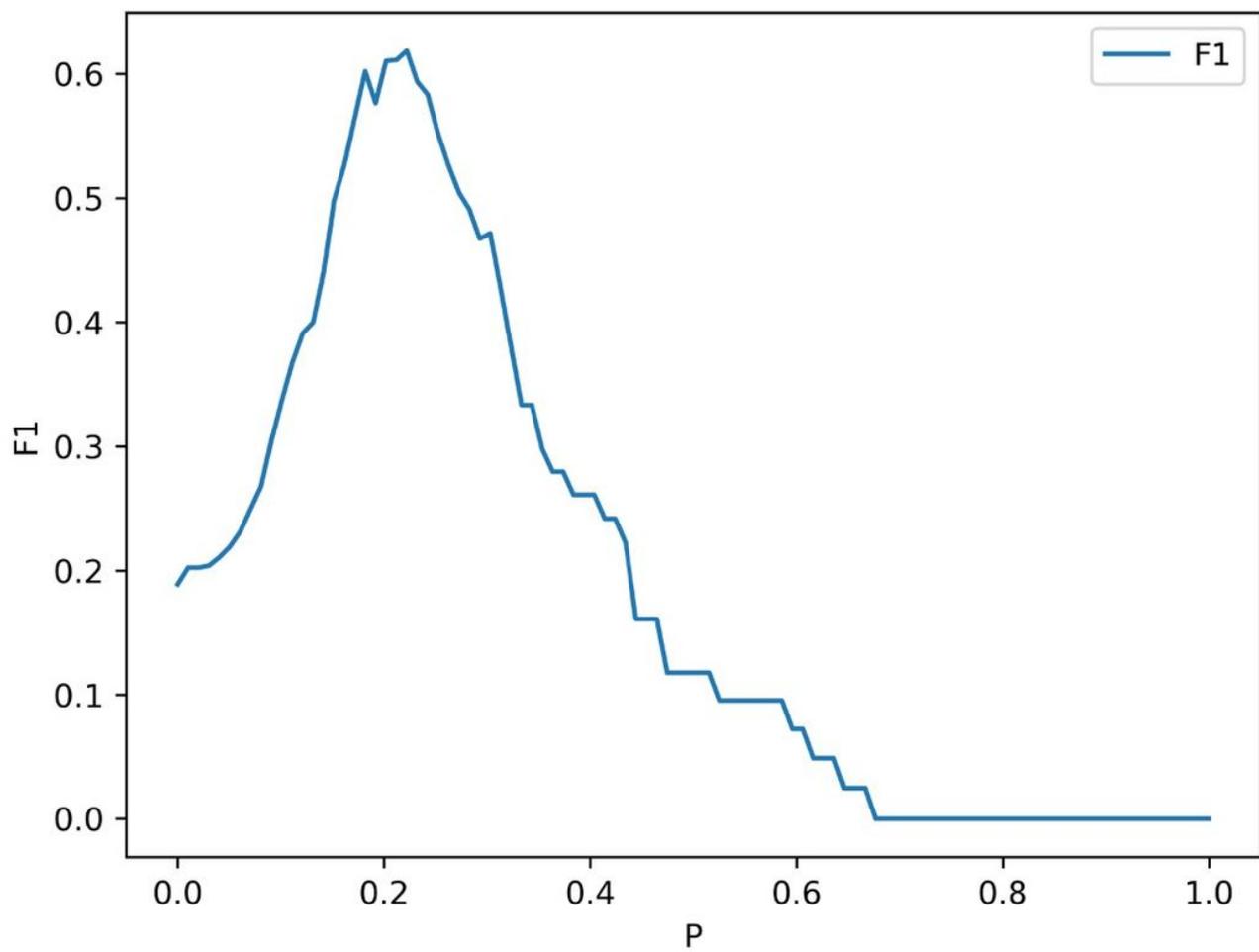
**Figure 4**

Location of submarine earthquakes used in the logit model. The earthquake catalogue was merged with tsunami database and filtered.



**Figure 5**

Distribution of binary variable that categorizes the tsunamigenic class of submarine earthquake on the 2D magnitude–depth plot according to the training dataset (see Sect. 2.5). Predictive models of tsunamigenic earthquakes: data-driven logit model (black solid line) and Iida-1963 (black dotted line).



**Figure 6**

F1 score against the threshold probability of the logit model.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Finalearthquakelistwithtsunamiindicator.csv](#)