

Statistica Sinica Preprint No: SS-2017-0403

Title	Detection and replenishment of missing data in marked point processes
Manuscript ID	SS-2017-0403
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202017.0403
Complete List of Authors	Jiancang Zhuang Ting Wang and Koji Kiyosugi
Corresponding Author	Jiancang Zhuang
E-mail	zhuangjc@ism.ac.jp

1 Detection and replenishment of missing data 2 in marked point processes

3 *Jiancang Zhuang¹⁾, Ting Wang²⁾, and Koji Kiyosugi³⁾*

¹⁾Institute of Statistical Mathematics

10-3 Midori-cho, Tachikawa, Tokyo 190-8562, Japan

²⁾Department of Mathematics and Statistics, University of Otago, Dunedin, New Zealand

³⁾Organization of Advanced Science and Technology, Kobe University

1-1 Rokkodai-cho, Nada-ku, Kobe 657-8501, Japan

4 August 13, 2019

5 Abstract

6 Records of geophysical events, such as earthquakes and volcanic eruptions,
7 are usually modeled as marked point processes. These records often suffer
8 from missing data, resulting in underestimations of the corresponding haz-
9 ards. We propose a computational approach for replenishing data missing
10 from the records of temporal point processes with time-separable marks. The
11 proposed method is based on the notion that if such a point process is com-
12 pletely observed, it can be transformed into a homogeneous Poisson process,
13 approximately on the unit square $[0, 1]^2$, by a biscale empirical transformation
14 (BEPIT). This approach includes three key steps: (1) transforming the pro-
15 cess onto $[0, 1]^2$ using the BEPIT, and finding a time–mark range that likely
16 contains missing events; (2) estimating a new empirical distribution function
17 based on the data in the time–mark range in which the events are supposed
18 to be completely observed; and (3) generating events in the missing region.
19 We test this method on a synthetic data set, and apply it to records of the

20 volcanic eruptions of the Hakone Volcano in Japan and the aftershock se-
21 quence following the 2008 Wenchuan Mw7.9 earthquake in Southwest China.
22 The results show that this algorithm provides a useful way to estimate miss-
23 ing data and to replenish incomplete records of marked point processes. In
24 addition, the replenished data provide estimates of the hazard function that
25 are more robust.

26 **1 Introduction**

27 Many geophysical processes, such as earthquakes and volcanic eruptions, oc-
28 cur at random times and/or locations, and, thus, are described naturally by
29 point-process models (e.g., Vere-Jones, 1970; Zhuang et al., 2002; Wang and
30 Bebbington, 2012, 2013). Point-process models and their related theories are
31 also widely used in fields such as crime, disease, and fire (Diggle and Rowl-
32 ingson, 1994; Schoenberg et al., 2007; Mohler et al., 2011). Furthermore,
33 advancements in the technology used to record these natural and social phe-
34 nomena are yielding significantly greater amounts of data. However, the de-
35 gree of completeness of these records varies, and in many cases, small events
36 are often missed in the early period of observation. For example, smaller
37 aftershocks are less likely to be recorded than are larger aftershocks during
38 the period immediately following a large earthquake (Ogata and Katsura,
39 1993; Omi et al., 2013). Other examples include missing data in volcanic
40 eruption records (Kiyosugi et al., 2015) and in the field of communication in
41 social networks (Zipkin et al., 2015). Missing data limit our efficient use of
42 these records, often resulting in biased estimates. However, statistical tools
43 for analyzing incomplete point-process data are not well developed.

44 Geophysicists have been searching for reliable methods of obtaining earth-
45 quake catalogs that are more complete. For example, waveform-based detec-
46 tion methods for small earthquakes within an aftershock sequence have been
47 proposed (e.g., Enescu et al., 2007, 2009; Peng et al., 2007; Marsan and
48 Enescu, 2012; Hainzl, 2016). However, even these methods cannot recover
49 all missing aftershocks. An alternative is to switch to energy-based descrip-
50 tions (Sawazaki and Enescu, 2014); that is, rather than viewing earthquake
51 occurrences as a process of events with different magnitudes, the process is
52 regarded as a stream of energies released by earthquakes. However, methods
53 related to such descriptions remain underdeveloped.

54 Based on the empirical law that the distribution of earthquake magni-
55 tudes follows the Gutenberg–Richter magnitude–frequency relation (Guten-
56 berg and Richter, 1944), Ogata and others investigated why events were
57 missing from earthquake catalogs (Ogata and Vere-Jones, 2003; Iwata, 2008,
58 2013, 2014). They used a Bayesian method to make probabilistic earthquake
59 forecasts, with missing earthquakes taken into account (Ogata, 2006; Omi
60 et al., 2013, 2014, 2015).

61 In most of the aforementioned studies, when dealing with missing events
62 in a point process, the full structure of the model, or at least the distribution
63 of marks, is assumed to be known. However, owing to incomplete records and
64 other reasons, on most occasions, the information available on the process
65 or the mark distribution is limited. Thus, a preferable method for evalu-
66 ating the missingness should be based on as few assumptions as possible,
67 especially when the temporal structure and the distribution of marks are
68 unknown. Zhuang et al. (2017) used a stochastic algorithm to restore miss-
69 ing aftershocks in the aftershock sequences following several earthquakes in

70 Kumamoto, Japan (April 14, 2016, $M6.5$; April 15, 2016, $M6.4$; April 16,
71 2016, $M7.3$). This method can be used to restore missing data in the records
72 of a more general temporal point process with time-separable marks, using
73 information from the parts of the process that are completely observed. In
74 Zhuang et al. (2017), the mathematical background is not well addressed. In
75 this study, we explain in detail the mathematics related to this fast algorithm
76 and discuss its asymptotic properties.

77 In the following sections, we first introduce the biscale empirical prob-
78 ability integral transformation (BEPIT), and then analyze the completely
79 observed process with time-separable marks after the transformation. Based
80 on the results of this transformation, we restore the empirical distributions
81 from an incomplete record using an iterative algorithm. We explain the
82 algorithm using a simulated data set. Finally, we apply the algorithm to in-
83 vestigate the incomplete eruption record of the Hakone volcano in Japan, and
84 the aftershock sequence of the Wenchuan $Mw7.9$ earthquake that occurred
85 in Southwest China on May 28, 2008. The proofs of the consistency and
86 asymptotic normality of the algorithm are given in Supplementary Material.

87 **2 Concepts, methodology, and illustration**

88 **2.1 Mark-separable temporal point process and BEPIT**

89 Mathematically, a marked temporal point process N is a random subset of
90 discrete points on the space $\mathbb{R} \times \mathbb{M}$, say $\{(t_i, m_i) : i = 1, 2, \dots, n\}$, which
91 includes a finite or countable number of elements, and satisfies the following
92 two conditions (Karr, 1991): (a) for any bounded subset $A \subset \mathbb{R}$, $\Pr\{N(A \times$
93 $\mathbb{M}) \equiv \#[N \cap (A \times \mathbb{M})] < \infty\} = 1$, where $\#[\]$ represents the number of
94 elements in a set; and (b) for each i , m_i is a random variable on \mathbb{M} . In

95 our study, we assume the following: (a) the marks are continuous random
96 variables, and (b) the point process is simple (i.e., $\Pr\{\max_{t \in \mathbb{R}} N(\{t\} \times \mathbb{M}) \leq$
97 $1\} = 1$), such that there are no overlapping events on the time axis.

98 A marked temporal point process is often specified by its conditional
99 intensity function, defined by

$$\lambda(t, m) dt dm = \mathbf{E} [N([t, t + dt) \times (m, m + dm) | \mathcal{H}_t)], \quad (1)$$

where \mathcal{H}_t denotes the history of N up to time t , but not including t . The conditional intensity can be decomposed as

$$\lambda(t, m) = \lambda_g(t) g(m|t),$$

100 where $\lambda_g(t) = \int_{\mathbb{M}} \lambda(t, m) dm$ is called the conditional intensity of the ground
101 point process N_g induced by N on \mathbb{R} , defined by $N_g(A) = N(A \times \mathbb{M})$, and
102 $g(m|t)$ is the probability density function of the event mark at time t . An
103 important property of the conditional intensity is that if a temporal point
104 process N has conditional intensity $\lambda(t)$, then the transformation

$$t_i \rightarrow \tau_i = \int_0^{t_i} \lambda(u) du \quad (2)$$

105 transforms N into a Poisson process $N' = \{\tau_i : i = 1, 2, \dots\}$ (see, e.g.,
106 Ogata, 1988; Schoenberg, 2003; Daley and Vere-Jones, 2003).

107 For the above conditional intensity, when the mark distribution is sepa-
108 rable from the occurrence times, that is,

$$\lambda(t, m) = \lambda_g(t) g(m), \quad (3)$$

109 the marks of this point process are said to be time separable. Point-process
110 models with time-separable marks are widely used in many research areas. In

111 seismology, most practical versions of earthquake forecasting models explic-
112 itly assume that the magnitude distribution is separable from time (see, e.g.,
113 Ogata and Zhuang, 2006; Zhuang et al., 2002, 2004; Zhuang, 2011; Werner
114 et al., 2011; Ogata et al., 2013). In volcanology, Bebbington (2014) suggested
115 that there is not enough evidence of a universal dependence of eruption size
116 on time. In forecasting, time-independent size distributions are used fre-
117 quently (e.g., Passarelli et al., 2010).

118 Other ways to specify point-process models include moment intensity
119 functions, Papangelou intensities, and Palm intensities. Traditionally, when
120 a point process is specified in one of these ways, it refers to a spatial point
121 process. A point process can be completely determined by its likelihood
122 (terminologically, the local Janossy density; see Daley and Vere-Jones, 2003,
123 2008). This gives the joint probability density/mass function of the total
124 number and each location of the particles in the process, assuming that the
125 particles are indistinguishable. The likelihood is also known (i.e., the point
126 process is completely determined) if one of the following three is known: (1)
127 the moment intensities of all orders, (2) the conditional intensity, and (3) the
128 Papangelou intensity. Here, we refer to Daley and Vere-Jones (2003, 2008)
129 and Møller and Waagepetersen (2003) for the relations between the Janossy
130 density and three other types of intensities. In this study, the method used
131 to replenish missing data in a marked point process does not depend on
132 any specific form of conditional intensity. Therefore, it can be applied to
133 spatial point processes as well if the ground space is one dimensional and the
134 conditional intensity is mark separable.

Before testing for missing data in a record of a marked point process and
replenishing the record, we need to know what a complete record looks like.

Given a series of independent and identically distributed (i.i.d.) observations on X , x_1, x_2, \dots, x_n , for a fixed x , the empirical cumulative distribution function (cdf)

$$\tilde{F}_X(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(x_i < x)$$

135 converges almost surely to $F_X(x)$ and, thus, $\tilde{F}_X(X_j)$, for $j = 1, 2, \dots, n$, con-
136 verges to a unit uniform distribution. We call transformation $x \rightarrow \tilde{F}_X(x)$ the
137 empirical probability integral transformation induced by $\{x_1, x_2, \dots, x_n\}$. In
138 a general marked point process N in $[0, T]$, the occurrence times of an ar-
139 bitrary event may depend on the occurrence times and/or marks of other
140 events. However, the empirical probability integral transformation still re-
141 sults in an approximate unit uniform distribution, because the transformation
142 does not require an explicit formulation of the conditional intensity.

143 Suppose $N = \{(t_i, m_i) : i = 1, 2, \dots, n\}$ is a realization of a temporal
144 marked point process in a time–mark domain $[0, T] \times \mathbb{M}$, where \mathbb{M} is the
145 space of marks. Consider the following BEPIT:

$$\begin{aligned} \Gamma_N : [0, T] \times \mathbb{M} &\rightarrow [0, 1] \times [0, 1] \\ (t, m) &\rightarrow (t', m') = (\tilde{F}(t), \tilde{G}(m)), \end{aligned} \quad (4)$$

146 where \tilde{F} and \tilde{G} are the empirical cdfs of $\{t_i : i = 1, 2, \dots, n\}$ and $\{m_i : i = 1, 2, \dots, n\}$, respectively. If the marks of the events in the process
147 are separable from the occurrence times, then $\{t'_i : i = 1, 2, \dots, n\}$ and
148 $\{m'_i : i = 1, 2, \dots, n\}$, which are the images of $\{t_i : i = 1, 2, \dots, n\}$ and $\{m_i : i = 1, 2, \dots, n\}$, respectively, form an approximately homogeneous Poisson
149 process on $[0, 1] \times [0, 1]$. It is straightforward to show the independence
150 between $\tilde{F}(t)$ and $\tilde{G}(m)$. Thus, given the total number of events N , the
151 number of events in a cell of area $s \subseteq [0, 1] \times [0, 1]$ is a random variable from
152
153

154 a binomial distribution $B(N, s)$, which can be approximated by a Poisson
155 distribution with mean Ns . The smaller s gets, the better this approximation
156 becomes.

157 In the following discussion, we consider only mark-separable Poisson
158 processes. This is because we can transform a more general process, say
159 N , with a conditional intensity $\lambda(t, m)$, into a Poisson process N' with a
160 constant intensity using the marked version of the transformation in (2),
161 $(t_i, m_i) \in N \rightarrow (\tau_i, m_i) \in N'$, where $\tau_i = \int_0^{t_i} \int_{\mathbb{M}} \lambda(t, m) dm dt$. Because such
162 a transformation does not change the chronological order of the events or
163 the mark-separable property of the process, the BEPIT transforms N and
164 N' into the same point patterns.

Example 1. *In Figure 1(a), we simulate a Poisson process N (the combination of dots and crosses) with a temporal rate $\lambda = 1$ on $[0, 2,000]$, and marks that follow an exponential distribution with mean one; that is,*

$$g(x) = \begin{cases} e^{-x}, & x > 0; \\ 0, & \text{otherwise.} \end{cases}$$

165 *Figure 1(b) shows that under transformation (4), N is transformed into an*
166 *approximately homogeneous Poisson process, say N' , which has rate $\lambda =$*
167 *2,000 and i.i.d. marks uniformly distributed in $[0, 1]$.*

168 2.2 Detection of missing data

169 When events in part of an observed time–mark range are missing, determin-
170 istically or in probability, the separability between the occurrence times and
171 the marks of the observed events is usually destroyed. In addition, the image
172 of the observed N_{obs} mapped by the above BEPIT $\Gamma_{N_{\text{obs}}}$, as defined in (4),
173 may not be a homogeneous process.

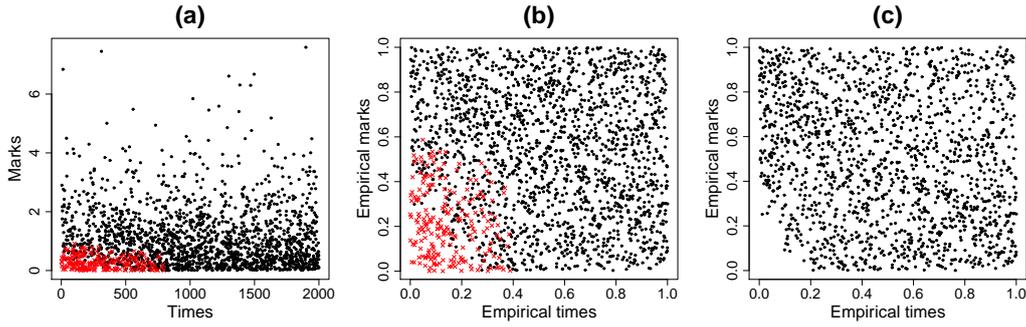


Figure 1: A synthetic data set of a marked point process. (a) Marks versus occurrence times. (b) Empirical marks versus empirical occurrence times of all synthetic events under the transformation Γ_N . (c) Empirical marks versus empirical occurrence times for the observed incomplete record under the transformation $\Gamma_{N_{\text{obs}}}$. The crosses in (a) and (b) represent the missing events.

174 **Example 2.** Consider the simulated data in Example 1 (Figures 1(a)). As-
175 sume the missing probability is

$$\begin{aligned}
 q(t, m) &= \Pr\{\text{an event occurring at } (t, m) \text{ is missing}\} \\
 &= \begin{cases} \min\left[1, \frac{(1000-t)(1-m)}{800}\right], & \text{if } 0 < t < 800, m < 0.3, \\ 0, & \text{otherwise.} \end{cases} \quad (5)
 \end{aligned}$$

176 If we thin the original process N (the combination of the crosses and dots)
177 in Figure 1(a) with this missing probability, then the crosses are deleted (i.e.,
178 they are missing from the record). Denote the remaining events (i.e., the
179 observed process) as N_{obs} . Figure 1(c) shows that the image of the observed
180 data of the process under the BEPIT $\Gamma_{N_{\text{obs}}}$ is not homogeneous.

181 In the above bi-scale transformation, we do not need to know the exact
182 forms of $g(m)$, λ_g , or q . This method relies only on the conditions that the
183 original process is mark separable, and that the process of missing events is
184 time- and mark-dependent. Thus, for a temporal point process N with time-
185 separable marks, we can test whether data are missing from its observed
186 record, N_{obs} , by testing the homogeneity of the image $\Gamma_{N_{\text{obs}}}(N_{\text{obs}})$ of the

187 observed data N_{obs} in the bi-scale transformed domain, when the missing
188 values are time- and mark-dependent. After using the BEPIT $\Gamma_{N_{\text{obs}}}$ to map
189 N_{obs} onto $[0, 1]^2$, we divide the overall area of $[0, 1]^2$ into L sub-regions of
190 equal areas, that is, $L = L_1 \times L_2$ cells. Here, L_1 is the number of cells
191 along the transformed time domain, and L_2 is the number of cells along the
192 transformed mark domain. Then, we calculate the following statistics:

$$R = \frac{\min\{C_1, C_2, \dots, C_L\}}{\max\{C_1, C_2, \dots, C_L\}}, \text{ and } D = \max\{C_1, C_2, \dots, C_L\} - \min\{C_1, C_2, \dots, C_L\},$$

(6)

193 where C_1, C_2, \dots, C_L are the numbers of events falling within each of the L
194 cells. These two statistics are analogous to the test statistics for homogeneous
195 multinomial distributions, where “homogeneous” means that each category
196 of the possible outputs has the same probability (Johnson, 1960; Johnson
197 and Young, 1960; Corrado, 2011).

198 Suppose that $[0, 1]^2$ is divided into $L = L_1 \times L_2$ cells with equal areas;
199 that is, $[0, 1]^2 = \bigcup_{j=1}^{L_2} \bigcup_{i=1}^{L_1} [(i-1)/L_1, i/L_1] \times [(j-1)/L_2, j/L_2]$, where L_1
200 and L_2 are positive integers. For any point process N on $[0, 1]^2$, if N is
201 a homogenous Poisson process, then the numbers of events in the above L
202 cells, C_1, C_2, \dots, C_L , form a homogeneous (n, \mathbf{p}) -multinomial random vector,
203 with $\mathbf{p} = (1/L, 1/L, \dots, 1/L)$. However, if N is obtained by applying the
204 BEPIT to a completely observed mark-separable point process, then the
205 row sum of C_i in the k th row ($1 \leq k \leq L_1$) and the column sum of C_i in
206 the j th column ($1 \leq j \leq L_2$) are fixed to $\lfloor kn/L_1 \rfloor - \lfloor (k-1)n/L_1 \rfloor$ and
207 $\lfloor jn/L_2 \rfloor - \lfloor (j-1)n/L_2 \rfloor$, respectively, where $\lfloor x \rfloor$ denotes the integer part of
208 x , and n is the total number of events in N . Such constraints do not hold for
209 the homogeneous multinomial distribution. Because the distributions of R

210 and D are complicated, we obtain them by simulation, as follows: (1) with n
211 fixed, simulate n events uniformly distributed in $[0, 1]^2$; (2) apply the BEPIT
212 to these n simulated events; (3) with the specified parameters, L_1 and L_2 ,
213 calculate R and/or D for the transformed points.

214 **Example 3.** *We use a simulation to test for missing data in the original and*
215 *the thinned point processes, as shown in Figures 1(a) and (c), respectively.*
216 *We simulate 500,000 sequences of the marked Poisson process defined in*
217 *Example 1, with the number of events in each simulation the same as those*
218 *in Figure 1(a). For each simulated sequence, we apply the BEPIT in (4),*
219 *which results in an image similar to the combination of the crosses and the*
220 *dots in Figure 1(b). Then, we divide the unit square image into five-by-five*
221 *cells with equal sizes, and calculate R and D , as defined in (6). Next, we plot*
222 *the empirical cdf of the 500,000 values of R and D , as shown in Figure 2(a).*
223 *To test the thinned process, we simulate further 500,000 sequences of the*
224 *marked point process, with the total number of events in each simulation the*
225 *same as those in Figure 1(c). The cumulative distributions of R and D*
226 *are shown in 2(b). We can see that the hypothesis of no missing data in*
227 *the observed (thinned) process can be rejected, with a significance level below*
228 *0.001 ($p \leq 2 \times 10^{-6}$, Figure 2(b)). Meanwhile, for the original process, the*
229 *p -values associated with R and D (0.396 and 0.700, respectively) provide no*
230 *evidence for rejection.*

231 **2.3 Imputation method and algorithm**

232 We start with a heuristic example to explain the algorithm. As shown in
233 Figure 3, suppose that N is a homogeneous point process on $[0, 1]^2$, and that
234 events in the domain S are completely unobservable. Let $N_{\text{obs}} = \{(x_i, y_i) :$

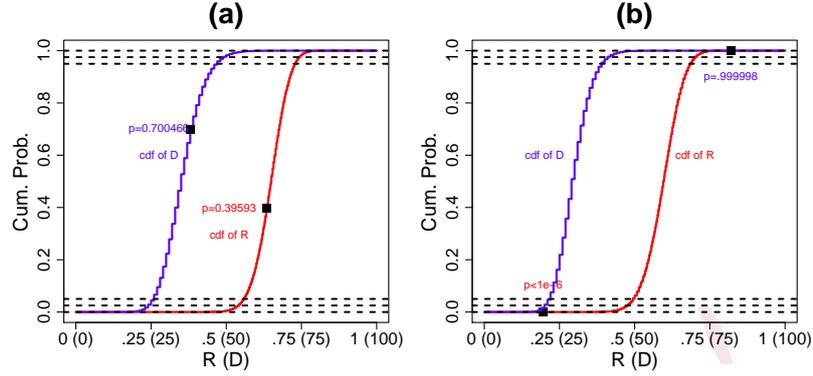


Figure 2: Statistical tests of the existence of missing data on (a) all events and (b) the observed events in the synthetic point process, with cdfs of R and D . R and D are defined in (6), with $L = L_1 \times L_2$, $L_1 = L_2 = 5$. The cdfs in (a) and (b) are obtained from 500,000 simulations with the same numbers of events as in Figures 1(a) and (c), respectively. The black dots in (a) and (b) are the statistics R and D , calculated for the original process in 1(a) and (c), respectively.

235 $(x_i, y_i) \in N \setminus S$. Then, the empirical distributions of the x - and y -coordinates
236 are, respectively,

$$\tilde{F}_X(x) = \frac{\sum_{i:(x_i, y_i) \in N \setminus S} w_{x,i} I(x_i \leq x)}{\sum_{i:(x_i, y_i) \in N \setminus S} w_{x,i}} \quad (7)$$

237 and

$$\tilde{F}_Y(y) = \frac{\sum_{i:(x_i, y_i) \in N \setminus S} w_{y,i} I(y_i \leq y)}{\sum_{i:(x_i, y_i) \in N \setminus S} w_{y,i}}, \quad (8)$$

238 where

$$w_{x,i} = \frac{1}{1 - \int_0^1 I((x_i, y) \in S) dy}, \quad w_{y,i} = \frac{1}{1 - \int_0^1 I((x, y_i) \in S) dx}. \quad (9)$$

239 In most cases, N is not homogeneous in $[0, 1]^2$, and the variation of the
240 event density in S should be considered. Equation (9) should then be

$$w_{x,i} = \frac{1}{1 - \int_0^1 I((x_i, y) \in S) dF_y(y)}, \quad w_{y,i} = \frac{1}{1 - \int_0^1 I((x, y_i) \in S) dF_x(x)}. \quad (10)$$

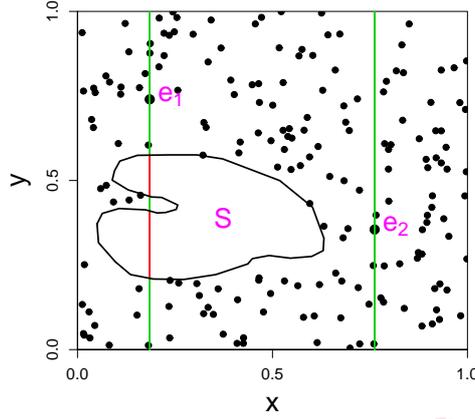


Figure 3: A heuristic estimation of the empirical distribution with missing points. Suppose that, among events $e_i = (x_i, y_i)$, for $i = 1, 2, \dots, N$, events that fall in S cannot be observed. To estimate the empirical distribution $\tilde{F}_X(x)$ of x_i , for $i = 1, 2, \dots, N$, weights need to be assigned to each observed point. That is, when N is uniform, $\tilde{F}_X(x) = \sum_{i=1}^N w_{x,i} I(x_i < x) / \sum_{i=1}^N w_{x,i}$, where $w_{y,i} = 1 - \int_0^1 I((x_i, y) \in S) dy$. In this figure, $w_{x,1}$ is the total length of the green part of the vertical line segments crossing over e_1 , and $w_{x,2} = 1$ because the vertical line segment crossing e_2 has no intersection with S .

241 Because F_Y and F_X are unknown, we replace them with \tilde{F}_Y and \tilde{F}_X , respec-
242 tively; that is,

$$w_{x,i} = \frac{1}{1 - \int_0^1 I((x_i, y) \in S) d\tilde{F}_y(y)}, \quad w_{y,i} = \frac{1}{1 - \int_0^1 I((x, y_i) \in S) d\tilde{F}_X(x)}. \quad (11)$$

243 The above equation, together with (7) and (8), form a solvable equation
244 system. Below we propose an algorithm to solve this equation system.

245 First, the missing region S needs to satisfy the following condition.

246 **Condition 1.** *The projections of $([0, T] \times M) \setminus S$ (i.e., the sub-region in*
247 *which no event is missing) on the t - and m -axes cover the entire observation*
248 *period and the entire range of possible marks, respectively.*

249 This requirement ensures that the empirical distributions of $\{t_i\}$ and $\{m_i\}$
250 can be restored. With Condition 1 satisfied, when a record is incomplete, we
251 can determine the area, say S , outside of which the record is complete. This

252 can be done either in the original time-mark plot, based on prior knowledge
253 of the data quality, or in the BEPIT domain, based on the statistics R or D .

254 The algorithm to replenish the record includes three key steps: (1) trans-
255 forming the process onto $[0, 1]^2$ using the BEPIT to find a time-mark range
256 that likely contains all missing events; (2) estimating a new empirical dis-
257 tribution function based on the data in the time-mark range, inside which
258 events are supposed to be completely observed; (3) generating events in the
259 missing region.

260 **Initial settings.** Given the data set $N_{\text{obs}} = \{(t_i, m_i) : i = 1, 2, \dots, n\}$
261 observed in $[0, T] \times M$ and a time-mark range S , known to include the
262 missing events, suppose that S satisfies Condition 1.

263 *Step 1.* We project the observed data and the range S that contains the
264 missing data onto $[0, 1]^2$ using the BEPIT in (4). Explicitly, set

$$(t_i^{(1)}, m_i^{(1)}) = \Gamma_{N_{\text{obs}}}^{(1)}(t_i, m_i), \quad (12)$$

265 where

$$\Gamma_{N_{\text{obs}}}^{(1)}(t, m) = \left(\tilde{F}^{(1)}(t), \tilde{G}^{(1)}(m) \right) = \left(\frac{1}{n} \sum_{j=1}^n \mathbf{1}(t_j < t), \frac{1}{n} \sum_{j=1}^n \mathbf{1}(m_j < m) \right). \quad (13)$$

266 Denote $S^{(1)}$ as the image of S under the transformation $\Gamma_{N_{\text{obs}}}^{(1)}$.

267 *Step 2.* Starting from $\ell = 1$, repeat the following iterative computation until
268 convergence (e.g., $\max\{|t_i^{(\ell+1)} - t_i^{(\ell)}|, |m_i^{(\ell+1)} - m_i^{(\ell)}|\} < \epsilon$), where ϵ is a
269 given small positive number.

$$(t_i^{(\ell+1)}, m_i^{(\ell+1)}) = \Gamma_{N_{\text{obs}}}^{(\ell+1)}(t_i^{(\ell)}, m_i^{(\ell)}; S^{(\ell)}), \quad i = 1, 2, \dots, n, \quad (14)$$

270

$$S^{(\ell+1)} = \Gamma_{N_{\text{obs}}}^{(\ell+1)}(S^{(\ell)}; S^{(\ell)}), \quad (15)$$

271 where

$$\Gamma_{N_{\text{obs}}}^{(\ell+1)}(t, m; A) = \left(\frac{\sum_{j=1}^n w_1^{(\ell)}(t_j^{(\ell)}, m_j^{(\ell)}, A) \mathbf{1}(t_j^{(\ell)} < t)}{\sum_{j=1}^n w_1^{(\ell)}(t_j^{(\ell)}, m_j^{(\ell)}, A)}, \frac{\sum_{j=1}^n w_2^{(\ell)}(t_j^{(\ell)}, m_j^{(\ell)}, A) \mathbf{1}(m_j^{(\ell)} < m)}{\sum_{j=1}^n w_2^{(\ell)}(t_j^{(\ell)}, m_j^{(\ell)}, A)} \right) \quad (16)$$

272 with the weights defined as

$$w_1^{(\ell)}(t, m, A) = \frac{\mathbf{1}((t, m) \notin A)}{1 - \int_0^1 \mathbf{1}((t, m') \in A) dG^{(\ell)}(m')} \quad (17)$$

$$w_2^{(\ell)}(t, m, A) = \frac{\mathbf{1}((t, m) \notin A)}{1 - \int_0^1 \mathbf{1}((t', m) \in A) dF^{(\ell)}(t')} \quad (18)$$

273 for any regular region $A \subset [0, 1]^2$. Denote the results upon convergence

274 as $N_{\text{obs}}^* = \{(t_i^*, m_i^*) : i = 1, 2, \dots, n\}$ and S^* .

Step 3. Generate a random number K from a negative binomial distribution, with parameters $(k, 1 - |S^*|)$, where $|S^*|$ is the area of S^* and

$$k = \sum_{i=1}^n \mathbf{1}((t_i^*, m_i^*) \notin S^*) = \#(N_{\text{obs}}^* \setminus S^*).$$

276 Step 4. Generate K random events independently, identically, and uni-
277 formly distributed in S^* . Denote these newly generated events as N_{rep}^* .

278 Step 5. For each event in N_{obs}^* , say, (t_j, m_j) , that falls in S^* , sequentially
279 remove from N_{rep}^* the event that is the closest to (t_j, m_j) .

280 Step 6. Convert the resulting N_{rep}^* from the last step to the original obser-
281 vation space $[0, T] \times M$ through linear interpolation:

$$s_j = \text{LI}(s_j^*; [0, t_1^*, t_2^*, \dots, t_n^*, 1], [0, t_1, t_2, \dots, T]), \quad (19)$$

$$v_j = \text{LI}(v_j^*; [0, m_1^*, m_2^*, \dots, m_n^*], [0, m_1, m_2, \dots, m_n]), \quad (20)$$

283 for each $(s_j^*, v_j^*) \in N_{\text{rep}}^*$, where $\text{LI}(x, A, B)$ represents the linear interpo-

284 lation value of x , conditional on the function values for each component

285 in A being locations corresponding to each component in B . Denote the
286 set consisting of all (s_j, v_j) as N_{rep} .

287 **Final output.** Return N_{rep} .

288 **Example 4.** Here we apply the above algorithm to the thinned data set in
289 Example 2. The output from Steps 4 to 6 is shown in Figures 4(b)–(c). The
290 final output for our simulation example is shown in Figure 4(d). The tests
291 using statistics R and D in (6) give p -values of 0.605 and 0.718, respectively,
292 providing no evidence to reject the hypothesis that the replenished data set is
293 complete (Figure 4(e)). Figure 4(f) compares the cumulative numbers of
294 events in the original, observed, and replenished processes, showing that the
295 replenishing algorithm recovers the missing data to some extent.

296 **Notes:**

297 (1) Equation (13) is the BEPIT mentioned in the previous section. If the
298 data are completely recorded, $\{(t_i^{(1)}, m_i^{(1)}), i = 1, 2, \dots, n\}$ form an ap-
299 proximately homogeneous process on $[0, 1]^2$. As shown in Figure 2(b),
300 the sparseness of the points around the lower, left corner implies that
301 smaller events are missing in the earlier period. Rather than choosing
302 S in Figure 1(a), it is more convenient to specify $S^{(1)}$ directly in Figure
303 2(a) or (b).

304 (2) Step 2 is carried out based on the fact that the transformation $\Gamma_{N_{\text{obs}}}$ and
305 $S^{(1)} = \Gamma_{N_{\text{obs}}}(S)$ can be quite different from Γ_N , owing to the missing
306 data. The iteration in this step helps us construct a bi-scale transfor-
307 mation as close as possible to the BEPIT yielded by the complete data

308 (i.e., $\Gamma_{N_{\text{obs}}}^* \approx \Gamma_N$). At the same time, the corresponding area that con-
309 tains the missing data, S^* , is restored. This can be seen by comparing
310 Figures 1(b) and 4(b).

311 Step 2 essentially solves F^* and G^* in the following equations:

$$F^*(t) = \frac{\sum_{j=1}^n w_1(t_j, m_j, S) \mathbf{1}(t_j < t)}{\sum_{j=1}^n w_1(t_j, m_j, S)}, \quad (21)$$

$$G^*(m) = \frac{\sum_{j=1}^n w_2(t_j, m_j, S) \mathbf{1}(m_j < m)}{\sum_{j=1}^n w_2(t_j, m_j, S)}, \quad (22)$$

313 where

$$w_1(t, m, S) = \frac{\mathbf{1}((t, m) \notin S)}{1 - \int_M \mathbf{1}((t, m') \in S) dG^*(m')} \quad (23)$$

$$w_2(t, m, S) = \frac{\mathbf{1}((t, m) \notin S)}{1 - \int_M \mathbf{1}((t', m) \in S) dF^*(t')}. \quad (24)$$

315 If we define $\Gamma_{N_{\text{obs}}}^*(t, m) = (F^*(t), G^*(m))$ as a mapping from $[0, T] \times M$
316 to $[0, 1]^2$, then $\Gamma_{N_{\text{obs}}}^*(t, m)$ directly maps N_{obs} to N_{obs}^* and S to S^* .

317 The existence of a solution in the iteration given by (21) to (24) and
318 the asymptotic property of the solution are given in the Supplementary
319 Material.

320 (3) Steps 3 and 4 are based on the following fact: given a homogeneous
321 Poisson process with an unknown occurrence rate, if there are k events
322 falling within an area of S_1 , then the number of events falling in the
323 complementary area, S_2 , follows a negative binomial distribution with
324 parameter $(k, |S_1|/(|S_1| + |S_2|))$ (e.g., DeGroot, 1986, 258–259).

325 (4) In Step 5, we should keep the existing events observed in S , and remove
326 the same number of simulated points.

327 One advantage of the algorithm is that if S is unknown, we can use the time-
328 mark plot of $N^{(1)}$, as in Figure 2(b), to determine $S^{(1)}$ by justifying which
329 region is likely to contain the missing events, and then continue with Step 2.
330 Once the replenishment is complete, S can be obtained by substituting the
331 coordinate of each point on the boundary of S^* into (19) and (20).

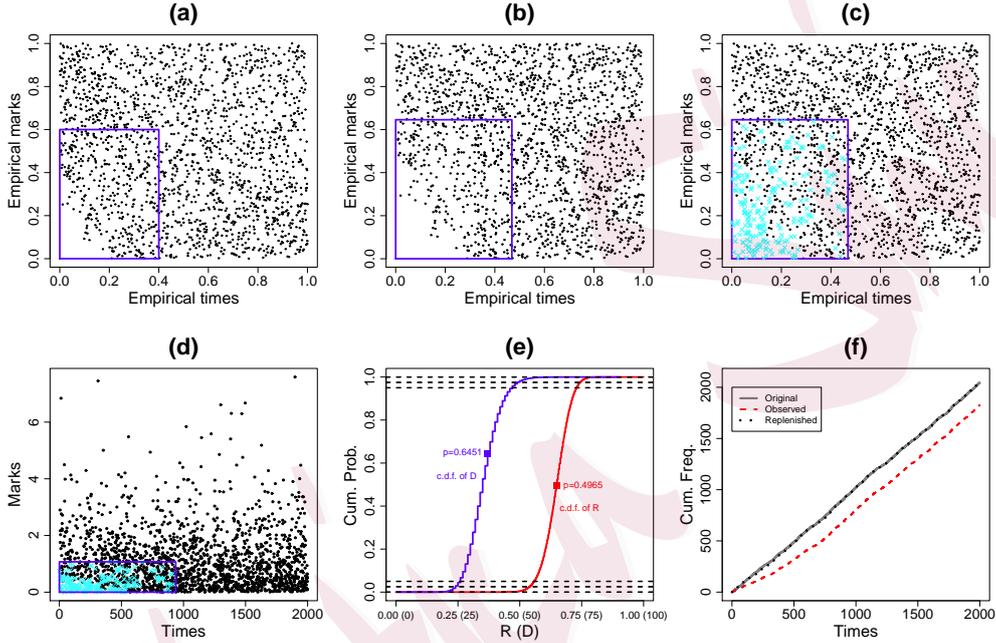


Figure 4: An application of the proposed replenishing algorithm to the synthetic data set. (a) Rescaled marks versus rescaled occurrence times of the observed events (dots), with the bi-scale transformation $\Gamma_{N_{\text{obs}}}$ based on the observed process. The polygon is the missing area, $S^{(1)}$. (b) Rescaled marks versus rescaled occurrence times of the observed events (dots), with the rescaling $\Gamma_{N_{\text{obs}}^*}$ based on the events outside of S . The polygon is the missing area after transformation $\Gamma_{N_{\text{obs}}^*}$, that is., S^* . (c) Rescaled marks versus rescaled occurrence times of the observed and replenished events (crosses) (i.e., newly generated events after removing events that are closest to any of those observed in S , with the rescaling $\Gamma_{N_{\text{obs}}^*}$ based on the empirical distributions of the events outside S . (d) Marks versus occurrence times of the observed synthetic events and the replenished events. (e) Cdfs of R and D for testing missing data in the replenished data set in (c). (f) Cumulative frequencies versus occurrence times for the original, observed, and replenished processes.

332 2.4 Additional simulations

333 To illustrate the overall behavior of the above replenishing algorithm, we
334 repeat the algorithm many times, with S fixed, for the following two cases:
335 (1) simulating a Poisson process with $\lambda = 2,000$; and (2) simulating Poisson
336 processes with rate λ , drawn from a uniform distribution within $[100, 3,000]$.
337 Both simulations have the same missing probability functions, as given by (5).
338 Figures 5(a) and (b) compare between the true numbers of missing events
339 and the numbers of replenished events for cases (1) and (2), respectively. In
340 Figure 5(a), since λ is fixed, the number of replenished events is independent
341 of the true number of missing events, and has a larger variance. Several
342 statistics related to these simulations are given in Table 1, including the
343 mean numbers and the variances of the missing and replenished points, the
344 mean of the relative differences, and the relative difference between the means
345 in 500 and 2,000 simulations. In particular, the near-zero relative deviation
346 of the mean number of replenished events shows that the proposed method
347 is consistent. Here, the larger values of the mean relative deviation of the
348 number of replenished events from the number of missing events illustrate
349 the nature of the uncertainty related to the problem. Such uncertainty is
350 produced not only by the randomness of the numbers of replenished and
351 missing events, but also by the uncertainty in the estimation of the occurrence
352 rate in the process from the events in the nonmissing part. In Figure 5(b),
353 the expected number of replenished events in many repeated simulations
354 is close to the number of missing events. Moreover, the relative deviation
355 decreases when the number of missing events (or λ) increases. These results
356 imply that this algorithm replenishes the missing events reasonably well. In
357 addition, when λ or the number of events in the process is quite small, some

Table 1: Statistics related to Figure 5(a). $\#m$: number of missing points; $\#r$: number of replenished points; $\bar{\cdot}$: mean value; $\sigma(\cdot)$: standard deviation.

#simu.	$\bar{\#m}$	$\sigma(\#m)$	$\bar{\#r}$	$\sigma(\#r)$	$\frac{ \#m-\#r }{\#m}$	$\frac{ \#m-\#r }{\#m}$
500	228.274	14.929	232.006	63.926	0.226	0.016
2000	227.712	14.719	230.860	62.145	0.224	0.014

358 outputs yield a negative number of replenished events (when the number of
359 missing events is less than 50 in Figure 5(b)). The number of replenished
360 events is calculated simply as the number of simulated events in S in Steps
361 3 and 4 minus the number of observed events in S . This finding indicates
362 that the existence of missing data in these situations cannot be quantified
363 probabilistically.

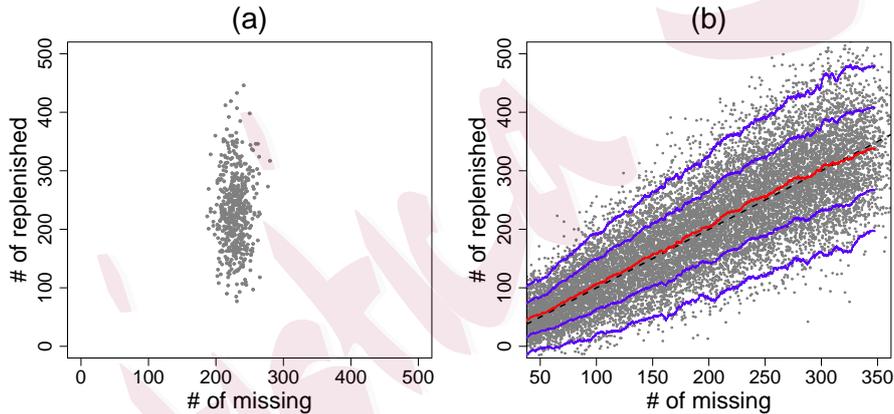


Figure 5: Comparison between the number of true missing events and the number of replenished events. (a) $\lambda = 2,000$, fixed. (b) λ is drawn from a uniform distribution between 100 and 3,000. The dashed line represents the case where the numbers of missing and replenished events are equal. The curves represent the running mean and the corresponding single and double standard deviation bands.

364 **3 Application**

365 **3.1 Volcanic eruption record**

366 In this example, we analyze a record of the eruptions from the Hakone
367 volcano, an active volcano located at the northern boundary zone of the

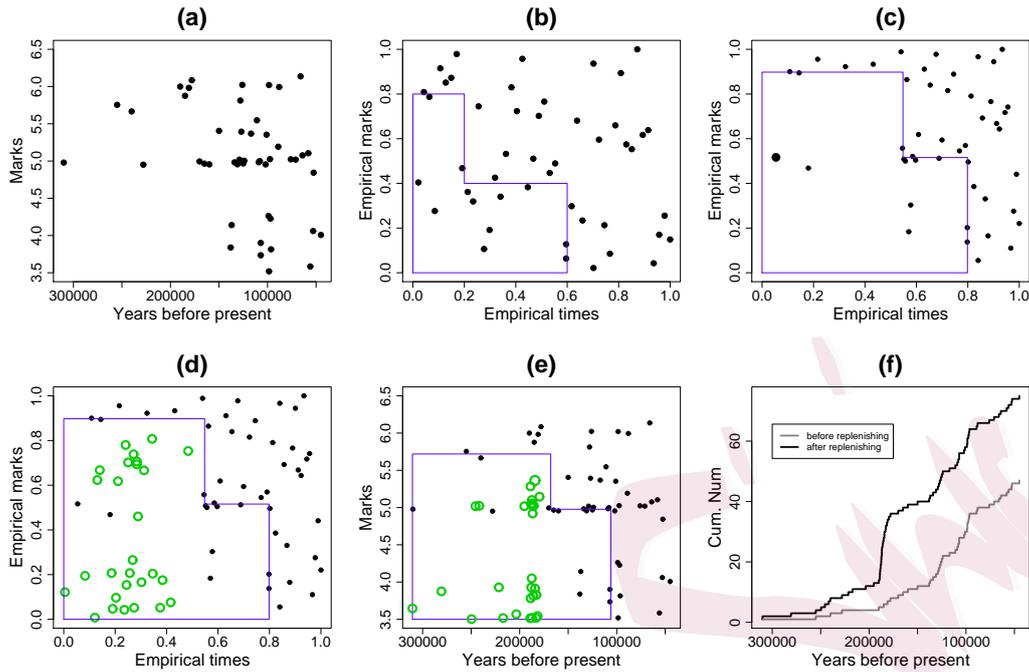


Figure 6: Results after applying the replenishment algorithm to volcanic eruption data. (a) Marks versus occurrence times of the eruption events. (b) Empirical distribution of marks versus that of occurrence times. (c) Rescaled marks versus rescaled occurrence times, with the rescaling based on the empirical distributions of the events outside of S . (d) Rescaled marks versus rescaled occurrence times of the observed and replenished events (i.e., newly generated events after removing events that are closest to any of those observed in S), with the rescaling based on the empirical distributions of the events outside of S . (e) Marks versus occurrence times of the observed and replenished events. (f) Cumulative numbers of events against occurrence times. The polygon is the area S and its corresponding mappings, in which the missing events fall. The green dots are the replenished events.

368 Izu-Mariana volcanic arc in central Japan (Yukutake et al., 2010; Honda
369 et al., 2014). The data on Japanese explosive eruptions are compiled from
370 the Smithsonian’s Global Volcanism Program database (Siebert and Simkin,
371 2002), the Large Magnitude Explosive Volcanic Eruptions database (LaMEVE
372 database, Croweller et al., 2012), and additional Japanese databases (Machida
373 and Arai, 2003; Committee for Catalog of Quaternary Volcanoes in Japan
374 (ed), 2000; Geological Survey of Japan, AIST (ed), 2013; Hayakawa, 2010).

375 For the Hakone volcano, 46 of the 54 compiled events have an eruption

376 magnitude ($M = \log_{10}[\text{erupted mass in kg}] - 7$; see Pyle (2015)) equal to or
377 larger than 4.0 (Table S1 in Supplementary material). Figure 6(a) shows
378 the eruption magnitudes versus occurrence times of these 46 events. Figure
379 6(b) shows the empirical distribution, transformed following Step 1 of the
380 algorithm. From this plot, the polygon boundaries of S are determined
381 based on the following assumptions. First, the events of empirical marks
382 < 0.8 ($M < 5.7$) are missing before the empirical time = 0.2 (165 ka).
383 Second, the recording of larger events improves after the empirical time =
384 0.2 (165 ka), although the events of empirical marks < 0.4 ($M < 5.0$) are
385 still missing. Third, the recording of events improves further and there are
386 no missing events after the empirical time = 0.6 (105 ka). The results of the
387 replenishing algorithm are shown in Figures 6(c) to 6(e).

388 The estimated cumulative number of events for the replenished data set
389 shows a remarkable jump of around 180 ka (Figure 6(f)). This jump is caused
390 by the replenished events synthesized around 180 ka (Figure 6(e)), based on
391 the cluster of four large events ($M \sim 6$) at 178 ka, 181 ka, 185 ka, and 190 ka
392 (Figure 6(a); Hayakawa, 2010). The ages of the events at the Hakone volcano
393 are still not fully agreed upon in the literature. For example, Yamamoto
394 (2015) assumed that the ages of the aforementioned eruptions are about 135
395 ka, 135 ka, 180 ka, and 215 ka, respectively. Therefore, the reliability of the
396 jump of the cumulative number of events (Figure 6(f)) might be problematic in
397 the volcanological dating of event ages, as in estimating the tephra volume
398 and rounded eruption magnitude in volcanology (Brown et al., 2014). For
399 example, the analyzed data set has clusters of events of magnitudes 4 and 5
400 (Figure 6(a)) and, therefore, the replenished events around 180 ka are also
401 clustered around magnitudes 4 and 5 (Figure 6(e)).

402 Note that it is difficult to determine the exact period of under-recording
403 in the eruption history of each volcano. Kiyosugi et al. (2015) showed that
404 many eruptions are still missing from the overall Japanese database, even for
405 the last 100,000 years. Therefore, the polygon shape (Figure 6(b)) we have
406 used suggests that our replenished data have the same completeness level
407 as that of the data outside the polygon. Our method is a way to consider
408 the under-recording of events in volcanic hazard assessments of explosive
409 eruptions using geological records.

410 **3.2 Earthquake catalog: Missing aftershocks**

411 It is well known that immediately after a large earthquake, many aftershocks
412 cannot be recorded, because the seismic waveforms generated by the after-
413 shocks, many of which occur in a short time after the mainshock, overlap with
414 each other and cannot be distinguished. In this section, we study the earth-
415 quake catalog from Southwest China for the period from January 1, 1990 to
416 April 20, 2013, in a space range of $26^{\circ} - 34^{\circ} N$ and $97^{\circ} - 107^{\circ} E$, with minimum
417 magnitude 3.0 (Figure S2 in the Supplementary Material). This data set is
418 selected from the Chinese Earthquake catalog compiled by the China Earth-
419 quake Data Center (CEDC) (URL: <http://data.earthquake.cn/index.html>).
420 The Wenchuan Mw 7.9 (Ms 8.0) earthquake, which occurred on May 12,
421 2008, was one of the two largest seismic events in China during the last 50
422 years. There are 6,249 events in the selected space and time range, of which
423 3,754 occurred after the Wenchuan earthquake, indicating a low seismicity
424 level above magnitude 3.0 in the study region prior to 2008. Many after-
425 shocks are missing immediately after the mainshock. In particular, events
426 of magnitudes between 3.0 and 4.0 are not properly recorded for a period

427 of about one-and-a-half months after the mainshock. The majority of the
428 events after May 12, 2008, can be taken as clustering events triggered by the
429 Wenchuan mainshock. When analyzing the seismicity in this area, Jia et al.
430 (2014) and Guo et al. (2015) chose a relatively high magnitude threshold of
431 4.0 to avoid biases in estimates caused by missing events. As a results, 5,217
432 of the 6,249 events had to be ignored.

433 This example is quite different from the previous example and that based
434 on the simulated data. The missing range can be well specified before replen-
435 ishment: the missing values are known immediately after the occurrence of
436 the mainshock, and the monitoring ability for events between magnitudes 3.0
437 and 4.0 are restored one-and-a-half months later. The results are illustrated
438 in Figure 7. We can see that missing events take up about half the total
439 number of events.

440 In seismology, the frequency of aftershock occurrences in an aftershock
441 sequence can be modeled by the empirical Omori–Utsu formula (e.g., Utsu
442 et al., 1995),

$$\lambda(t) = \frac{K}{(t + c)^p}, \quad (25)$$

443 where K is an index proportional to the number of earthquakes excited by
444 the mainshock, c is related to the period after the mainshock from which the
445 aftershock rate drops slowly, and p is the power related to the decay rate
446 of the aftershocks. Utsu et al. (1995) discussed how the parameters c and
447 p change with the cut-off magnitude threshold, and hypothesized that such
448 changes occur because small aftershocks in an early stage of the sequence
449 are missing from the catalog. We fit the above Omori–Utsu formula to both
450 the original and the replenished catalogs (Table 2), and obtain maximum

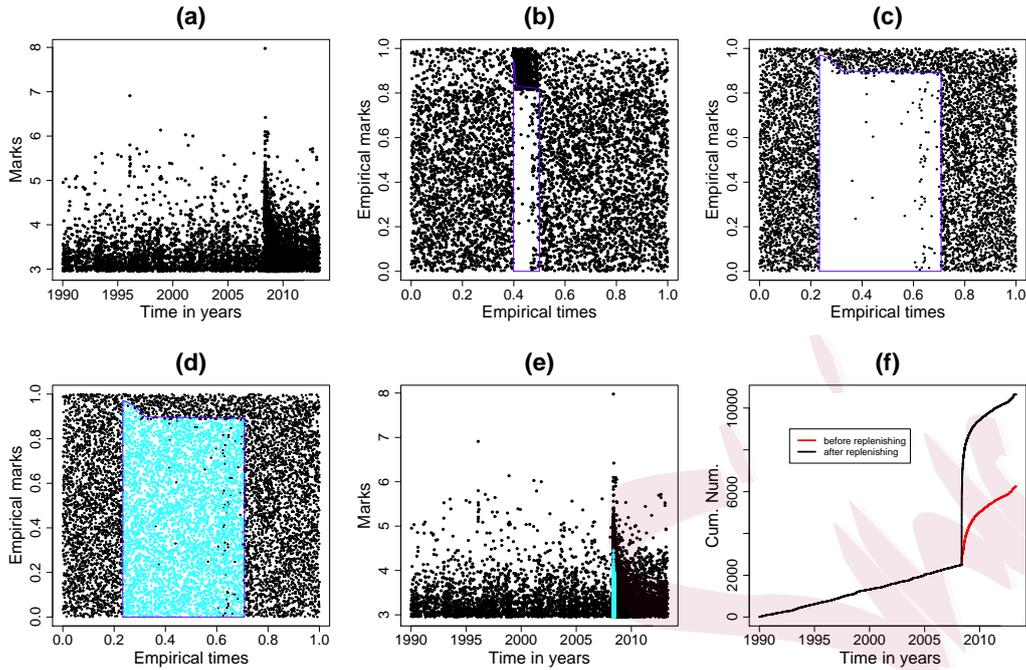


Figure 7: Results from applying the replenishment algorithm to earthquake data from Southwest China. (a) Marks versus occurrence times of the earthquake events. (b) Empirical distribution of marks versus that of occurrence times. (c) Rescaled marks versus rescaled occurrence times, with the rescaling based on the empirical distributions of the events outside of S . (d) Rescaled marks versus rescaled occurrence times of the observed and replenished events (i.e., newly generated events after removing events that are closest to any of the observed in S), with the rescaling based on the empirical distributions of the events outside S . (e) Marks versus occurrence times of the observed and replenished events. (f) Cumulative numbers of events against occurrence times. The polygon is the area S and its corresponding mappings, in which the missing events fall.

451 likelihood estimates of the parameters. The results show that after the re-
452 plenishment, the Omori parameters c and p no longer change. We also fit the
453 Omori–Utsu formula to the original data set, but only consider earthquakes
454 that occurred at least 54 days after the mainshock. In this case, although c
455 and p are slightly different from the estimates for the replenished data from
456 the starting time, they do not change much when the magnitude threshold
457 changes from 2.95 to 4.15 (Table S2 in the Supplementary Material). These
458 results numerically confirm the hypothesis of Utsu et al. (1995) that missing

Magnitude threshold	Replenished dataset [t_{main}, T]			Orig. dataset [t_{main}, T]		
	\hat{K}	\hat{c}	\hat{p}	\hat{K}	\hat{c}	\hat{p}
2.95	804.4	.1140	1.003	82.29	.0553	.6205
3.05	639.2	.1131	1.003	80.31	.0596	.6547
3.15	511.5	.1134	1.001	79.25	.0660	.6872
3.25	412.9	.1110	.9965	79.04	.0737	.7185
3.35	327.3	.1067	.9926	78.80	.0825	.7555
3.45	260.3	.1141	.9925	80.67	.0991	.7986
3.55	213.8	.1142	.9953	83.33	.1177	.8407
3.65	171.6	.1135	.9907	85.73	.1360	.8799
3.75	135.9	.1132	.9911	90.18	.1642	.9278
3.85	111.2	.1029	.9941	95.17	.1935	.9708
3.95	100.0	.1241	1.015	103.2	.2383	1.023
4.05	74.12	.1082	1.013	79.20	.1938	1.027
4.15	60.65	.1266	1.026	62.92	.1690	1.034

Table 2: Results from fitting the Omori–Utsu formula to the original and the replenished data sets of earthquakes from Southwest China, with different magnitude thresholds. t_{main} : occurrence time of the mainshock; T : end of the time interval.

459 small events in the early stage of an aftershock sequence causes the instability
460 of the estimate of the Omori–Utsu formula.

461 4 Conclusion

462 In this study, we proposed a method for replenishing missing data in marked
463 temporal point processes, based only on the assumption that the marks of the
464 events are separable from the occurrence times, regardless of how the events
465 interact on the time axis. The key point of this method is an algorithm that
466 iteratively estimates the missing area in the transformed domain, based on
467 the parts where the data are completely recorded. We applied the proposed
468 method to an eruption record of the Hakone volcano in Japan and to an
469 earthquake catalog from Southwest China, which includes the aftershock
470 zone of the 2008 Mw7.9 Wenchuan earthquake. The results show that the

471 proposed method helps to both evaluate the influence of missing data and
472 correct the bias caused by such data.

473 **Detection of the missing area** In our two examples, the missing area
474 is determined by visual inspection of the bi-scale transformed data for the
475 historical records of the Hakone volcano, and by prior information on the
476 seismic network for the Wenchuan aftershock sequence. In most cases, the
477 missing area is determined by the experience of data analysts or by infor-
478 mation on the data from other sources. However, it is possible to turn the
479 replenishing algorithm into an automated algorithm.

480 Starting from $S' = \emptyset$, we divide the unit square into small cells in the bi-
481 scale transformed domain, obtained by applying the transformation defined
482 in (9) to (13). Then, we carry out the statistical tests based on the statistics
483 R or D on the cells that do not intersect S' , as discussed in Section 3. If the
484 test shows that missing cells exist, then we merge these cells into S' . These
485 steps are iterated until no further cells are added to S' . Note that because
486 this topic belongs within the scope of data processing algorithms, we did not
487 include it in this paper.

488 **Separability of marks** As discussed earlier, the applicability of this al-
489 gorithm depends on whether the mark distribution is separable from the
490 occurrence time. If such dependence is known explicitly as a probability
491 density function, say $g(m | t)$, we can directly use the cdf that corresponds
492 to f in Steps 1 and 2 in the algorithm (i.e., $m_i^{(\ell)} = G(m_i | t_i)$, for $\ell \geq 1$).
493 Of course, such dependence should also be considered when transforming
494 the marks of replenished events from $[0, 1]$ to the original mark space. If

495 the mark is dependent on time, but we do not know how, together with the
496 existence of missing events, the replenishment/imputation problem becomes
497 unidentifiable.

498 Another case worth discussing is when the mark distribution is known
499 and does not depend on time. We can again use the cdf of the marks in
500 Steps 1 and 2 directly in the algorithm (i.e., by setting $m_i^{(\ell)} = G(m_i)$, for
501 $\ell \geq 1$). Such missing data can also be estimated using Bayesian methods, as
502 in Ogata and Katsura (1993), and then replenished by direct simulation.

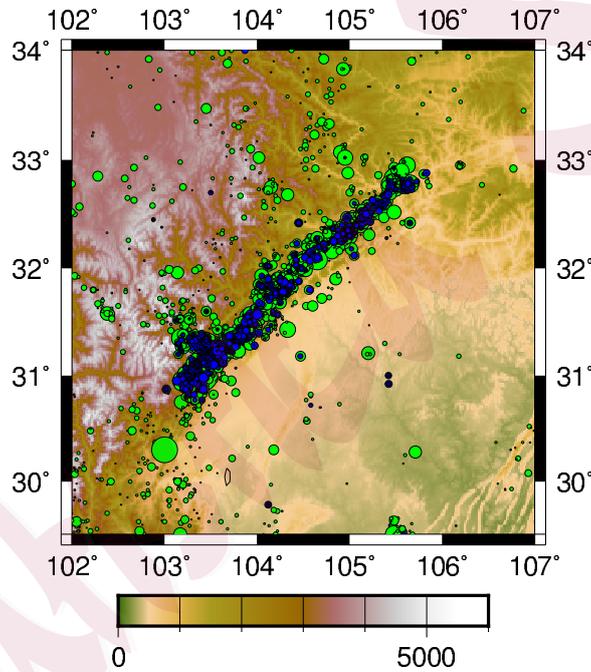


Figure 8: Epicenter map of imputed earthquakes (solid blue circles) for the Wenchuan aftershock sequence.

503 **Imputation of locations** This method is powerful for marked temporal
504 point processes, but it cannot be extended easily to high-dimensional or spa-
505 tiotemporal cases because, in most cases, the process is not homogeneous in
506 space. However, it is still possible case by case. For example, to replenish the
507 Wenchuan aftershock sequence, we can use the clustering feature of earth-

508 quakes. A simple replenishing algorithm is as follows. For each simulated
509 event, find a fixed number (e.g., 50) of events closest to it in time in the
510 observed process. Then, construct a Delaunay tessellation network for these
511 50 events, and select with equal probability one of the Delaunay triangles.
512 Lastly, place the the simulated event randomly and uniformly in this selected
513 triangle. An example of the imputed locations of the missing aftershocks of
514 the Wenchuan earthquake is shown in Figure 8. For a spatially inhibitive
515 process, different methods should be used.

516 In summary, the proposed method is useful when dealing with the missing
517 data problem in point-process observations, such as volcano eruption records
518 and historical or short-term earthquake catalogs.

519 **Supplementary Material**

520 The online Supplementary Material includes the following topics: (1) a proof
521 of the existence of a solution to the equation system in (21) to (24); (2)
522 the asymptotic properties of the solution; (3) additional simulations for the
523 case in which the missing region is wrongly specified; (4) list of the history
524 record of the Hakone volcano; and (5) comments on the Wenchuan aftershock
525 sequence.

526 **Acknowledgements**

527 This project was supported by the Royal Society of New Zealand Marsden
528 Fund (contact UOO1419). JZ was also partially supported by Grants-in-Aid
529 No. 2530052 for Scientific Research (C) from the Japan Society for the Pro-
530 motion of Science. Helpful discussions with David Harte from GNS Science

531 New Zealand and Boris Baeumer from the University of Otago are grate-
532 fully acknowledged. The authors also thank the AE and three anonymous
533 reviewers for their encouragement and constructive comments.

534 **References**

535 Bebbington, M. S. (2014). Long-term forecasting of volcanic explosivity.
536 *Geophysical Journal International*, 197:1500–1515.

537 Brown, S. K., Crosweller, H. S., Sparks, R. S. J., Cottrell, E., Deligne, N. I.,
538 Guerrero, N. O., Hobbs, L., Kiyosugi, K., Loughlin, S. C., Siebert, L., and
539 Takarada, S. (2014). Characterisation of the quaternary eruption record:
540 analysis of the large magnitude explosive volcanic eruptions (LaMEVE)
541 database. *J Appl Volcanol*, 3:5.

542 Committee for Catalog of Quaternary Volcanoes in Japan (ed) (2000). Cat-
543 alog of quaternary volcanoes in Japan (in Japanese).

544 Corrado, C. J. (2011). The exact distribution of the maximum, minimum
545 and the range of multinomial/dirichlet and multivariate hypergeometric
546 frequencies. *Statistics and Computing*, 21(3):349–359.

547 Crosweller, H., Arora, B., Brown, S. K., Cottrell, E., Deligne, N., Guerrero,
548 N., Hobbs, L., Kiyosugi, K., C., L. S., Lowndes, J., Nayembil, M., Siebert,
549 L., Sparks, R. S. J., Takarada, S., and Venzke, E. (2012). Global database
550 on large magnitude explosive volcanic eruptions (LaMEVE). *Journal of*
551 *Applied Volcanology*, 1:4.

552 Daley, D. D. and Vere-Jones, D. (2003). *An Introduction to Theory of Point*
553 *Processes – Volume 1: Elementary Theory and Methods (2nd Edition)*.
554 Springer, New York, NY.

- 555 Daley, D. D. and Vere-Jones, D. (2008). *An Introduction to Theory of*
556 *Point Processes – Volume II: General Theory and Structure (2nd Edition)*.
557 Springer, New York, NY.
- 558 DeGroot, M. H. (1986). *Probability and Statistics (Second ed.)*. Addison-
559 Wesley.
- 560 Diggle, P. J. and Rowlingson, B. S. (1994). A conditional approach to point
561 process modelling of elevated risk. *Journal of the Royal Statistical Society.*
562 *Series A (Statistics in Society)*, 157(3):433–440.
- 563 Enescu, B., Mori, J., and Miyazawa, M. (2007). Quantifying early aftershock
564 activity of the 2004 mid-Niigata prefecture earthquake ($M_w6.6$). *Journal*
565 *of Geophysical Research: Solid Earth*, 112(B4):B004629.
- 566 Enescu, B., Mori, J., Miyazawa, M., and Kano, Y. (2009). Omori-Utsu law
567 c -values associated with recent moderate earthquakes in Japan. *Bulletin*
568 *of the Seismological Society of America*, 99(2A):884–891.
- 569 Geological Survey of Japan, AIST (ed) (2013). Catalog of eruptive events
570 during the last 10,000 years in japan, version 2.1 (in japanese). Technical
571 report.
- 572 Guo, Y., Zhuang, J., and Zhou, S. (2015). An improved space-time
573 ETAS model for inverting the rupture geometry from seismicity trigger-
574 ing. *Journal of Geophysical Research: Solid Earth*, 120(5):3309–3323.
575 2015JB011979.
- 576 Gutenberg, B. and Richter, C. F. (1944). Frequency of earthquakes in Cali-
577 fornia. *Bull. Seis. Soc. Am.*, 34:184–188.

- 578 Hainzl, S. (2016). Rate-dependent incompleteness of earthquake catalogs.
579 *Seismological Research Letters*, 87(2A):337–344.
- 580 Hayakawa, Y. (2010). Hayakawafs 2000-year eruption database and one
581 million-year tephra database, <http://www.hayakawayukio.jp/database/>.
- 582 Honda, R., Yukutake, Y., Yoshida, A., Harada, M., Miyaoka, K., and Sato-
583 mura, M. (2014). Stress-induced spatiotemporal variations in anisotropic
584 structures beneath hakone volcano, japan, detected by s wave splitting:
585 A tool for volcanic activity monitoring. *Journal of Geophysical Research:*
586 *Solid Earth*, 119(9):7043–7057.
- 587 Iwata, T. (2008). Low detection capability of global earthquakes after the
588 occurrence of large earthquakes: Investigation of the Harvard CMT cata-
589 logue. *Geophysical Journal International*, 174(3):849–856.
- 590 Iwata, T. (2013). Estimation of completeness magnitude considering daily
591 variation in earthquake detection capability. *Geophysical Journal Interna-*
592 *tional*, 194(3):1909–1919.
- 593 Iwata, T. (2014). Decomposition of seasonality and long-term trend in seis-
594 mological data: A Bayesian modelling of earthquake detection capability.
595 *Australian & New Zealand Journal of Statistics*, 56(3):201–215.
- 596 Jia, K., Zhou, S., Zhuang, J., and Jiang, C. (2014). Possibility of the inde-
597 pendence between the 2013 Lushan earthquake and the 2008 Wenchuan
598 earthquake on Longmen Shan Fault, Sichuan, China. *Seismological Re-*
599 *search Letters*, 85(1):60–67.
- 600 Johnson, N. L. (1960). An approximation to the multinomial distribution
601 some properties and applications. *Biometrika*, 47(1-2):93–102.

- 602 Johnson, N. L. and Young, D. H. (1960). Some applications of two approxi-
603 mations to the multinomial distribution. *Biometrika*, pages 463–469.
- 604 Karr, A. (1991). *Point Processes and Their Statistical Inference*. Marcel
605 Dekker, Inc., New York and Basel.
- 606 Kiyosugi, K., Connor, C. B., Sparks, R. S. J., Crossweller, H. S., Brown,
607 S. K., Siebert, L., Wang, T., and Takarada, S. (2015). How many explosive
608 eruptions are missing from the geologic record? analysis of the quaternary
609 record of large magnitude explosive eruptions in japan. *Journal of Applied*
610 *Volcanology*, 4:17.
- 611 Machida, H. and Arai, F. (2003). *Atlas of Tephra in and around Japan*,
612 *revised edition*. University of Tokyo Press, Japan (in Japanese).
- 613 Marsan, D. and Enescu, B. (2012). Modeling the foreshock sequence prior
614 to the 2011, MW9.0 Tohoku, Japan, earthquake. *Journal of Geophysical*
615 *Research: Solid Earth*, 117(B6):B06316.
- 616 Mohler, G. O., Short, M. B., Brantingham, P. J., Schoenberg, F. P., and
617 Tita, G. E. (2011). Self-exciting point process modeling of crime. *Journal*
618 *of the American Statistical Association*, 106(493):100–108.
- 619 Møller, J. and Waagepetersen, R. P. (2003). *Statistical Inference and Simu-*
620 *lation for Spatial Point Processes*. Chapman and Hall/CRC.
- 621 Ogata, Y. (1988). Statistical models for earthquake occurrences and residual
622 analysis for point processes. *Journal of the American Statistical Associa-*
623 *tion*, 83(401):9–27.

- 624 Ogata, Y. (2006). Monitoring of anomaly in the aftershock sequence of the
625 2005 earthquake of M7.0 off coast of the western Fukuoka, Japan, by the
626 ETAS model. *Geophysical Research Letters*, 33:L01303.
- 627 Ogata, Y. and Katsura, K. (1993). Analysis of temporal and spatial het-
628 erogeneity of magnitude frequency distribution inferred from earthquake
629 catalogues. *Geophysical Journal International*, 113(3):727–738.
- 630 Ogata, Y., Katsura, K., Falcone, G., Nanjo, K., and Zhuang, J. (2013).
631 Comprehensive and topical evaluations of earthquake forecasts in terms of
632 number, time, space, and magnitude. *Bulletin of the Seismological Society
633 of America*, 103(3):1692–1708.
- 634 Ogata, Y. and Vere-Jones, D. (2003). Examples of statistical models and
635 methods applied to seismology and related earth physics. In Lee, W. H.,
636 Kanamori, H., Jennings, P. C., and Kisslinger, C., editors, *International
637 Handbook of Earthquake and Engineering Seismology, Vol.81B*, chapter 82.
638 International Association of Seismology and Physics of Earth’s Interior.
- 639 Ogata, Y. and Zhuang, J. (2006). Space-time ETAS models and an improved
640 extension. *Tectonophysics*, 413(1-2):13–23.
- 641 Omi, T., Ogata, Y., Hirata, Y., and Aihara, K. (2013). Forecasting large
642 aftershocks within one day after the main shock. *Scientific Reports*, 3:2218.
- 643 Omi, T., Ogata, Y., Hirata, Y., and Aihara, K. (2014). Estimating the ETAS
644 model from an early aftershock sequence. *Geophysical Research Letters*,
645 41(3):850–857.

- 646 Omi, T., Ogata, Y., Hirata, Y., and Aihara, K. (2015). Intermediate-term
647 forecasting of aftershocks from an early aftershock sequence: Bayesian and
648 ensemble forecasting approaches. *Journal of Geophysical Research: Solid*
649 *Earth*, 120(4):2561–2578.
- 650 Passarelli, L., Sandri, L., Bonazzi, A., and Marzocchi, W. (2010). Bayesian
651 hierarchical time predictable model for eruption occurrence: an application
652 to kilauea volcano. *Geophysical Journal International*, 181(3):1525–1538.
- 653 Peng, Z., Vidale, J. E., Ishii, M., and Helmstetter, A. (2007). Seismic-
654 ity rate immediately before and after main shock rupture from high-
655 frequency waveforms in Japan. *Journal of Geophysical Research: Solid*
656 *Earth*, 112(B3):B03306.
- 657 Pyle, D. M. (2015). Sizes of volcanic eruptions. In Sigurdsson, H., editor,
658 *The Encyclopedia of Volcanoes (Second Edition)*, chapter 13, pages 257 –
659 264. Academic Press, Amsterdam, second edition edition.
- 660 Sawazaki, K. and Enescu, B. (2014). Imaging the high-frequency energy radi-
661 ation process of a main shock and its early aftershock sequence: The case of
662 the 2008 Iwate-Miyagi Nairiku earthquake, Japan. *Journal of Geophysical*
663 *Research: Solid Earth*, 119(6):4729–4746.
- 664 Schoenberg, F. P. (2003). Multidimensional residual analysis of point process
665 models for earthquake occurrences. *J. Am. Stat. Assoc.*, 98:789–795(7).
- 666 Schoenberg, F. P., Chang, C., Keeley, J., Pompa, J., Woods, J., and H., X.
667 (2007). A critical assessment of the burning index in Los Angeles County,
668 California. *International Journal of Wildland Fire*, 16:473–483.

- 669 Siebert, L. and Simkin, T. (2002). Volcanoes of the world: an illustrated
670 catalog of holocene volcanoes and their eruptions, smithsonian institution,
671 global volcanism program digital information series, gvp-3.
- 672 Utsu, T., Ogata, Y., and Matsu'ura, R. S. (1995). The centenary of the
673 Omori formula for a decay law of aftershock activity. *Journal of Physics*
674 *of the Earth*, 43(1):1–33.
- 675 Vere-Jones, D. (1970). Stochastic models for earthquake occurrence. *J. Roy.*
676 *Stat. Soc. Series B (Methodological)*, 32(1):1–62 (with discussion).
- 677 Wang, T. and Bebbington, M. (2012). Estimating the likelihood of an erup-
678 tion from a volcano with missing onsets in its record. *Journal of Volcanol-*
679 *ogy and Geothermal Research*, 243–244:14–23.
- 680 Wang, T. and Bebbington, M. (2013). Robust estimation for the weibull
681 process applied to eruption records. *Mathematical Geosciences*, 45(7):851–
682 872.
- 683 Werner, M. J., Helmstetter, A., Jackson, D. D., and Kagan, Y. Y. (2011).
684 High-resolution long-term and short-term earthquake forecasts for Califor-
685 nia. *Bulletin of the Seismological Society of America*, 101(4):1630–1648.
- 686 Yamamoto, T. (2015). Cumulative volume step-diagrams for eruptive mag-
687 mas from major quaternary volcanoes in japan. Technical Report GSJ
688 Open-File Report, No.613, Geological Survey of Japan, AIST.
- 689 Yukutake, Y., Tanada, T., Honda, R., Harada, M., Ito, H., and Yoshida,
690 A. (2010). Fine fracture structures in the geothermal region of Hakone
691 volcano, revealed by well-resolved earthquake hypocenters and focal mech-
692 anisms. *Tectonophysics*, 489:104–118.

693 Zhuang, J. (2011). Next-day earthquake forecasts by using the ETAS model.
694 *Earth, Planet, and Space*, 63:207–216.

695 Zhuang, J., Ogata, Y., and Vere-Jones, D. (2002). Stochastic declustering
696 of space-time earthquake occurrences. *Journal of the American Statistical*
697 *Association*, 97(3):369–380.

698 Zhuang, J., Ogata, Y., and Vere-Jones, D. (2004). Analyzing earthquake clus-
699 tering features by using stochastic reconstruction. *Journal of Geophysical*
700 *Research*, 109(3):B05301.

701 Zhuang, J., Ogata, Y., and Wang, T. (2017). Data completeness of the Ku-
702 mamoto earthquake sequence in the JMA catalog and its influence on the
703 estimation of the ETAS parameters. *Earth, Planets and Space*, 69(1):36.

704 Zipkin, J. R., Schoenberg, F. P., Coronges, K., and Bertozzi, A. L. (2015).
705 Point-process models of social network interactions: Parameter estimation
706 and missing data recovery. *European Journal of Applied Mathematics*,
707 FirstView:1–28.