

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
Notice: Accepted version subject to English editing.	

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 January 16, 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 have
8 missing data, resulting in underestimation of the corresponding hazards.
9 We propose a computational approach for replenishing missing data in the
10 records of temporal point processes with time-separable marks. The basis of
11 this method is that, if such a point process is completely observed, it can
12 be transformed into a homogeneous Poisson process approximately on the
13 unit square $[0, 1]^2$ by a biscale empirical probability integral 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; (3) Generating events in the missing region. We

19 test this method on a synthetic dataset, and apply it to the records of vol-
20 canic eruptions of the Hakone Volcano in Japan and the aftershock sequence
21 following the 2008 Wenchuan Mw7.9 earthquake in Southwest China. The
22 results show that this algorithm provides a useful way to estimate missing
23 data and to replenish incomplete records of marked point processes. The
24 replenished data provide more robust estimates of the hazard function.

25 **1 Introduction**

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

43 Geophysicists have been searching for reliable methods to obtain more
44 complete earthquake catalogs. For example, waveform-based detection meth-
45 ods for small earthquakes within an aftershock sequence have been proposed
46 (e.g., Enescu et al., 2007, 2009; Peng et al., 2007; Marsan and Enescu, 2012;
47 Hainzl, 2016). However, even these methods cannot recover all missing after-
48 shocks. An alternative is to switch to energy-based descriptions (Sawazaki
49 and Enescu, 2014); that is, instead of regarding it as a process of events with
50 different magnitudes, the process of earthquake occurrences is regarded as
51 a stream of energies released by earthquakes. However, methods related to
52 such descriptions remain underdeveloped.

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

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

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

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

86 **2 Concepts, methodology, and illustration**

87 **2.1 Mark-separable temporal point process and biscale** 88 **empirical probability integral transformation**

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

98 A marked temporal point process is often specified by its conditional
99 intensity function, which is 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, i.e.,

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

109 the marks of this point process is said time-separable. Point-process models
110 with time-separable marks have been 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 these ways, it refers to a spatial point process.
121 A point process can be completely determined by its likelihood (termino-
122 logically, the local Janossy density, see Daley and Vere-Jones, 2003, 2008).
123 This gives the joint probability density/mass function of the total number
124 and each location of the particles in the process, assuming that the particles
125 are indistinguishable. If one of the following three is known: (1) the moment
126 intensities of all orders, (2) the conditional intensity, and (3) the Papangelou
127 intensity, then the likelihood is also known (i.e., the point process is com-
128 pletely determined). 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, as we see in the fol-
131 lowing sections, the method for replenishing missing data in a marked point
132 process does not depend on any specific form of the conditional intensity.
133 Therefore, it can be applied to spatial point processes as well if the ground
134 space is one-dimensional and the 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 i.i.d. observations on X , x_1, x_2, \dots, x_n , for a fixed x , the empirical 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, $\tilde{F}_X(X_j)$, $j = 1, 2, \dots, n$, 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. But the empirical probability integral transformation still results in
 141 an approximate unit uniform distribution since the transformation does not
 142 require the 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 space
 145 of marks. Consider the following biscale empirical transformation (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 are
 147 $i = 1, 2, \dots, n\}$, respectively. If the marks of the events in the process are
 148 separable from the occurrence times, then $\{t'_i : i = 1, 2, \dots, n\}$ and $\{m'_i : i = 1, 2, \dots, n\}$, which are the images of $\{t_i : i = 1, 2, \dots, n\}$ and $\{m_i : i = 1, 2, \dots, n\}$, respectively, approximately form a homogeneous Poisson process
 150 on $[0, 1] \times [0, 1]$. It is straightforward to show the independence between $\tilde{F}(t)$
 151 and $\tilde{G}(m)$ and, thus, given the total number of events N , the number of
 152 events in a cell of area $s \subseteq [0, 1] \times [0, 1]$ is a random variable from a binomial
 153 distribution $B(N, s)$, which can be approximated by a Poisson distribution
 154 with mean Ns . The smaller s gets, the better this approximation.
 155

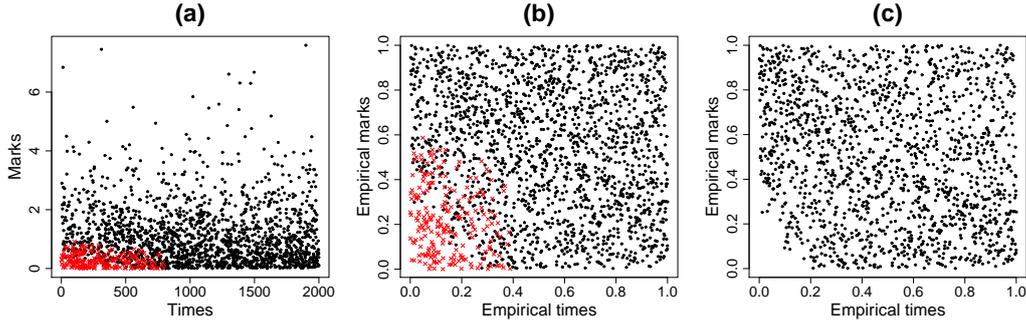


Figure 1: A synthetic dataset of a marked point process. (a) Marks versus occurrence times. (b) Empirical marks versus empirical occurrence times of all the synthetic events under the transformation Γ_N . (c) Empirical marks versus empirical occurrence times for the observed incomplete record under the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

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

Example 1. In Figure 1(a), we simulate a Poisson process N (the combi-
 nation of black and red points) with a temporal rate $\lambda = 1$ on $[0, 2000]$, and
 marks following an exponential distribution with mean 1, i.e.,

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

164 Figure 1(b) shows that under transformation (4), N is transformed into an
 165 approximately homogeneous Poisson process, say N' , which has rate $\lambda = 2000$
 166 and i.i.d. marks uniformly distributed in $[0, 1]$.

167 2.2 Detection of missing data

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

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

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

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

180 In the above biscale transformation, we do not need to know the exact
181 forms of $g(m)$, λ_g , or q . This method only uses the conditions that the
182 original process is mark-separable, and that the process of missing events
183 is time- and mark-dependent. Thus, for a temporal point process N with
184 time-separable marks, we can test whether there are data missing from its
185 observed record, N_{obs} , by testing the homogeneity of the image $\Gamma_{N_{\text{obs}}}(N_{\text{obs}})$ of
186 the observed data N_{obs} in the biscale transformed domain, when the missing
187 values are time- and mark-dependent. After using the BEPIT $\Gamma_{N_{\text{obs}}}$ to map
188 N_{obs} onto $[0, 1]^2$, we divide the overall area of $[0, 1]^2$ into L sub-regions of

189 equal areas, $L = L_1 \times L_2$ cells. Here, L_1 is the number of cells along the
190 transformed time domain and L_2 is the number of cells along the transformed
191 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)

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

197 Suppose that $[0, 1]^2$ is divided into $L = L_1 \times L_2$ cells with equal ar-
198 eas, i.e., $[0, 1]^2 = \bigcup_{j=1}^{L_2} \bigcup_{i=1}^{L_1} [(i-1)/L_1, i/L_1] \times [(j-1)/L_2, j/L_2]$, L_1 and
199 L_2 being positive integers. For any point process N on $[0, 1]^2$, if N is a
200 homogenous Poisson process, then the numbers of events in the above L
201 cells, C_1, C_2, \dots, C_L , form a homogeneous (n, \mathbf{p}) -multinomial random vec-
202 tor, with $\mathbf{p} = (1/L, 1/L, \dots, 1/L)$. However, if N is obtained by applying
203 the BEPIT to a completely observed mark-separable point process, then the
204 row sum of C_i in the k th row ($1 \leq k \leq L_1$), and the column sum of C_i
205 in 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
206 $\lfloor jn/L_2 \rfloor - \lfloor (j-1)n/L_2 \rfloor$, respectively, where $\lfloor x \rfloor$ denotes the integer part
207 of x , and n is the total number of events in N . Such constraints do not
208 hold for the homogeneous multinomial distribution. Since the distributions
209 of R and D are complicated, we obtain them by simulation: (1) with n fixed,
210 simulating n events uniformly distributed in $[0, 1]^2$; (2) applying the BEPIT
211 to these n simulated events; (3) with the specified parameters, L_1 and L_2 ,
212 calculating R and/or D for the transformed points.

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

230 2.3 Imputation method and algorithm

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

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

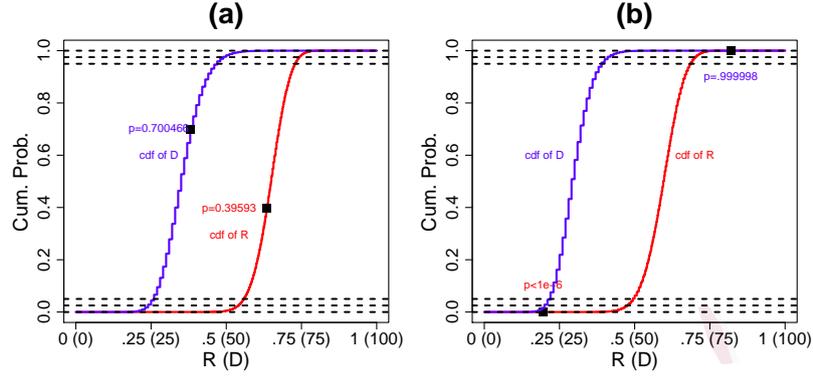


Figure 2: Statistical tests of the existence of missing data on (a) all the events and (b) the observed events in the synthetic point process, with cumulative distribution functions of R (red curve) and D (blue curve). R and D are defined in (6), with $L = L_1 \times L_2$, $L_1 = L_2 = 5$. The cumulative distribution functions 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.

236 and

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

237 where

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

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

$$w_{x,i} = \frac{1}{1 - \int_0^1 I((x_i, y) \in S) d\tilde{F}_y(y)}, \quad w_{y,i} = \frac{1}{1 - \int_0^1 I((x, y_i) \in S) d\tilde{F}_x(x)}. \quad (10)$$

240 Since F_Y and F_X are unknown, we replace them by \tilde{F}_Y and \tilde{F}_X , respectively,

241 i.e.,

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

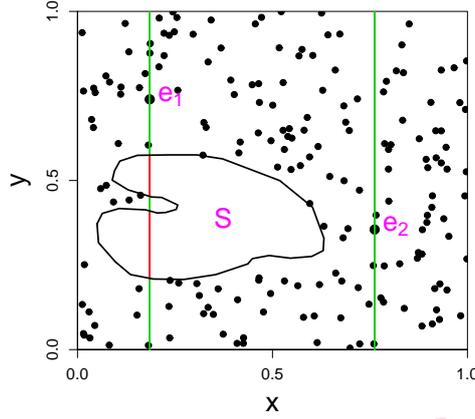


Figure 3: A heuristic estimation of the empirical distribution with missing points. Suppose that, among events $e_i = (x_i, y_i)$, $i = 1, 2, \dots, N$, events that fall in S cannot be observed. To estimate the empirical distribution $\tilde{F}_X(x)$ of x_i , $i = 1, 2, \dots, N$, weights need to be 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 $w_{x,2} = 1$ since the vertical line segment crossing e_2 has no intersection with S .

242 The above equation, together with (7) and (8), form a solvable equation
 243 system. We introduce below an algorithm to solve this equation system.

244 Firstly, the missing region S needs to satisfy the following condition:

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

248 This requirement is to ensure that the empirical distributions of $\{t_i\}$
 249 and $\{m_i\}$ can be restored. With Condition 1 satisfied, when a record is
 250 incomplete, we can determine the area, say S , outside of which the record is
 251 complete. This can be done either in the original mark-time plot based on
 252 prior knowledge of the data quality or in the BEPIT domain based on the
 253 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 the missing events; (2) estimating a new empirical
 257 distribution 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 dataset $N_{\text{obs}} = \{(t_i, m_i) : i = 1, 2, \dots, n\}$
 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 by

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

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

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

275 by $N_{\text{obs}}^* = \{(t_i^*, m_i^*) : i = 1, 2, \dots, n\}$ 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 by 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) by N_{rep} .

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

288 **Example 4.** *The above algorithm is applied to the thinned dataset in Ex-*
289 *ample 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). Tests us-*
291 *ing 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 dataset*
293 *is complete (Figure 4(e)). Figure 4(f) compares the cumulative numbers of*
294 *events in the original, the observed, and the replenished processes, showing*
295 *that the replenishing algorithm recovers the missing data to some extent.*

296 **Notes:**

297 (1) Equation (13) is the BEPIT that we mentioned in the previous section. If
298 the data are completely recorded, $\{(t_i^{(1)}, m_i^{(1)}), i = 1, 2, \dots, n\}$ form an
299 approximately homogeneous process on $[0, 1]^2$. As we can see in Figure
300 2(b), the sparseness of 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 biscale 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 materials.

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, given the existing events observed in S , we should keep them
 326 and remove the same number of simulated points.

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

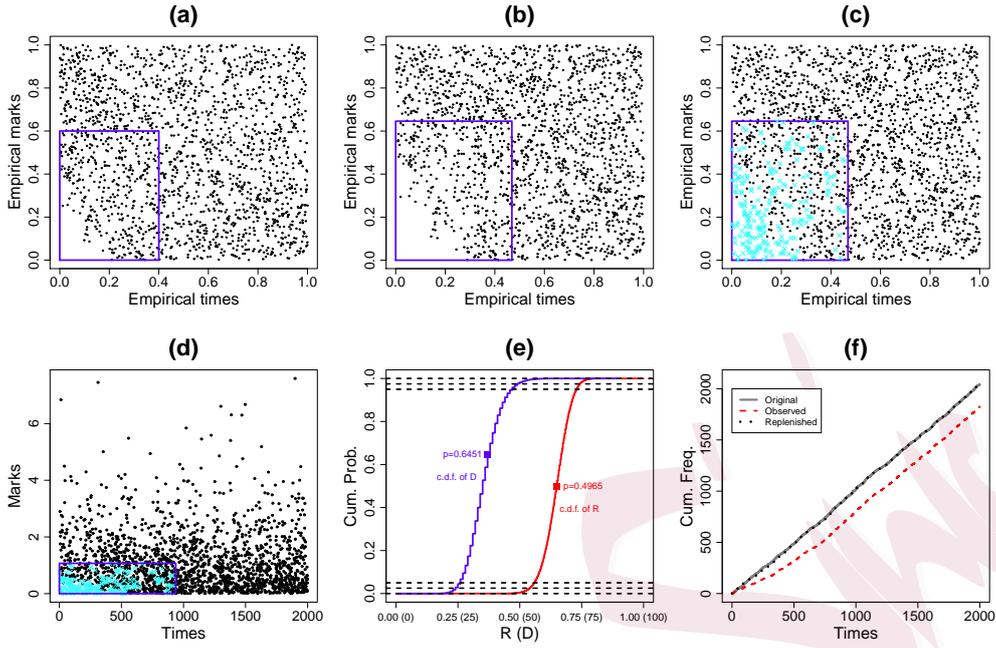


Figure 4: An application of the proposed replenishing algorithm to the synthetic dataset. (a) Rescaled marks versus rescaled occurrence times of the observed events (black dots), with the biscale transformation $\Gamma_{N_{\text{obs}}}$ based on the observed process. The blue polygon is the missing area, $S^{(1)}$. (b) Rescaled marks versus rescaled occurrence times of the observed events (black dots), with the rescaling $\Gamma_{N_{\text{obs}}}^*$ based on the events outside of S . The blue polygon is the missing area after transformation $\Gamma_{N_{\text{obs}}}^*$, i.e., S^* . (c) Rescaled marks versus rescaled occurrence times of the observed and replenished events (blue dots) (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 the empirical distributions of the events outside S . (d) Marks versus occurrence times of the observed synthetic events and the replenished events. (e) Cumulative distribution functions of R (red curve) and D (blue curve) for testing missing data in the replenished dataset in (c). (f) Cumulative frequencies versus occurrence times for the original, observed, and replenished processes.

332 2.4 More 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 $\lambda=2000$; (2) Simulating Poisson pro-
 336 cesses with rate λ drawn from a uniform distribution within $[100, 3000]$.

337 Both simulations have the same missing probability functions, as given by

338 (5). Figures 5(a) and (b) give the comparison between the true numbers of
339 missing events and the number of the replenished events for cases (1) and
340 (2), respectively. In Figure 5(a), since λ is fixed, the number of replenished
341 events is independent of the true number of missing events, and has a larger
342 variance. Some statistics related to these simulations, including the mean
343 numbers and variances of the missing and the replenished points, the mean
344 of relative differences, and the relative difference of means in 500 and 2000
345 simulations are given in Table 1. In particular, the near-zero relative deviation
346 of the mean number of the replenished events shows that the proposed
347 method is consistent. Here, the larger values of the mean relative deviation of
348 the number of replenished events from the number of missing events illustrate
349 the nature of the uncertainty related to the problem. Such uncertainty is pro-
350 duced not only by the randomness of the numbers of replenished and missing
351 events, but also by the uncertainty in the estimation of the occurrence rate
352 in the process from the events in the non-missing part. In Figure 5(b), the
353 expected number of replenished events in many repeated simulations is close
354 to the number of missing events. Moreover, the relative deviation decreases
355 when the number of missing events (or λ) increases. These results imply that
356 this algorithm replenishes the missing events reasonably well. Also, when λ
357 or the number of events in the process is quite small, there are some outputs
358 that the number of replenished events (when the number of missing events
359 is less than 50 in Figure 5(b)), which is simply calculated by the number of
360 simulated events in S in Steps 3 and 4 minus the number of observed events
361 in S , is negative. This indicates that the existence of missing data in these
362 situations cannot be quantified probabilistically.

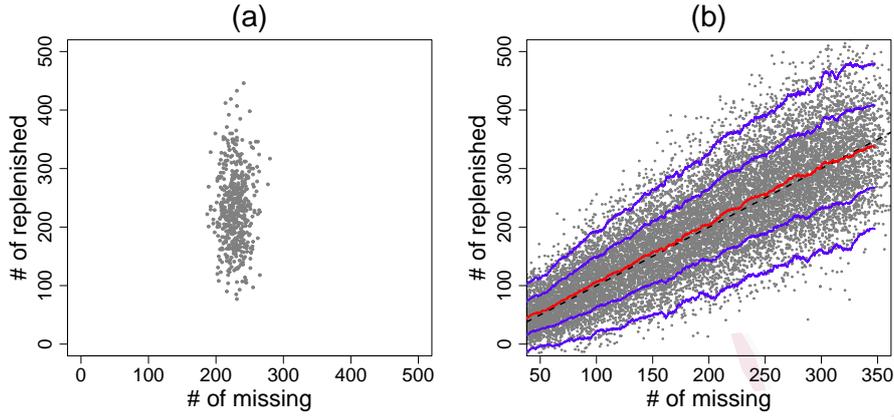


Figure 5: Comparison between the number of true missing events and the number of replenished events. (a) $\lambda = 2,000$ fixed. (b) λ is drawn from a uniform distribution between 100 and 3000. The dashed line represents the case where the numbers of missing and replenished events are equal. The blue and red curves represent the running mean and the corresponding single and double standard deviation bands.

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.	$\bar{\#m}$	$\sigma(\#m)$	$\bar{\#r}$	$\sigma(\#r)$	$\left[\frac{ \#m - \#r }{\#m} \right]$	$\frac{ \#m - \#r }{\#m}$
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

363 **3 Application**

364 **3.1 Volcanic eruption record**

365 In this example, we analyze the record of eruptions from the Hakone vol-
 366 cano. The Hakone volcano is an active volcano located at the northern
 367 boundary zone of the Izu-Mariana volcanic arc in central Japan (Yukutake
 368 et al., 2010; Honda et al., 2014). Data on Japanese explosive eruptions
 369 are compiled from the Smithsonian’s Global Volcanism Program database
 370 (Siebert and Simkin, 2002), the Large Magnitude Explosive Volcanic Erup-
 371 tions database (LaMEVE database, Croweller et al., 2012), [and additional](#)
 372 [Japanese databases](#) (Machida and Arai, 2003; Committee for Catalog of Qua-
 373 ternary Volcanoes in Japan (ed), 2000; Geological Survey of Japan, AIST

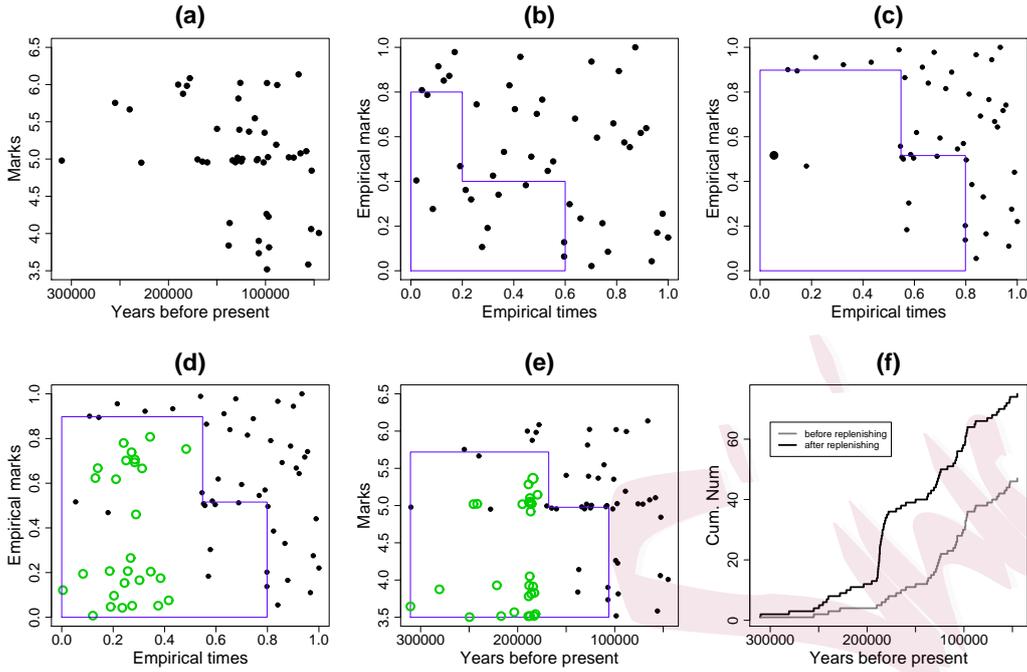


Figure 6: Results from applying the replenishment algorithm to volcanic eruption data. (a) Marks versus occurrence times of the eruption events. (b) Empirical distribution of marks versus that of occurrence times. (c) Rescaled marks versus rescaled occurrence times, with the rescaling based on the empirical 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 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 blue polygon is the area S and its corresponding mappings in which the missing events fall. The green dots are the replenished events.

374 (ed), 2013; Hayakawa, 2010).

375 For the Hakone volcano, 46 of 54 compiled events have an eruption mag-
 376 nitude ($M = \log_{10}[\text{erupted mass in kg}] - 7$; see Pyle (2015)) equal to or
 377 larger than 4 (Table S1 in the supplementary materials). 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. Based on this plot, the polygon boundaries of S are determined
 381 based on the following assumptions. First, events of empirical marks < 0.8

382 ($M < 5.7$) are missing before the empirical time = 0.2 (165 ka). Second,
383 the recording of larger events improves after the empirical time = 0.2 (165
384 ka), though events of empirical marks < 0.4 ($M < 5.0$) are still missing.
385 Third, the recording of events improves further, and there are no missing
386 events after the empirical time = 0.6 (105 ka). The results from running the
387 replenishing algorithm are shown in Figures 6(c) to 6(e).

388 The estimated cumulative number of events for the replenished dataset
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
392 ka (Figure 6(a); Hayakawa, 2010). The ages of the events at the Hakone
393 volcano are still not fully agreed in the literature. For example, Yamamoto
394 (2015) assumed that the ages of these eruptions are about 135 ka, 135 ka,
395 180 ka and 215 ka, respectively. Therefore, the reliability of the jump of
396 the cumulative number of events (Figure 6(f)) is a problem in volcanological
397 dating of event ages. In addition, estimating the tephra volume and rounded
398 eruption magnitude is also a problem in volcanology (Brown et al., 2014).
399 For example, the analyzed dataset has clusters of events with magnitude 4
400 and 5 (Figure 6(a)) and, therefore, the replenished events around 180 ka are
401 also 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 there are still a lot of eruptions missing in the overall Japanese database,
405 even for the last 100,000 years. Therefore, the polygon shape (Figure 6(b))
406 that we used suggests that our replenished data have the same completeness
407 level as the data outside the polygon. Our method is one possibility of

408 considering the under-recording of events in volcanic hazard assessments of
409 explosive 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 cannot be distinguished from the overlapping waveforms generated
414 by the mainshock on seismographs. In this section, we study the earthquake
415 catalog from Southwest China, between January 1, 1990, and April 20, 2013,
416 in a space range of $26^\circ - 34^\circ N$ and $97^\circ - 107^\circ E$ with minimum magnitude
417 3.0 (Figure S2 in supplementary materials). This dataset is selected from the
418 Chinese Earthquake catalog compiled by the China Earthquake Data Cen-
419 ter (CEDC) (URL: <http://data.earthquake.cn/index.html>). The Wenchuan
420 Mw 7.9 (Ms 8.0) earthquake, which occurred on May 12, 2008, was one
421 of the two largest seismic events in China during the last 50 years. There
422 are 6,249 events in the selected space and time range, among which 3,754
423 events occurred after the Wenchuan earthquake, indicating low seismicity
424 level above magnitude 3 in the study region before 2008. There are many
425 aftershocks missing immediately after the mainshock. In particular, events of
426 magnitudes between 3 and 4 are not properly recorded for a period of about
427 one-and-a-half months after the mainshock. The majority of the events after
428 May 12, 2008, can be taken as clustering events triggered by the Wenchuan
429 mainshock. When analyzing seismicity in this area, Jia et al. (2014) and Guo
430 et al. (2015) adopted a relatively high magnitude threshold of 4.0 to avoid
431 biases in estimates caused by missing events, with 5,217 of the 6,249 events
432 being ignored.

433 This example is quite different from the previous example and the sim-
 434 ulated data. The missing range can be well specified before replenishment:
 435 the missing values are known immediately after the occurrence of the main-
 436 shock, and the monitoring ability for events between magnitudes 3 and 4
 437 are restored one and half months later. The results are illustrated in Figure
 438 7. We can see that missing events take up about half the total number of
 439 events.

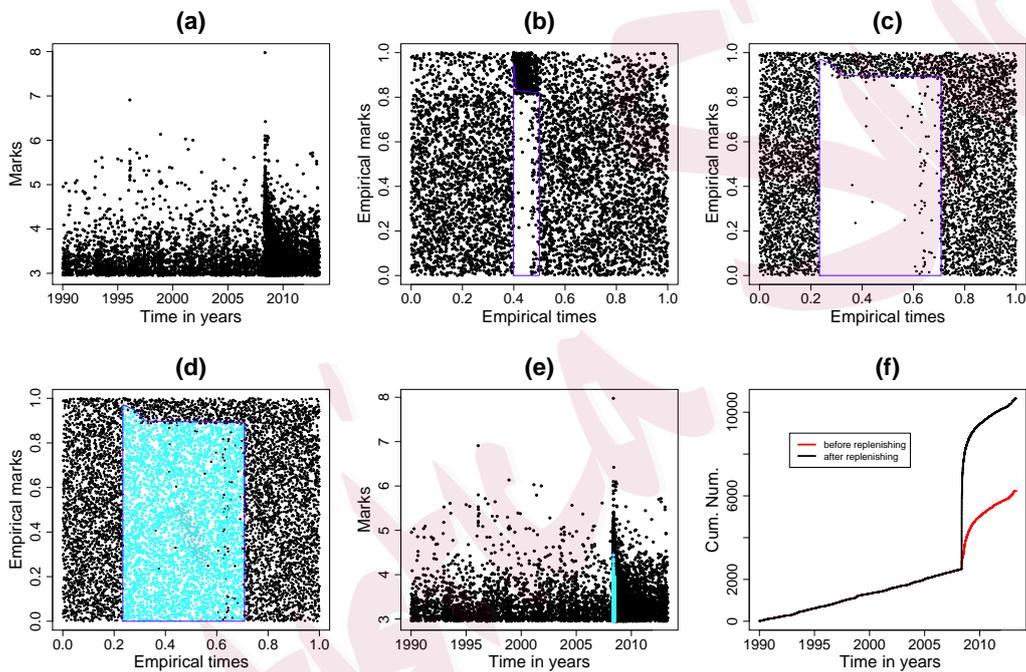


Figure 7: Results from applying the replenishment algorithm to the earthquake data from Southwest China. (a) Marks versus occurrence times of the earthquake events. (b) Empirical distribution of marks versus that of occurrence times. (c) Rescaled marks versus rescaled occurrence times, with the rescaling 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 blue polygon is the area S and its corresponding mappings in which the missing events fall. The blue dots are replenished 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
445 the aftershock rate drops slowly, and p is the power related to the decay
446 rate of aftershocks. Utsu et al. (1995) discussed how the parameters c and
447 p change with the cutoff magnitude threshold, and hypothesized that such
448 changes are caused by the fact that small aftershocks in an early stage of the
449 sequence are missing from the catalog. We fit the above Omori-Utsu formula
450 to both the original and the replenished catalogs (Table 2) and obtain the
451 maximum likelihood estimates of parameters. The results show that after
452 the replenishment, the Omori parameters c and p no longer change. We also
453 fit the Omori formula to the original dataset, but only consider earthquakes
454 that occurred at least 54 days after the mainshock. In this case, though 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 materials). These
458 results confirm numerically Utsu et al. (1995)'s hypothesis that missing small
459 events in the early stage of an aftershock sequence causes the instability of
460 the estimate of the Omori-Utsu formula.

461 **4 Conclusions and Discussions**

462 In this study, we proposed a method for replenishing missing data in marked
463 temporal point processes, based on only the assumption that the marks of the

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 datasets of earthquakes from Southwest China, with different magnitude thresholds. t_{main} : occurrence time of the mainshock; T : end of the time interval.

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 according
 467 to the parts where data are completely recorded. This method is applied
 468 to the eruption record of the Hakone volcano in Japan and the earthquake
 469 catalog from Southwest China, including the aftershock zone of the 2008
 470 Mw7.9 Wenchuan earthquake. The results show that the proposed method
 471 helps us evaluate the influence of missing data and correct the bias caused
 472 by missing data in our conclusion.

473 **Detection of the missing area** In our two examples, the missing area
 474 is determined by visual inspection of the biscale 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, such
477 missing area needs to be determined by the experience of data analysts or
478 information on the data from other sources. However, it is possible to turn
479 the replenishing algorithm into an automated algorithm.

480 Starting from $S' = \emptyset$, we divide the unit square into small cells in the
481 biscale 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
484 the test shows that missing cells exist, then we merge these cells into S' .
485 Such steps are iterated until no more cells are added to S' . Since this topic
486 belongs to the scope of data processing algorithms, we did not include it in
487 this statistical article.

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 the
494 marks of replenished events from $[0, 1]$ to the original mark space. If the mark
495 is dependent on the time, but we do not know how it depends on the time,
496 together with the existence of missing events, the replenishment/imputation
497 problem becomes unidentifiable.

498 Another case that is worth discussing is when the mark distribution is
499 known and does not depend on time. We can again use the cdf of the marks
500 in 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.

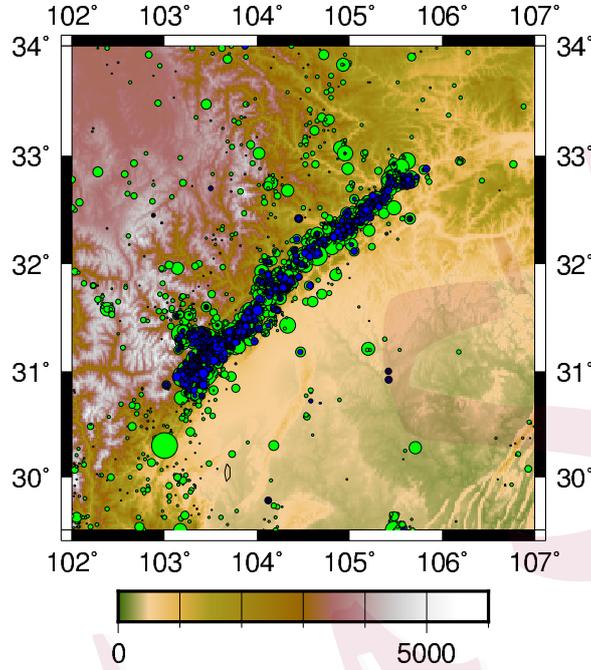


Figure 8: Epicenter map of imputed earthquakes (solid blue circles) for the Wenchuan aftershock sequence.

503 **Imputation of locations** This method is powerful for marked tempo-
504 ral point processes, but it cannot be extended easily to high-dimension or
505 spatiotemporal cases. This is because in most cases, the process is not ho-
506 mogeneous in space. However, it is still possible case by case. For example,
507 in replenishing the Wenchuan aftershock sequence, we can use the clustering
508 feature of earthquakes. A simple replenishing algorithm is as follows. For
509 each simulated event, find a fixed number, for example, 50, of events closest
510 to it in time in the observed process. Then we construct a Delaunay tessella-
511 tion network for these 50 events and select with equal probabilities one of the
512 Delaunay triangles, and put this simulated event randomly and uniformly in
513 this selected triangle. An example of the imputed locations of the missing

514 aftershocks of the Wenchuan earthquake is shown in Figure 8. For a spatially
515 inhibitive process, different methods should be used.

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

519 **Acknowledgement**

520 This project is supported by the Royal Society of New Zealand Marsden
521 Fund (contact UOO1419). JZ is also partially supported by Grants-in-Aid
522 No. 2530052 for Scientific Research (C) from the Japan Society for the Pro-
523 motion of Science. Helpful discussions with David Harte from GNS New
524 Zealand and Boris Baeumer from Otago University are gratefully acknowl-
525 edged. The authors also thank the AE and three anonymous reviewers for
526 their encouragement and constructive comments.

527 **References**

- 528 Bebbington, M. S. (2014). Long-term forecasting of volcanic explosivity.
529 *Geophysical Journal International*, 197:1500–1515.
- 530 Brown, S. K., Croweller, H. S., Sparks, R. S. J., Cottrell, E., Deligne, N. I.,
531 Guerrero, N. O., Hobbs, L., Kiyosugi, K., Loughlin, S. C., Siebert, L., and
532 Takarada, S. (2014). Characterisation of the quaternary eruption record:
533 analysis of the large magnitude explosive volcanic eruptions (LaMEVE)
534 database. *J Appl Volcanol*, 3:5.
- 535 Committee for Catalog of Quaternary Volcanoes in Japan (ed) (2000). Cat-
536 alog of quaternary volcanoes in Japan (in Japanese).

- 537 Corrado, C. J. (2011). The exact distribution of the maximum, minimum
538 and the range of multinomial/dirichlet and multivariate hypergeometric
539 frequencies. *Statistics and Computing*, 21(3):349–359.
- 540 Crossweller, H., Arora, B., Brown, S. K., Cottrell, E., Deligne, N., Guerrero,
541 N., Hobbs, L., Kiyosugi, K., C., L. S., Lowndes, J., Nayembil, M., Siebert,
542 L., Sparks, R. S. J., Takarada, S., and Venzke, E. (2012). Global database
543 on large magnitude explosive volcanic eruptions (LaMEVE). *Journal of*
544 *Applied Volcanology*, 1:4.
- 545 Daley, D. D. and Vere-Jones, D. (2003). *An Introduction to Theory of Point*
546 *Processes – Volume 1: Elementary Theory and Methods (2nd Edition)*.
547 Springer, New York, NY.
- 548 Daley, D. D. and Vere-Jones, D. (2008). *An Introduction to Theory of*
549 *Point Processes – Volume II: General Theory and Structure (2nd Edition)*.
550 Springer, New York, NY.
- 551 DeGroot, M. H. (1986). *Probability and Statistics (Second ed.)*. Addison-
552 Wesley.
- 553 Diggle, P. J. and Rowlingson, B. S. (1994). A conditional approach to point
554 process modelling of elevated risk. *Journal of the Royal Statistical Society.*
555 *Series A (Statistics in Society)*, 157(3):433–440.
- 556 Enescu, B., Mori, J., and Miyazawa, M. (2007). Quantifying early aftershock
557 activity of the 2004 mid-Niigata prefecture earthquake ($M_w6.6$). *Journal*
558 *of Geophysical Research: Solid Earth*, 112(B4):B004629.

- 559 Enescu, B., Mori, J., Miyazawa, M., and Kano, Y. (2009). Omori-Utsu law
560 c -values associated with recent moderate earthquakes in Japan. *Bulletin*
561 *of the Seismological Society of America*, 99(2A):884–891.
- 562 Geological Survey of Japan, AIST (ed) (2013). Catalog of eruptive events
563 during the last 10,000 years in japan, version 2.1 (in japanese). Technical
564 report.
- 565 Guo, Y., Zhuang, J., and Zhou, S. (2015). An improved space-time
566 ETAS model for inverting the rupture geometry from seismicity trigger-
567 ing. *Journal of Geophysical Research: Solid Earth*, 120(5):3309–3323.
568 2015JB011979.
- 569 Gutenberg, B. and Richter, C. F. (1944). Frequency of earthquakes in Cali-
570 fornia. *Bull. Seis. Soc. Am.*, 34:184–188.
- 571 Hainzl, S. (2016). Rate-dependent incompleteness of earthquake catalogs.
572 *Seismological Research Letters*, 87(2A):337–344.
- 573 Hayakawa, Y. (2010). Hayakawafs 2000-year eruption database and one
574 million-year tephra database, <http://www.hayakawayukio.jp/database/>.
- 575 Honda, R., Yukutake, Y., Yoshida, A., Harada, M., Miyaoka, K., and Sato-
576 mura, M. (2014). Stress-induced spatiotemporal variations in anisotropic
577 structures beneath hakone volcano, japan, detected by s wave splitting:
578 A tool for volcanic activity monitoring. *Journal of Geophysical Research:*
579 *Solid Earth*, 119(9):7043–7057.
- 580 Iwata, T. (2008). Low detection capability of global earthquakes after the
581 occurrence of large earthquakes: Investigation of the Harvard CMT cata-
582 logue. *Geophysical Journal International*, 174(3):849–856.

- 583 Iwata, T. (2013). Estimation of completeness magnitude considering daily
584 variation in earthquake detection capability. *Geophysical Journal Interna-*
585 *tional*, 194(3):1909–1919.
- 586 Iwata, T. (2014). Decomposition of seasonality and long-term trend in seis-
587 mological data: A Bayesian modelling of earthquake detection capability.
588 *Australian & New Zealand Journal of Statistics*, 56(3):201–215.
- 589 Jia, K., Zhou, S., Zhuang, J., and Jiang, C. (2014). Possibility of the inde-
590 pendence between the 2013 Lushan earthquake and the 2008 Wenchuan
591 earthquake on Longmen Shan Fault, Sichuan, China. *Seismological Re-*
592 *search Letters*, 85(1):60–67.
- 593 Johnson, N. L. (1960). An approximation to the multinomial distribution
594 some properties and applications. *Biometrika*, 47(1-2):93–102.
- 595 Johnson, N. L. and Young, D. H. (1960). Some applications of two approxi-
596 mations to the multinomial distribution. *Biometrika*, pages 463–469.
- 597 Karr, A. (1991). *Point Processes and Their Statistical Inference*. Marcel
598 Dekker, Inc., New York and Basel.
- 599 Kiyosugi, K., Connor, C. B., Sparks, R. S. J., Crossweller, H. S., Brown,
600 S. K., Siebert, L., Wang, T., and Takarada, S. (2015). How many explosive
601 eruptions are missing from the geologic record? analysis of the quaternary
602 record of large magnitude explosive eruptions in japan. *Journal of Applied*
603 *Volcanology*, 4:17.
- 604 Machida, H. and Arai, F. (2003). *Atlas of Tephra in and around Japan*,
605 *revised edition*. University of Tokyo Press, Japan (in Japanese).

- 606 Marsan, D. and Enescu, B. (2012). Modeling the foreshock sequence prior
607 to the 2011, MW9.0 Tohoku, Japan, earthquake. *Journal of Geophysical*
608 *Research: Solid Earth*, 117(B6):B06316.
- 609 Mohler, G. O., Short, M. B., Brantingham, P. J., Schoenberg, F. P., and
610 Tita, G. E. (2011). Self-exciting point process modeling of crime. *Journal*
611 *of the American Statistical Association*, 106(493):100–108.
- 612 Møller, J. and Waagepetersen, R. P. (2003). *Statistical Inference and Simu-*
613 *lation for Spatial Point Processes*. Chapman and Hall/CRC.
- 614 Ogata, Y. (1988). Statistical models for earthquake occurrences and residual
615 analysis for point processes. *Journal of the American Statistical Associa-*
616 *tion*, 83(401):9–27.
- 617 Ogata, Y. (2006). Monitoring of anomaly in the aftershock sequence of the
618 2005 earthquake of M7.0 off coast of the western Fukuoka, Japan, by the
619 ETAS model. *Geophysical Research Letters*, 33:L01303.
- 620 Ogata, Y. and Katsura, K. (1993). Analysis of temporal and spatial het-
621 erogeneity of magnitude frequency distribution inferred from earthquake
622 catalogues. *Geophysical Journal International*, 113(3):727–738.
- 623 Ogata, Y., Katsura, K., Falcone, G., Nanjo, K., and Zhuang, J. (2013).
624 Comprehensive and topical evaluations of earthquake forecasts in terms of
625 number, time, space, and magnitude. *Bulletin of the Seismological Society*
626 *of America*, 103(3):1692–1708.
- 627 Ogata, Y. and Vere-Jones, D. (2003). Examples of statistical models and
628 methods applied to seismology and related earth physics. In Lee, W. H.,

- 629 Kanamori, H., Jennings, P. C., and Kisslinger, C., editors, *International*
630 *Handbook of Earthquake and Engineering Seismology, Vol.81B*, chapter 82.
631 International Association of Seismology and Physics of Earth's Interior.
- 632 Ogata, Y. and Zhuang, J. (2006). Space-time ETAS models and an improved
633 extension. *Tectonophysics*, 413(1-2):13–23.
- 634 Omi, T., Ogata, Y., Hirata, Y., and Aihara, K. (2013). Forecasting large
635 aftershocks within one day after the main shock. *Scientific Reports*, 3:2218.
- 636 Omi, T., Ogata, Y., Hirata, Y., and Aihara, K. (2014). Estimating the ETAS
637 model from an early aftershock sequence. *Geophysical Research Letters*,
638 41(3):850–857.
- 639 Omi, T., Ogata, Y., Hirata, Y., and Aihara, K. (2015). Intermediate-term
640 forecasting of aftershocks from an early aftershock sequence: Bayesian and
641 ensemble forecasting approaches. *Journal of Geophysical Research: Solid*
642 *Earth*, 120(4):2561–2578.
- 643 Passarelli, L., Sandri, L., Bonazzi, A., and Marzocchi, W. (2010). Bayesian
644 hierarchical time predictable model for eruption occurrence: an application
645 to kilauea volcano. *Geophysical Journal International*, 181(3):1525–1538.
- 646 Peng, Z., Vidale, J. E., Ishii, M., and Helmstetter, A. (2007). Seismic-
647 ity rate immediately before and after main shock rupture from high-
648 frequency waveforms in Japan. *Journal of Geophysical Research: Solid*
649 *Earth*, 112(B3):B03306.
- 650 Pyle, D. M. (2015). Sizes of volcanic eruptions. In Sigurdsson, H., editor,
651 *The Encyclopedia of Volcanoes (Second Edition)*, chapter 13, pages 257 –
652 264. Academic Press, Amsterdam, second edition edition.

- 653 Sawazaki, K. and Enescu, B. (2014). Imaging the high-frequency energy radi-
654 ation process of a main shock and its early aftershock sequence: The case of
655 the 2008 Iwate-Miyagi Nairiku earthquake, Japan. *Journal of Geophysical*
656 *Research: Solid Earth*, 119(6):4729–4746.
- 657 Schoenberg, F. P. (2003). Multidimensional residual analysis of point process
658 models for earthquake occurrences. *J. Am. Stat. Assoc.*, 98:789–795(7).
- 659 Schoenberg, F. P., Chang, C., Keeley, J., Pompa, J., Woods, J., and H., X.
660 (2007). A critical assessment of the burning index in Los Angeles County,
661 California. *International Journal of Wildland Fire*, 16:473–483.
- 662 Siebert, L. and Simkin, T. (2002). *Volcanoes of the world: an illustrated*
663 *catalog of holocene volcanoes and their eruptions*, smithsonian institution,
664 global volcanism program digital information series, gvp-3.
- 665 Utsu, T., Ogata, Y., and Matsu'ura, R. S. (1995). The centenary of the
666 Omori formula for a decay law of aftershock activity. *Journal of Physics*
667 *of the Earth*, 43(1):1–33.
- 668 Vere-Jones, D. (1970). Stochastic models for earthquake occurrence. *J. Roy.*
669 *Stat. Soc. Series B (Methodological)*, 32(1):1–62 (with discussion).
- 670 Wang, T. and Bebbington, M. (2012). Estimating the likelihood of an erup-
671 tion from a volcano with missing onsets in its record. *Journal of Volcanol-*
672 *ogy and Geothermal Research*, 243–244:14–23.
- 673 Wang, T. and Bebbington, M. (2013). Robust estimation for the weibull
674 process applied to eruption records. *Mathematical Geosciences*, 45(7):851–
675 872.

- 676 Werner, M. J., Helmstetter, A., Jackson, D. D., and Kagan, Y. Y. (2011).
677 High-resolution long-term and short-term earthquake forecasts for Califor-
678 nia. *Bulletin of the Seismological Society of America*, 101(4):1630–1648.
- 679 Yamamoto, T. (2015). Cumulative volume step-diagrams for eruptive mag-
680 mas from major quaternary volcanoes in japan. Technical Report GSJ
681 Open-File Report, No.613, Geological Survey of Japan, AIST.
- 682 Yukutake, Y., Tanada, T., Honda, R., Harada, M., Ito, H., and Yoshida,
683 A. (2010). Fine fracture structures in the geothermal region of Hakone
684 volcano, revealed by well-resolved earthquake hypocenters and focal mech-
685 anisms. *Tectonophysics*, 489:104–118.
- 686 Zhuang, J. (2011). Next-day earthquake forecasts by using the ETAS model.
687 *Earth, Planet, and Space*, 63:207–216.
- 688 Zhuang, J., Ogata, Y., and Vere-Jones, D. (2002). Stochastic declustering
689 of space-time earthquake occurrences. *Journal of the American Statistical*
690 *Association*, 97(3):369–380.
- 691 Zhuang, J., Ogata, Y., and Vere-Jones, D. (2004). Analyzing earthquake clus-
692 tering features by using stochastic reconstruction. *Journal of Geophysical*
693 *Research*, 109(3):B05301.
- 694 Zhuang, J., Ogata, Y., and Wang, T. (2017). Data completeness of the Ku-
695 mamoto earthquake sequence in the JMA catalog and its influence on the
696 estimation of the ETAS parameters. *Earth, Planets and Space*, 69(1):36.
- 697 Zipkin, J. R., Schoenberg, F. P., Coronges, K., and Bertozzi, A. L. (2015).
698 Point-process models of social network interactions: Parameter estimation

699 and missing data recovery. *European Journal of Applied Mathematics*,
700 FirstView:1–28.

Statistica Sinica