# Probability distributions of COVID-19 tweet posted trends use a nonhomogeneous Poisson process

Devi Munandar[1,2,*], Sudradjat Supian[1], Subiyanto[3]

[1]*Department of Mathematics, Faculty of Mathematics and Natural Sciences, Universitas Padjadjaran, Indonesia*
[2]*Research Center for Informatics, Indonesian Institute of Sciences, Indonesia*
[3]*Department of Marine Sciences, Faculty of Fishery and Marine Sciences, Universitas Padjadjaran, Indonesia*

*Corresponding author e-mail address:devi19010@mail.unpad.ac.id*

**Abstract**

The influence of social media in disseminating information, especially during the COVID-19 pandemic, can be observed with time interval, so that the probability of number of tweets discussed by netizens on social media can be observed. The nonhomogeneous Poisson process (NHPP) is a Poisson process dependent on time parameters and the exponential distribution having unequal parameter values and, independently of each other. The probability of no occurrence an event in the initial state is one and the probability of an event in initial state is zero. Using of non-homogeneous Poisson in this paper aims to predict and count the number of tweet posts with the keyword coronavirus, COVID-19 with set time intervals every day. Posting of tweets from one time each day to the next do not affect each other and the number of tweets is not the same. The dataset used in this study is crawling of COVID-19 tweets three times a day with duration of 20 minutes each crawled for 13 days or 39 time intervals. The result of this study obtained predictions and calculated for the probability of the number of tweets for the tendency of netizens to post on the situation of the COVID-19 pandemic.

*Keywords:* Nonhomogeneous Poisson Process, Tweet, COVID-19, Estimated Parameter, Prediction

## 1. Introduction

The dissemination of information has a positive impact in the midst of a pandemic that occurs due to the spread of the Coronavirus in the world community (Severo et al., 2020). The dissemination of information through social media, especially Twitter, is very fast, especially conversations between netizens. Use of a non-homogeneous Poisson process (NHPP) to define the retweet hierarchy of the original tweet data (Gu and Kurov, 2020) set where all retweet processes are placed on a single hierarchical model so that information can be collect to estimate a parameter (Lee and Wilkinson, 2020). In determining the estimate for the application for non-life insurance data, NHPP is used as a method of solution to analyze data within a certain period of time (Vedyushenko, 2018). Even the nonhomogeneous Poisson process (NHPP) hierarchical model can also be used to disseminate information on online social

media, especially Twitter retweets. The retweet of each original tweet modeled by the NHPP, which has a function of intensity, is the product of a component that is built up by time and other components that are assigned to the number of followers of the original tweet's author.

In the Poisson process, it is provided the process of counting for each number of events at a certain time interval with λ parameters. The Poisson process of calculation does not depend on the previous interval process or depends on one another or is stationary and have relationship with exponential distribution process consisting of homogeneous Poisson and nonhomogeneous Poisson (Ross, 2014).

In its application, the Poisson process is very widely used in the field of statistics, as well as in calculations for prediction and implementation of other problems. In fact, other developments for nonhomogeneous are linked to the Poisson process and produce a nonhomogeneous compound Poisson process which is described mathematically to anticipate the number of occurrences (Grabski, 2019). In the process of calculating the process of marine ship accidents in the Baltic sea and ports, the nonhomogeneous Poisson process involved an important role in calculating and providing a model that can anticipate marine ship accidents (Franciszek, 2018). In medicine, to model seasonal events that occur due to dengue fever, using a nonhomogeneous Poisson process, which combines seasonal factors for the number of disease sufferers with analysis to improve the function of NHPP in daily cases (Cifuentes-Amado and Cepeda-Cuervo, 2015). In terms of climatology, modeling of extreme rainfall is also useful for studying the effects of seasonality and trends on modeling of extreme events of daily rainfall that exceed predetermined threshold values (Ngailo et al., 2016). In geostatistical modeling, the process of calculating the space-time approach data uses a nonhomogeneous Poisson process, involving two components: the Gaussian spatial component and the accounting component for its temporal effects, the objective is data suitability and identification of areas with the highest levels of pollution, namely the Southwest, Central and Northwest of Mexico City (Morales et al., 2017). In order to maximize the performance of NHPP, keep using the traditional Poisson base for software performance optimization by presenting the process in detail to prove that the resulting model is considered effective for improving and optimizing the performance of the traditional NHPP model and distribution function (Wang et al. 2016; Kim et al., 2010). In fact, in addition to testing software failures, from the hardware development side, NHPP is used to predict how long the development testing process will be terminated (Yu et al., 2007), So we also utilize of NHPP for the tweet dataset. In this study, we tried to model the set of Twitter crawling data with the keywords coronavirus and (Covid19 or COVID-19 as we call COVID-19), the process of calculating each keyword that appears in each tweet at defined time intervals. With the number of tweets obtained, the number of tweets can be modeled at a certain time using the nonhomogeneous Poisson process.

## 2. Materials and Methods

### 2.1. Materials

In this study, we used a dataset obtained from crawling tweet data about COVID-19, netizen posts, taken by filtering using predetermined keywords. We use coronavirus and COVID-19 as keywords and use of Twitter developer Application Programming Interface (API) which can be accessed to facilitate the crawling process. The data retrieval process is carried out in 3 time durations each day for 20 minutes for each time (See Table 1). This process is carried out for 13 days.

**Table 1.** Statistics of tweet datasets with coronavirus and COVID-19 keywords

| | Time to | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | ... | 39 |
| Number of tweet keywords | 8,994 | 31,514 | 33,237 | 9,232 | ... | 26,423 |
| Total of crawling data | 16,616 | 58,074 | 58,139 | 15,913 | ... | 58,139 |
| Ratio of tweet keywords | 0.541 | 0.543 | 0.572 | 0.580 | ... | 0.454 |

## 2.2. Methods

### 2.2.1. Counting Process

The stochastic process $\{N(t); t \geq 0\}$ is defined as counting process if $N(t)$ or $N_t$ states the number of events that occurred during time $t$ (Osaki, 1992). If it satisfies;

(i).   $N(t) > 0$

(ii).   $N(t)$ is integer number,

(iii).   If $s < t$ then $N(s) \leq N(t)$

(iv).   For $s < t, N(t) - N(s)$ denotes the number of events that occur at an interval time $(s, t]$.

The counting process is called a process with independent increments if the number of events that occurs in separate time intervals is mutually independent (Santitissadeekorn et al., 2020; Grabski., 2019). That is, the number of events that happened to time $t$, (i.e. $N(t)$), is independent of the number time events between t and $t + s$, (i.e. $N(t + s) - N(t)$.).

The counting process is named a process with stationary increments if the distribution of the number of events is occurs at certain time intervals only depending on the length of the interval, not depending on the location of the interval. It mean, the number of events in the time interval $(t_1 + s, t_2 + s]$ (i.e. $N(t_2 + s) - N(t_1 - s)$) has the same distribution as the number of events in the time intervals $(t_1, t_2)$ (i.e. $N(t_2) - N(t_1)$), for all $t_1 < t_2, s >$ (Kenney and Keeping, 1962).

### 2.2.2. Homogeneous Poisson Process

In the counting process with $\{N(t); t \geq 0\}$ is named Poisson process with a parameter rate $\lambda > 0$ as well $N(0) = 0$ and process has stationary independent increments as satisfies $P(N(h) = 1) = \lambda h + o(h)$ and $P(N(h) \geq 2) = o(h)$, for $t \geq 0$ as Poisson stationer process, then: $P(N(s + t) - N(s) = k)$ $= P(N(t) = k \mid N(0) = 0) = P_k(t)$ for any $s \geq 0, t \geq 0$, represents the probability that $k$ events occur at the interval $(0, t]$.

The Poisson process is a plain stochastic process and is highly used for modeling the time at which appearance put in a system (Mingola, 2013). A counting process $\{N(t); t \geq 0\}$ is named to be a Poisson process with rate (parameter) $\lambda$ if:

(i).   $N(0) = 0$

(ii).   Process has Independent Increment.

(iii).   The probability of there being $k$ events in the time interval $t$

$$P_k(t) = P(N(t+s) - N(s) = k) = \frac{(\lambda t)^k}{k!} e^{-\lambda t}, k = 0,1,2,... \quad (1)$$

$$\forall s,t > 0, \Rightarrow N(s+t) - N(s) \approx POI(\lambda t).$$

where:

$$E[N(t)] = \lambda t, \ Var[N(t)] = \lambda t. \quad (2)$$

$\lambda = \dfrac{E[N(t)]}{t} =$ rate or the average number of events that occur per time $t$.

## 2.2.3. Nonhomogeneous Poisson Process

In the counting process $\{N(t); t \geq 0\}$ is named nonhomogeneous Poisson process with intensity function $\lambda t \geq 0$ for $t \geq 0$, if;

(i). $P(N(0) = 0) = 1$,

(ii). The stochastic process with independent increments for the process $\{N(t); t \geq 0\}$.

(iii). $P(N(t+s) - N(t) = k = \dfrac{\left(\int_t^{t+s} \lambda x(dx)\right)^k}{k!} e^{-\int_t^{t+s} \lambda x(dx)};$ \quad (3)

Based on the above statement, it can be determined for one dimension of the nonhomogeneous Poisson process, is:

$$P(N(t) = k = \frac{\left(\int_0^t \lambda x(dx)\right)^k}{k!} e^{-\int_0^t \lambda x(dx)} \text{ with } k = 0,1,2,.... \quad (4)$$

Meanwhile, the probabilities of no occur an event in the initial state is one and the number of events occurring at an interval of time is independent of each other.

Suppose: $\quad \Lambda(t) = \int_0^t \lambda(x)dx$

Where: $P\{N(t+s) - N(t) = k\} = \dfrac{\left(\Lambda(t+s) - \Lambda(t)^k\right)}{k!} e^{-\left(\Lambda(t+s) - \Lambda(t)^k\right)}$

The distributions of nonhomogeneous Poisson processes have expectation and variance function:

$$\Lambda(t) = E[N(t)] = \int_0^t \lambda(x)dx \quad (5)$$

$$V(t) = V[N(t)] = \int_0^t \lambda(x)dx \quad (6)$$

And standard deviation is:

$$D(t) = \sqrt{V[N(t)]} = \sqrt{\int_0^t \lambda(x)dx}, t \geq 0 \quad (7)$$

The increment expected value for $N(s+t) - N(t)$ is

$$\Delta(t;s) = E(N(s+t) - N(t) = \int_t^{t+s} \lambda x(dx) \quad (8)$$

And appropriate to standard deviation is:

$$\sigma(t;s) = \sqrt{\int_{t}^{t+s} \lambda(x)dx}. \tag{9}$$

A nonhomogeneous Poisson process with $\lambda(t) = \lambda, t \geq 0$ for each $t \geq 0$, is a regular Poisson process. The increments of a nonhomogeneous Poisson process are independent, but not necessarily stationary.

### 2.2.4. Model Parameter Estimation

Tweet posts with the topic of coronavirus, COVID-19 within 13 days and 3 time intervals every day provide information on the number of tweets that are still being discussed by netizens. From Figure 1, it can be seen the number of tweets obtained every 20 minutes. Then it will be approximated to calculate the intensity of $\lambda(t)$ with an estimate of the simple linear regression function, namely $y = \alpha + \beta x$ with satisfied (Sumiati et al., 2019):

$$T(\alpha, \beta) = \sum_{i=1}^{n} \left( y_i - (\alpha + \beta x_i) \right)^2 \tag{10}$$

To increase this equation to find $\alpha$ and $\beta$ in quadratic expression with derive value minimum the objective function $T$ and denote $\hat{\alpha}$ and $\hat{\beta}$ (Kenney and Keeping, 1962).

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} \tag{11}$$

$$\hat{\beta} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \tag{12}$$

$$\bar{x} = \frac{1}{n}\sum_{i=1}^{n} x_i \, , \ \bar{y} = \frac{1}{n}\sum_{i=1}^{n} y_i$$

with
$\bar{x}$ and $\bar{y}$ as the average $x_i$ and $y_i$, respectively.
$\hat{\alpha}$ and $\bar{\beta}$ as approximate parameters.

## 3. Results and Discussion

The numbers of tweets obtained through crawling three times per day to retrieve post tweet data according to the desired keywords are coronavirus and COVID-19, and then the empirical intensity of tweets per crawling duration can be obtained as shown in Table 2.

**Table 2.** Empirical intensity of trending tweets coronavirus and COVID-19 keywords

| Dates | Interval | Center of interval | Number of Tweet | Intensity [tweet per minute] |
|---|---|---|---|---|
| 10 March 2020 11.00-11.20 | [0, 20) | 10 | 8,994 | 449.700 |
| 10 March 2020 16.00-16.20 | [20, 40) | 30 | 31,514 | 1,575.700 |
| 10 March 2020 21.00-21.20 | [40, 60) | 50 | 33,237 | 1,661.850 |
| 11 March 2020 11.00-11.20 | [60, 80) | 70 | 9,232 | 461.600 |
| 11 March 2020 16.00-16.20 | [80, 100) | 90 | 18,068 | 903.400 |
| 11 March 2020 21.00-21.20 | [100, 120) | 110 | 33,693 | 1,684.650 |
| 12 March 2020 11.00-11.20 | [120, 140) | 130 | 20,569 | 1,028.450 |
| 12 March 2020 16.00-16.20 | [140, 160) | 150 | 32,577 | 1,628.850 |
| 12 March 2020 21.00-21.20 | [160, 180) | 170 | 31,481 | 1,574.050 |
| 13 March 2020 11.00-11.20 | [180, 200) | 190 | 19,296 | 964.800 |
| 13 March 2020 16.00-16.20 | [200, 220) | 210 | 8,184 | 409.200 |
| 13 March 2020 21.00-21.20 | [220, 240) | 230 | 31,785 | 1,589.250 |
| 14 March 2020 11.00-11.20 | [240, 260) | 250 | 32,223 | 1,611.150 |
| 14 March 2020 16.00-16.20 | [260, 280) | 270 | 31,438 | 1,571.900 |
| 14 March 2020 21.00-21.20 | [280, 300) | 290 | 14,428 | 721.400 |
| 15 March 2020 11.00-11.20 | [300, 320) | 310 | 31,105 | 1,555.250 |
| 15 March 2020 16.00-16.20 | [320, 340) | 330 | 29,279 | 1,463.950 |
| 15 March 2020 21.00-21.20 | [340, 360) | 350 | 2,129 | 106.450 |
| 16 March 2020 11.00-11.20 | [360, 380) | 370 | 30,266 | 1,513.300 |
| 16 March 2020 16.00-16.20 | [380, 400) | 390 | 13,331 | 666.550 |
| 16 March 2020 21.00-21.20 | [400, 420) | 410 | 7,038 | 351.900 |
| 17 March 2020 11.00-11.20 | [420, 440) | 430 | 5,784 | 289.200 |
| 17 March 2020 16.00-16.20 | [440, 460) | 450 | 30,119 | 1,505.950 |
| 17 March 2020 21.00-21.20 | [460, 480) | 470 | 27,970 | 1,398.500 |
| 18 March 2020 11.00-11.20 | [480, 500) | 490 | 1,131 | 56.550 |
| 18 March 2020 16.00-16.20 | [500, 520) | 510 | 30,241 | 1,512.050 |
| 18 March 2020 21.00-21.20 | [520, 540) | 530 | 29,730 | 1,486.500 |
| 19 March 2020 11.00-11.20 | [540, 560) | 550 | 29,050 | 1,452.500 |
| 19 March 2020 16.00-16.20 | [560, 580) | 570 | 28,947 | 1,447.350 |
| 19 March 2020 21.00-21.20 | [580, 600) | 590 | 27,750 | 1,387.500 |
| 20 March 2020 11.00-11.20 | [600, 620) | 610 | 28,769 | 1,438.450 |
| 20 March 2020 16.00-16.20 | [620, 640) | 630 | 26,652 | 1,332.600 |
| 20 March 2020 21.00-21.20 | [640, 660) | 650 | 23,821 | 1,191.050 |
| 21 March 2020 11.00-11.20 | [660, 680) | 670 | 27,190 | 1,359.500 |
| 21 March 2020 16.00-16.20 | [680, 700) | 690 | 27,198 | 1,359.900 |
| 21 March 2020 21.00-21.20 | [700, 720) | 710 | 26,413 | 1,320.650 |
| 22 March 2020 11.00-11.20 | [720, 740) | 730 | 27,125 | 1,356.250 |
| 22 March 2020 16.00-16.20 | [740, 760) | 750 | 27,282 | 1,364.100 |
| 22 March 2020 21.00-21.20 | [760, 780) | 770 | 26,423 | 1,321.150 |

As shown in Table 2 with $n =$ number of data set, it can be computationally calculated by finding the values of the following variables

$$T_x = \bar{x} = 390, \quad T_y = \bar{y} = 1,181.362, \qquad\qquad T_{xx} = T_{xy} - (T_x * T_y) = 13,992.05128$$

$$T_{yy} = \left(\frac{1}{n}\sum_{i=1}^{n} x_i^2\right) - T_x^2 = 50,666.6667, \qquad\qquad T_{xy} = \frac{1}{n}\sum x_i y_i = 474,723.1$$

This calculation will be used to calculate the estimated regression coefficient:

$$\beta = \frac{T_{xx}}{T_{yy}} = 0.276158907, \quad \alpha = T_y - (\beta * T_x) = 1,073.659564777,$$

And then application of equations (11) and (12) to obtain the linear regression intensity of tweet dataset for Table 2

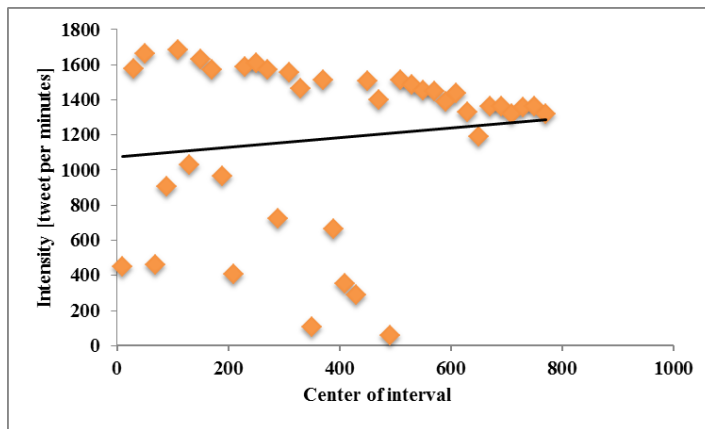$$\lambda(x) = 1,073.659564777 + 0.276158907x, x \geq 0 \tag{13}$$



**Figure 1.** Plot of center of interval with Intensity [tweet per minutes]

In Figure 1, it can be seen that to get a linear regression equation $(\lambda(x))$ as in equation (13), a plot is carried out against the center of interval data which is represented as an independent variable with the intensity of tweets per minutes represented by the dependent variable,

Using the equation (5), then obtained:

$$\Lambda(t) = \int_0^t (1,073.659564777 + 0.276158907x)dx$$
$$= 1,073.659564777t + 0.13807945t^2, t \geq 0 \tag{14}$$

Based on equations of (4), (5) and (14), we can obtain a one dimensional distribution of nonhomogeneous Poisson process for tweet dataset above is:

$$P(N(t) = k) = \frac{(\Lambda(t))^k}{k!} e^{-(\Lambda(t))}, k = 0,1,2,... \tag{15}$$

$$= \frac{(1,073.659564777t + 0.13807945t^2)^k}{k!} e^{-(1,073.659564777t + 0.13807945t^2)} \tag{16}$$

By utilizing $\Lambda(t) = (1,073.659564777t + 0.13807945t^2)$ where $t \geq 0$, we can anticipate number of tweet at certain interval $s$ with the time interval increase $N(t+s) - N(t)$ according to equations (8) and (9) namely the increment expected value and appropriate to standard deviation and can be calculated nonhomogeneous Poisson process with that time interval increase, like this equation:

$$P(N(t+s) - N(t) = k) = \frac{\left(\Lambda(t+s) - \Lambda(t)\right)^k}{k!} e^{-(\Lambda(t+s) - \Lambda(t))}. \tag{17}$$

**Example**

Suppose we will predict the number of coronavirus and COVID-19 tweets occurring 10 May 2020:21.00-21.20, So that we have a time interval [2920, 2940) and Probability that the number of coronavirus and COVID-19 tweets in considered interval of time is not greater than g=38,000 and not less that h=36,000.

Before calculate number of tweets that occurred on that date, it the first doing to determine the $t$ and $s$ parameter, the length of the interval $s = 1$ and $t = 2920$. With equations (8) and (9), we can predict that interval and standart deviation:

$$\Delta(2920,1) = E(N(2940) - N(2920)$$

$$\int_{2920}^{2940} (1,073.659564777 + 0.276158907 x) dx = 37,656.10284$$

and

$$\sigma(2920,1) = \sigma\left((N(2940) - N(2920)\right) = \sqrt{V(2920,1)} = 194.0518$$

The prediction results of the calculations above indicate the $147^{th}$ interval on that date. The predicted tweets discussing coronavirus and COVID-19 by netizens are 37,656 tweets and a standard deviation is 194.

Meanwhile, we can calculate of probability distribution the number of tweets posted by netizens following this:

$$P_{38,000 \leq k \leq 36,000} = P(38,000 \leq N(t+s) - N(t) \leq 36,000$$

$$= \sum_{k=36,000}^{k=38,000} \frac{37,656.10284^k}{k!} e^{-37,656.10284}$$

Using approximation to find out probability by normal distribution, we can calculate:

$$P_{36,000 \leq k \leq 38,000} = \Phi\left(\frac{38,000 - 37,656.10284}{194.0518}\right) - \Phi\left(\frac{36,000 - 37,656.10284}{194.0518}\right)$$

$$= \Phi(0.981819) - \Phi(7.04828E - 18) = 0.8315.$$

In other words the probability of the number of tweets in interval of time 36,000 and 38,000 at the interval [2920, 2940) is 83.15%.

## 4. Conclusion

Using the calculation process theory by utilizing the nonhomogeneous Poisson process concept, it is possible to build a stochastic model of the number of tweets with the keywords coronavirus and COVID-19 obtained by crawling the tweet data with a specified duration of time every day continuously. The counting process with an independent increment is a reasonable model for counting the number of tweets in a certain time period. The use of linear regression is one of is one of options to determine estimation parameters in approximation of number of tweets. With numbers of tweet data obtained is different for each defined time duration, namely three times taken in one day, and does not affect each other, so that the number of tweets can be marked as nonhomogeneous Poisson process. The application of the model obtained can be used to count the number of tweets at certain intervals, as well as the probability for prediction at time intervals outside the defined dataset.

## References

Cifuentes-Amado, M. V., and Cepeda-Cuervo, E. (2015). Non-homogeneous poisson process to model seasonal events: Application to the health diseases. *International Journal of Statistics in Medical Research*, *4*(4), 337-346.

Franciszek, G. (2018). Nonhomogeneous compound Poisson process application to modeling of random processes related to accidents in the Baltic Sea waters and ports. *Journal of Polish Safety and Reliability Association*, *9*(3), 21-29.

Grabski, F. (2018). Nonhomogeneous stochastic processes connected to Poisson process. *Scientific Journal of Polish Naval Academy*, *213*(2), 5-15.

Gu, C., and Kurov, A. (2018). Informational role of social media: evidence from twitter sentiment. *Journal of Banking & Finance*, *121*, 105969.

Kenney J. F. and Keeping, E. S. (1962). *Linear Regression and Correlation.* Ch. 15 in. Pt. 1, 3rd ed. Princeton, NJ: Van Nostrand.

Kim, D. K., Park, D. H., and Yeo, I. K. (2010). Mean functions based on meta-mixtures in nonhomogeneous Poisson processes. *Journal of the Korean Statistical Society*, *39*, 237-244.

Lee, C., and Wilkinson, D. J. (2020). A Hierarchical model of nonhomogeneous Poisson processes for twitter retweets. *Journal of the American Statistical Association*, *115*(529), 1-15.

Mingola, P. (2013). *A Study of Poisson and Related Processes with Applications*. University of Tennessee, Knoxville.

Morales, F. E. C., Vicini, L., Hotta, L. K., and Achcar, J. A. (2017). A nonhomogeneous Poisson process geostatistical model. *Stochastic Environmental Research and Risk Assessment*, *31*(2), 493-507.

Ngailo, T., Shaban, N., Reuder, J., Rutalebwa, E., and Mugume, I. (2016). Non homogeneous poisson process modelling of seasonal extreme rainfall events in Tanzania. *International journal of science and research (IJSR)*, *5*(10), 1858-1868.

Osaki, S. (1992). *Applied Stochastic System Modeling*. New York: Springer-Verlag.

Ross, S. (2014). *The Exponential Distribution and the Poisson Process*. In *Introduction to Probability Models (Eleventh Edition)*, ed. Sheldon Ross. Boston: Academic Press, 277–356.

Santitissadeekorn, N., Lloyd, D. J., Short, M. B., and Delahaies, S. (2020). Approximate filtering of conditional intensity process for Poisson count data: Application to urban crime. *Computational Statistics & Data Analysis*, *144*, 106850.

Severo, E. A., De Guimarães, J. C. F., and Dellarmelin, M. L. (2020). Impact of the COVID-19 pandemic on environmental awareness, sustainable consumption and social responsibility: Evidence from generations in Brazil and Portugal. *Journal of Cleaner Production*, 124947.

Sumiati, I., Rahmani, U., Supian, S., & Subiyanto, S. (2019). Application of the Nonhomogeneous Poisson Process for Counting Earthquakes. *World Scientific News*, *127*(3), 163-176.

Vedyushenko, B. A. (2018). *Non-Homogeneous Poisson Process - Estimation and Simulation*. Faculty Of Mathematics and Physics, Charles University, Prague.

Wang, J., Wu, Z., Shu, Y., and Zhang, Z. (2016). An optimized method for software reliability model based on nonhomogeneous Poisson process. *Applied Mathematical Modelling*, *40*(13-14), 6324-6339.

Yu, J. W., Tian, G. L., and Tang, M. L. (2007). Predictive analyses for nonhomogeneous Poisson processes with power law using Bayesian approach. *Computational Statistics & Data Analysis*, *51*(9), 4254-4268.