

Article

# Use of Hazard Functions for Determining Power-Law Behaviour in Data

Joseph D. Bailey 

School of Mathematics, Statistics and Actuarial Science, University of Essex, Colchester CO4 3SQ, UK;  
jbailef@essex.ac.uk

**Abstract:** Determining the ‘best-fitting’ distribution for data is an important problem in data analysis. Specifically, observing how the distribution of data changes as values below (or above) a threshold are omitted from analyses can be of use in various applications, from animal movement to the modelling of natural phenomena. Such truncated distributions, known as hazard functions, are widely studied and well understood in survival analysis, although rarely widely used in data analysis. Here, by considering the hazard and reverse-hazard functions, we demonstrate a qualitative assessment of the ‘best-fit’ distribution of data. Specifically, we highlight the potential advantages of this method when determining whether power-law behaviour may or may not be present in data. Finally, we demonstrate this approach using some real-world datasets.

**Keywords:** hazard function; truncated distribution; log-normal; power law

## 1. Introduction

Truncated datasets are widely encountered in data analysis, such as in mathematical biology, when determining between purposeful and non-purposeful movement [1–3]; in economic theory, when modelling risk of insurance data [4]; and in medical sciences, when considering infection and incubation time in infectious diseases [5] (see Ch.1 of [6] for more examples in various fields). In general, truncated distributions are well understood in how they relate to the initial distribution [7–9], with much work focused on methods for simulating such distributions [8,10,11] and understanding the properties of distributions after truncation [10,12–14].

A major aspect of data analysis involves inferring the distribution which most accurately describes a given dataset. In particular, there has been much debate regarding the presence of power-law (PL) behaviour in datasets [15–21]. Whilst recently many advances have been made in developing techniques for identifying PL behaviour [22–24], there are still occasions when it is difficult to establish whether data are best described by a PL. Specifically, determining between data which are drawn from a PL and that from a log-normal distribution is often intractable without a large number of data points [18,21,22,25,26].

The importance of being able to obtain the best descriptive distribution for modelling data is clear, as different but similar distributions can lead to large differences in predicted outcomes. For example, when considering random walk models in animal movements, the distributions describing the underlying characteristics of step lengths and turning angles have been shown to give observably different qualitative and quantitative results, which have large impacts in the use and predictive power of such models [27,28].

Here, we consider the hazard and reverse-hazard distribution function of a given distribution, which is found by continually truncating a distribution (from either the left or the



Academic Editor: Murilo da Silva  
Baptista

Received: 8 November 2024  
Revised: 22 December 2024  
Accepted: 7 January 2025  
Published: 9 January 2025

**Citation:** Bailey, J.D. Use of Hazard Functions for Determining Power-Law Behaviour in Data. *Analytics* 2025, 4, 2. <https://doi.org/10.3390/analytics4010002>

**Copyright:** © 2025 by the author. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

right) to produce a characteristic curve which is unique to the initial distribution. We show that, given this characteristic curve, the initial underlying distribution can be obtained. We derive expressions relating this characteristic curve to the initial distribution and demonstrate how this links the exponential, uniform, and PL distributions. Finally, we give examples of where calculating this characteristic curve is beneficial for data analysis and, in particular, highlight how this can help determine between a PL and log-normal distribution.

## 2. Background Information

### 2.1. Review of Hazard and Reverse-Hazard Functions

The hazard function is a well-understood and widely used property of distributions, which gives information regarding the probability of events above a certain threshold happening, given that the threshold value has been reached. As such, hazard functions have widespread applications in fields such as medical data [29] and insurance/risk management [30], among many others. The reverse-hazard function considers the probability of events happening below a given threshold, given that threshold has been reached [31], and has found applications in fields from astronomy to psychology [32].

Essentially, the hazard function considers data that have been truncated from the left, where values below a threshold have been removed from the dataset, whereas the reverse hazard considers right truncated data, where values above a threshold have been removed. Here, we briefly introduce the hazard and reverse-hazard functions, as well as an important property connecting these functions with the initial probability function.

**Definition 1.** Let  $f(x)$  be a continuous probability density function (PDF) on the interval  $[a, b]$  (note, we allow  $a = -\infty$  and  $b = \infty$ ), define  $f_k^{\{R\}}(x)$  as the distribution found by truncating  $f$  on the right at the value  $k$ , and  $f_k^{\{L\}}(x)$  as the distribution when  $f$  is truncated on the left at  $k$ , for all  $x, k \in (a, b)$ .

Using Definition 1, it can be easily seen that distributions  $f_k^{\{L\}}(x)$  and  $f_k^{\{R\}}(x)$  are given by

$$f_k^{\{L\}}(x) = \frac{f(x)}{\int_k^b f(x) dx} = \frac{f(x)}{F(b) - F(k)} = \frac{f(x)}{1 - F(k)}$$

$$f_k^{\{R\}}(x) = \frac{f(x)}{\int_a^k f(x) dx} = \frac{f(x)}{F(k) - F(a)} = \frac{f(x)}{F(k)}$$

where  $F(x)$  is the CDF of  $f(x)$  with  $F(b) = 1$  and  $F(a) = 0$ .

If we now consider how the value of the truncated distribution changes at the point of the truncation,  $k$ , we can define the hazard,  $h^{\{L\}}$ , and reverse-hazard functions,  $h^{\{R\}}$ .

**Definition 2.** Let  $f$  be a continuous PDF on  $[a, b]$ , define the hazard and reverse-hazard functions for all  $x \in (a, b)$ , respectively, as

$$h^{\{L\}}(x) = \frac{f(x)}{1 - F(x)} \tag{1}$$

$$h^{\{R\}}(x) = \frac{f(x)}{F(x)}. \tag{2}$$

Essentially,  $h^{\{L\}}, h^{\{R\}}$  describe the curve found by considering the value of the distribution,  $f$ , as it is continuously truncated from either the right or the left, at the point of truncation,  $x$ . Note, the reverse hazard,  $h^{\{R\}}$ , is often denoted as  $r(x)$ ; however, we use the  $h^{\{R\}}$  notation throughout to make it clear when we are truncating from the right.

Equations (1) and (2) show how the characteristic curves,  $h^{\{L\}}$  and  $h^{\{R\}}$ , can be found when the initial distribution  $f$  is known; however, a potentially beneficial question when attempting to discern ‘best-fit’ distributions from datasets would be “Given that we have the form of the curve found by truncating some distribution, can we determine the initial distribution?”

We now show how this can be carried out by deriving commonly used expressions for  $f(x)$  which depend only on  $h^{\{L\}}(x)$  and  $h^{\{R\}}(x)$ , showing that any function can be uniquely defined by its hazard or reverse-hazard function.

**Theorem 1.** *Let  $f(x)$  be a PDF defined on  $[a, b]$  and  $h^{\{L\}}, h^{\{R\}}$  be defined as in Definition 2; then, for all  $x \in (a, b)$ , we have*

$$f(x) = h^{\{L\}}(x) \exp\left(-\int_a^x h^{\{L\}}(\tilde{x}) d\tilde{x}\right), \tag{3}$$

$$f(x) = h^{\{R\}}(x) \exp\left(-\int_x^b h^{\{R\}}(\tilde{x}) d\tilde{x}\right). \tag{4}$$

**Proof.** We start by considering  $h^{\{L\}}$ , found by taking the left truncation. Writing  $f(x) = F'(x)$  in Equation (1) and using the anti-derivative, we obtain

$$h^{\{L\}}(x) = \frac{F'(x)}{1 - F(x)} = -\frac{d}{dx} \log |1 - F(x)|. \tag{5}$$

As  $F(x)$  is the CDF for  $f$  defined on the domain  $[a, b]$ , then  $0 < F(x) < 1$  for all  $x \in (a, b)$ ; hence, we can ignore the absolute sign in the logarithm. Now, we define the integrated hazard function (also known as the cumulative hazard function),  $H^{\{L\}}(x)$ , by  $\frac{d}{dx} H^{\{L\}}(x) = h^{\{L\}}(x)$ . Therefore, we have

$$H^{\{L\}}(x) = \int_a^x h^{\{L\}}(\tilde{x}) d\tilde{x},$$

noting that  $H^{\{L\}}(a) = \log(1 - F(a)) = \log(1 - 0) = 0$ . This implies that

$$h^{\{L\}}(x) = \frac{d}{dx} H^{\{L\}}(x) = \frac{d}{dx} \int_a^x h^{\{L\}}(\tilde{x}) d\tilde{x}.$$

Using this and Equation (5) gives

$$\begin{aligned} \frac{d}{dx} \int_a^x h^{\{L\}}(\tilde{x}) d\tilde{x} &= -\frac{d}{dx} \log(1 - F(x)) \\ \Rightarrow -\int_a^x h^{\{L\}}(\tilde{x}) d\tilde{x} &= \log(1 - F(x)) \\ \Rightarrow \exp\left(-\int_a^x h^{\{L\}}(\tilde{x}) d\tilde{x}\right) &= 1 - F(x) \end{aligned}$$

Finally, substituting this expression for  $1 - F(x)$  back into Equation (1) and rearranging gives

$$f(x) = h^{\{L\}}(x) \exp\left(-\int_a^x h^{\{L\}}(\tilde{x}) d\tilde{x}\right), \tag{6}$$

as required.

We can now apply similar reasoning for  $h^{\{R\}}$ , found by truncating the distribution from the right. Starting with Equation (2), this gives

$$h^{\{R\}}(x) = \frac{F'(x)}{F(x)} = \frac{d}{dx} \log |F(x)|. \tag{7}$$

Denoting  $H^{\{R\}}(x)$  as the integrated reverse-hazard function for  $h^{\{R\}}(x)$ , we have

$$h^{\{R\}}(x) = \frac{d}{dx} H^{\{R\}}(x) = -\frac{d}{dx} \int_x^b h^{\{R\}}(\tilde{x}) d\tilde{x},$$

noting that  $H^{\{R\}}(b) = \log(F(b)) = \log(1) = 0$ . Using this and Equation (7) gives

$$\begin{aligned} -\frac{d}{dx} \int_x^b h^{\{R\}}(\tilde{x}) d\tilde{x} &= \frac{d}{dx} \log(F(x)) \\ \Rightarrow -\int_x^b h^{\{R\}}(\tilde{x}) d\tilde{x} &= \log(F(x)) \\ \Rightarrow \exp\left(-\int_x^b h^{\{R\}}(\tilde{x}) d\tilde{x}\right) &= F(x) \end{aligned}$$

Finally, substituting this back into Equation (2) gives

$$f(x) = h^{\{R\}}(x) \exp\left(-\int_x^b h^{\{R\}}(\tilde{x}) d\tilde{x}\right), \tag{8}$$

as required.  $\square$

This demonstrates that whatever the given form of  $h^{\{L\}}(x)$  or  $h^{\{R\}}(x)$ , we can find the initial corresponding distribution by using Equations (3) and (4).

We now consider some examples and possible applications.

### 2.2. Examples

The following section demonstrates the connections between the exponential distribution, uniform distribution, and the simple PL through the truncation of data.

Let us consider that after right-truncating some unknown distribution,  $f(x)$ , a simple PL curve is found on the interval  $(a, b]$ ,  $h^{\{R\}}(x) = c(x - a)^{-1}$  where  $c > 0$ . Then, we can calculate the initial distribution using Equation (4)

$$\begin{aligned} f(x) &= \frac{c}{x - a} \exp\left(-\int_x^b \frac{c}{\tilde{x} - a} d\tilde{x}\right) \\ &= \frac{c}{x - a} \exp[c \log(x - a) - c \log(b - a)] \\ &= \frac{c}{x - a} \frac{(x - a)^c}{(b - a)^c} = \frac{c(x - a)^{c-1}}{(b - a)^c} \end{aligned}$$

considering the case for  $c = 1$  and  $b$  as a finite real value, then we have  $f(x) = \frac{1}{b-a}$  the uniform distribution on the domain  $[a, b]$ .

Now, let us consider the function,  $g(x)$ , which returns the uniform distribution when truncated from the left on the domain  $[a, b]$ . We have

$$\begin{aligned} g(x) &= \frac{1}{b - a} \exp\left(-\int_a^x \frac{1}{b - a} d\tilde{x}\right) \\ &= \frac{1}{b - a} \exp\left(\frac{a}{b - a} - \frac{x}{b - a}\right) \\ &= \frac{1}{b - a} \exp\left(\frac{1}{b - a}[a - x]\right). \end{aligned}$$

Letting  $\lambda = \frac{1}{b-a} \in \mathbb{R}$  a constant, we see  $g(x)$  is of the form  $g(x) = \lambda \exp(\lambda[a - x])$ , which is the exponential distribution with rate  $\lambda$  and minimum value  $a$ .

Hence, we have a path from the exponential distribution to a simple PL by first truncating the exponential distribution from the left to give the uniform distribution, then truncating the uniform distributions from the right to obtain the PL (Equation (9)).

$$\text{Exp} \xrightarrow{h^{L}} \text{Unif} \xrightarrow{h^{R}} (x - a)^{-1} \tag{9}$$

Table 1 shows some examples of the general form of the resulting distributions when truncated from the left (hazard) and from the right (reverse hazard), for selected initial distributions [31].

**Table 1.** Showing the form of the curves found by truncating various distributions from both the left and the right, where  $\phi$  and  $\Phi$  are the PDF and the CDF for the normal distribution, respectively. Note the PL case gives the general qualitative behaviour of the hazard and reverse-hazard curves, not the precise formulations.

Initial distribution, $f(x)$ .	$h^{L}$	$h^{R}$
Exponential, $\text{Exp}(x; \lambda)$	$= \lambda$	$= \frac{\lambda}{\exp(\lambda x) - 1}$
Power law, $\frac{\alpha}{x_{\min}} \left(\frac{x}{x_{\min}}\right)^{-(\alpha+1)}, \alpha > 1$	$\sim x^{-1}$	$\sim x^{-(\alpha+1)}$
Uniform, $U(x; a, b)$	$= (b - x)^{-1}$	$= (x - a)^{-1}$
log-normal, $LN(x; \mu, \sigma^2)$	$= \frac{1}{x\sigma} \frac{\phi\left(\frac{\log x - \mu}{\sigma}\right)}{1 - \Phi\left(\frac{\log x - \mu}{\sigma}\right)}$	$= \frac{1}{x\sigma} \frac{\phi\left(\frac{\log x - \mu}{\sigma}\right)}{\Phi\left(\frac{\log x - \mu}{\sigma}\right)}$

### 3. Applications

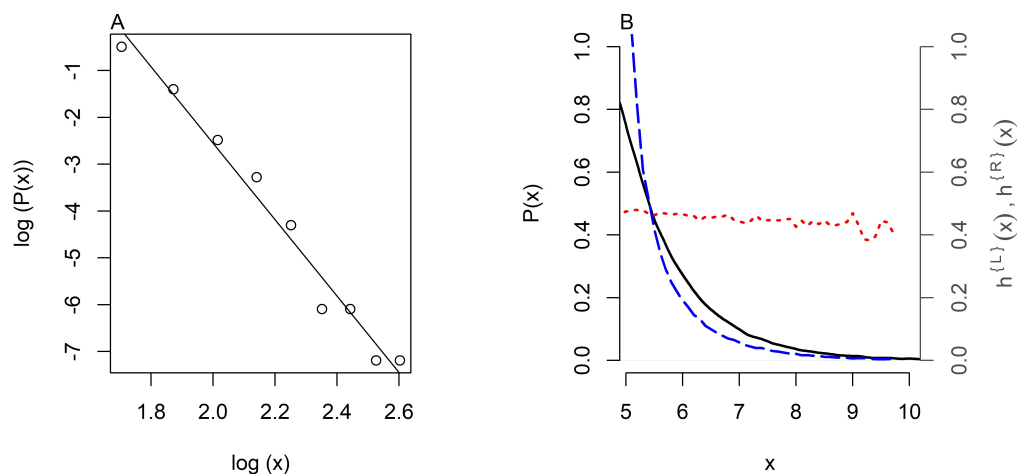
As discussed in the introduction, determining the distribution which best describes observed data is often important, as differing underlying distributions can have profound qualitative and quantitative differences. Specifically, accurately determining PL behaviour is important, as its properties, such as being scale-free, can be used or assumed in subsequent analyses and predictions [21]. Whilst there are many methods to ascertain the likely presence of PL behaviour [21,22,33], there are still occasions when it is difficult to decide between PL and other distributions. Here, we outline two such examples, and demonstrate how calculating the hazard and reverse-hazard functions can be beneficial in these cases. In all cases, the curves for  $h^{L}(x)$  and  $h^{R}(x)$  were found by removing data points above (below) the value  $x$  from the initial dataset and using the *denisty()* function in 'R' [34] to estimate the value of the distribution of the remaining data points at the value  $x$ .

#### 3.1. Tails of Distributions

One of the most common methods for determining PL behaviour is the log transform method (LTM), which considers whether a linear relation exists in the log–log plot of the distribution [22,33,35–38]. However, this is not infallible as other distributions can give straight line behaviour in log–log plots (Chapter 2 in [33]). Similarly, linear relationships can appear when only the tail of data or data points beyond a threshold value are used in the analysis. This truncating of data is often a necessity, and is seen in various areas of data analysis from risk modelling [4,39] to movement ecology [1,3].

Figure 1A demonstrates the log–log plot for an exponential distribution with rate  $\lambda = 1$  when only considering data points above 5. The log–log plot appears to demonstrate a linear relationship ( $R^2=0.9895$ ), which could therefore indicate PL behaviour. However, Figure 1A shows the results of calculating the hazard and reverse-hazard func-

tions. The hazard function (left truncation—red) gives an almost flat, horizontal curve, and the reverse-hazard function (right truncation—blue) describes a curve of the form,  $h^{\{R\}}(x) \propto x^{-k}$ .



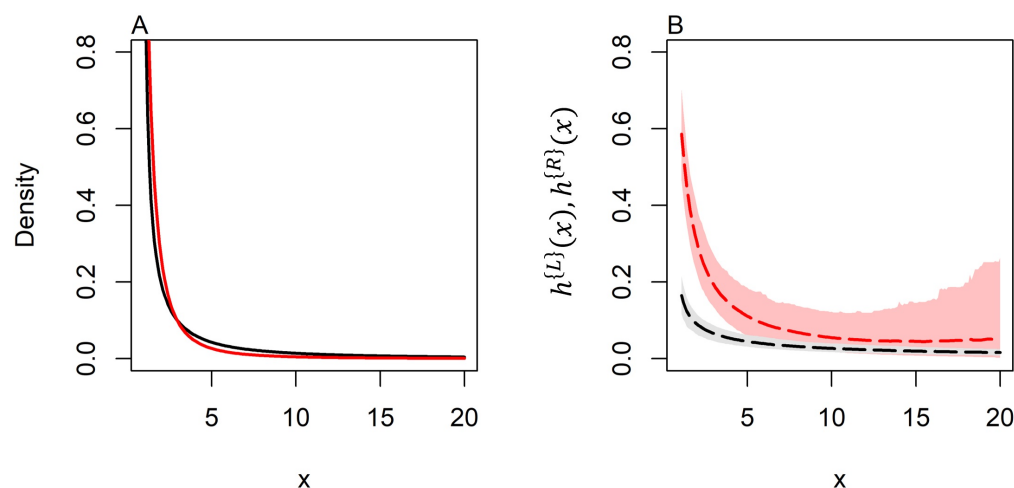
**Figure 1.** (A) log–log plot of data formed by drawing 100,000 data points from an exponential distribution with the rate  $\lambda = 1$  after a threshold value of  $x_{\min} = 5$  was used. The line of best fit was calculated by minimising the sum of squares, with an  $R^2$  value of 0.9895 indicating a good linear fit. (B) shows the density curve of the data after thresholding (black), with the hazard function given in red (truncation from the left),  $h^{\{L\}}$ , and blue for the reverse-hazard function (truncation from the right),  $h^{\{R\}}$ .

For a PL, we would expect a curve of the form  $h^{\{L\}}(x) \propto x^{-1}$  for truncating from the left, and  $h^{\{R\}}(x) \propto x^{-(\alpha+1)}$  for truncating on the right (Table 1). As this is clearly not the case, one can infer that the PL here is not a good fit for the data, and hence, other distributions should be considered, despite the indication from the LTM that a PL does potentially fit the data. Indeed, by considering Table 1, we note that an exponential distribution is a likely candidate, given the uniform distribution found from the left truncation and the inverse power in the right truncation.

### 3.2. Determining Between Power-Law and Log-Normal Distribution

The log-normal and PL distributions are known to closely resemble each other for certain parameter values and distinguishing between them is generally difficult [22]. Figure 2A shows the PDFs of a PL with  $\alpha = 2.5$  (red) and a log-normal distribution with mean = 0.3 and s.d. = 2 (black); both distributions have  $x_{\min} = 1$ . These were compared in [22], where it was noted that a large number of data points would be required to distinguish between them. Here, we consider datasets consisting of 300 random variables drawn from each distribution and compare the curves given when truncating the data from the left. Figure 2B depicts the expected shape of these characteristic curves depending on the distribution from which the data was drawn. The results indicate that despite only using a relatively small dataset of size 300, there is a noticeable difference in the observed curves when considering truncating the data from the left, particularly for values less than 10. Here, the truncated log-normal distribution gives a curve shallower than the initial distribution, whereas the truncated PL curve is much closer to the initial dataset.

This demonstrates that calculating the hazard function (truncating data from the left) could help give an indication of the better descriptive distribution when comparing between a log-normal and PL distribution even when the number of data points is small; a problem which is known to be a hard task.



**Figure 2.** (A) PDFs of a PL with  $\alpha = 2.5$  (red) and a log-normal distribution with mean = 0.3 and s.d. = 2 (black); both distributions have  $x_{\min} = 1$ . (B) The shape of the hazard function curves when datasets consisting of 300 random variables are truncated from the left. The results were found by running 1000 simulations, with dashed lines depicting the mean values and the shaded areas indicating 95% CIs.

### 3.3. Real-World Data

We now provide examples for when this qualitative method for determining the likely presence of PL behaviour in data can be used in real-world observed data. Each of these examples have been analysed previously, with a focus on determining whether a PL model is appropriate for the data.

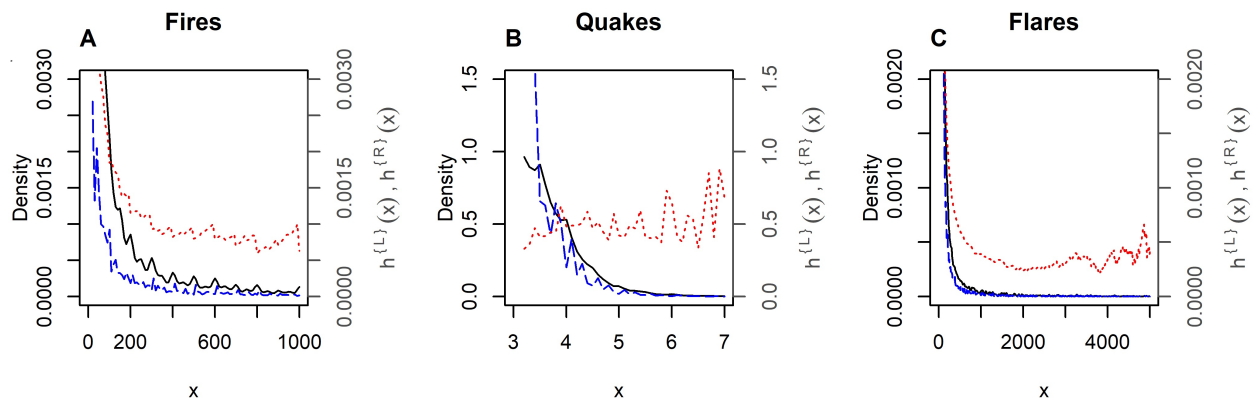
#### 3.3.1. Geophysical Data

Here, we consider three datasets which were analysed in Clauset et al. [22] and chosen as their properties all could indicate power-law behaviour. The three datasets as described in Clauset et al. [22] correspond to: (i) the acreage of wildfire on Federal Land between 1986 and 1996; (ii) earthquake intensity occurring in California between 1910 and 1992; (iii) the intensity of solar flares between 1980 and 1989, measured by the peak gamma-ray intensity. In all cases, more information about the datasets is provided in Newman [40].

In Clauset et al. [22], by deriving maximum likelihood estimators for both the scaling parameters,  $\alpha$ , and the  $x_{\min}$  value of a PL distribution, along with deriving a goodness-of-fit test based on the Kolmogorov–Smirnov test, it was determined that neither the wildfire data nor the earthquake data displayed power-law behaviour ( $p = 0.05$  and  $p = 0.00$ , respectively—Table 3 in Clauset et al. [22]), whereas the intensity in solar flares did demonstrate power-law behaviour ( $p = 1.00$ ). Although, it should be noted that the authors did report support for a power law with cut-off for wild fires.

Instead, by truncating these data and considering the form of the (reverse-) hazard function curves, we can quickly qualitatively determine whether a power law could potentially be a good fit to the observed data. Specifically, if we consider the hazard function (truncation from the left—red dotted line) for each dataset as shown in Figure 3, then for wildfires and solar flares (panels Figure 3A,C), the resulting curve gives something approximating a power law, which therefore indicates that the initial dataset (black solid line) could also be described by a power law. However, for the earthquake data (Figure 3B), the curve does not qualitatively resemble a power-law curve, appearing closer to a uniform curve, which indicates that the initial data are not well described by a power law (Table 1).

This corresponds to the findings in Clauset et al. [22]; however, this was achieved through a more simple and straight-forward approach.



**Figure 3.** Plots of the geophysical data as detailed in Section 3.3.1, with (A) presenting the acreage of wildfire, (B) the intensity of Californian earthquakes, and (C) the intensity of solar flares. In each case, the blue dashed line corresponds to the right truncated distribution ( $h^{\{R\}}(x)$ ), the red dotted line is the left truncated distribution ( $h^{\{L\}}(x)$ ), and the solid line is the initial distribution of the data.

### 3.3.2. Butterfly Movement

Here, we consider data describing the movement behaviour of butterflies, which was presented in Breed et al. [19]. The authors consider the step-length (distance between successive locations in recorded telemetry data) distribution of three datasets (Figure 4A), formed from the movement of populations of two species of butterfly: *Euphydryas editha taylori* and *Euphydryas phaeton*, featuring movement trajectories of one *E. e. taylori* colony from Corvallis, Oregon and two colonies of *E. phaeton* from Stevens–Coolidge Place (SCP) and Bullitt Meadow (BUL), respectively.

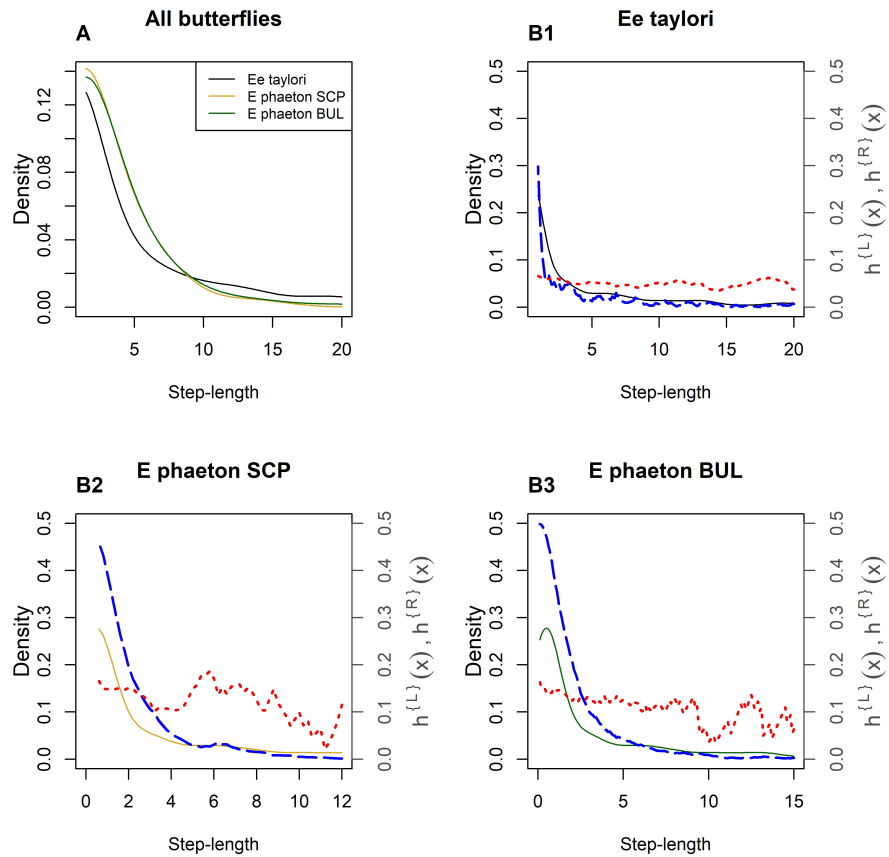
In their analysis, Breed et al. [19] highlighted that these data can be incorrectly identified as exhibiting power-law behaviour. They considered six distributions to model the step-length distributions: power law, exponential, bounded power law, bounded exponential, mixed distribution (formed of two exponentials), and Weibull. Using numerical MLE approaches following Edwards et al. [15,41], the authors identified that a power law was not as appropriate as the exponential-based distributions (either a mixture of exponential or Weibull) for all three datasets.

Similarly, by considering the truncated data distributions (Figure 4B1–B3), we note that the resulting distributions formed by truncating from the left (red dotted line) are not similar to the required power-law distribution, as would be expected (Table 1), therefore indicating that a power law is not apposite in these cases. Rather, given that in all cases the left truncation gives a distribution that is relatively flat and close to a uniform distribution, with the right truncation (blue dashed line) giving a PL-type distribution, Table 1 indicates that an exponential distribution may be most suitable for the observed data. In this case, this was also the finding of the authors in Breed et al. [19].

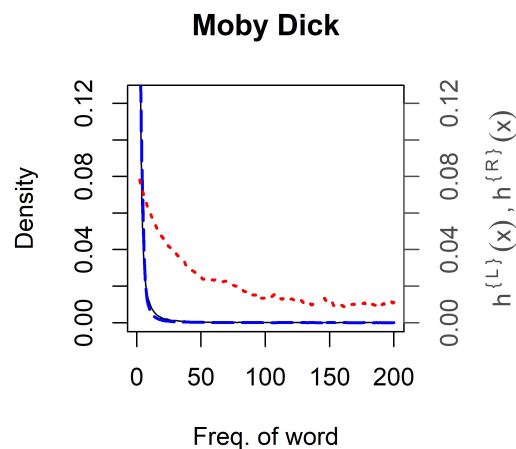
### 3.3.3. Moby Dick

Finally, we consider the dataset consisting of the number of times individual words are repeated in Herman Melville’s *Moby Dick*. This data are available as part of the *powerLaw* package in ‘R’ [23] and has been widely analysed and used as an example when analysing power-law or potential power-law behaviour [22,40,42]. In Clauset et al. [22], the authors demonstrated that there is evidence for the power law being a “good” fit (Table 4 in [22]); however, it has also been demonstrated that the data are best described by a log-normal distribution [40,42]. This confusion between a log-normal and power law is described in Section 3.2, with the hazard function for log-normal data giving a shallower curve than would be expected for a power law (Figure 2). Figure 5 demonstrates that the reverse-hazard function (blue dashed line) gives behaviour qualitatively similar to a power law, but the hazard function

(red dotted) gives a curve that could be a power law but appears as a much shallower curve, from Section 3.2. This could indicate that a log-normal is a better fit than a power law (as was found to be the case in Limpert et al. [42]), and that further analysis to directly compare the two distributions should be considered.



**Figure 4.** Plots showing the distribution of step-lengths (straight-line distance between successive locations) for butterfly data (see Section 3.3.2). (A) shows the distributions for all populations. (B1–B3) show the results found by calculating the hazard and reverse-hazard function distributions for each population. In each case, the blue dashed line corresponds to the right truncated distribution ( $h^{\{R\}}(x)$ ), the red dotted line is the left truncated distribution ( $h^{\{L\}}(x)$ ), and the solid line is the initial distribution of step-lengths.



**Figure 5.** Plot showing the distribution of the occurrence of unique words in the novel *Moby Dick*. The figure depicts the distribution of words with a frequency of below 200. The blue dashed line corresponds to the right truncated distribution ( $h^{\{R\}}(x)$ ), the red dotted line is the left truncated distribution ( $h^{\{L\}}(x)$ ), and the solid line is the initial distribution.

## 4. Discussion

Here, we have demonstrated how the truncating of data and the employment of the hazard and reverse-hazard function can be used to help identify the distribution which best describes data, specifically aiding in the complex task of identifying the potential presence of power-law behaviour. Our results show that using this simple approach, a qualitative result can be provided as to whether observed data may be well described by a power-law-type distribution, including when datasets feature relatively small number of samples.

Whilst it should be noted that no statistical test has been derived, and the method discussed here is purely indicative, we have however shown examples of how the hazard and reverse-hazard functions could be used in data analysis and demonstrated their benefits in specific examples such as when comparing between the similar distributions of a PL and log-normal when only a small number of data points is available.

As has been discussed, there are many approaches for statistically testing for the presence of a power law; however, they often come with a variety of caveats that need to be considered both a priori and in post hoc tests ([22,23]). The identification of power-law behaviour has been fraught in recent times, with various datasets analysed and reanalysed, either confirming or debating the proposed presence [19,20,43]. Due to this, a simple qualitative method to give researchers an indication as to whether the more statistically advanced and robust methodologies should be explored is a useful tool.

Overall, we have shown that using standard indicative features of datasets to identify PL distributions is not always sufficient, and that using a simple qualitative check concerning the truncation of data can give an indication as to when a deeper analysis should be carried out when attempting to ‘best-fit’ data to distributions.

**Funding:** This research received no external funding.

**Institutional Review Board Statement:** Not applicable

**Informed Consent Statement:** Not applicable

**Data Availability Statement:** All datasets used in this paper are referenced within the manuscript.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. Mashanova, A.; Oliver, T.H.; Jansen, V.A. Evidence for intermittency and a truncated power law from highly resolved aphid movement data. *J. R. Soc. Interface* **2010**, *7*, 199–208. [[CrossRef](#)] [[PubMed](#)]
2. Dahmen, H.; Wahl, V.L.; Pfeffer, S.E.; Mallot, H.A.; Wittlinger, M. Naturalistic path integration of *Cataglyphis* desert ants on an air-cushioned lightweight spherical treadmill. *J. Exp. Biol.* **2017**, *220*, 634–644. [[CrossRef](#)] [[PubMed](#)]
3. Bailey, J.D.; Benefer, C.M.; Blackshaw, R.P.; Codling, E.A. Walking behaviour in the ground beetle, *Poecilus cupreus*: Dispersal potential, intermittency and individual variation. *Bull. Entomol. Res.* **2021**, *111*, 200–209. [[CrossRef](#)] [[PubMed](#)]
4. Chernobai, A.; Burnecki, K.; Rachev, S.; Trück, S.; Weron, R. Modelling catastrophe claims with left-truncated severity distributions. *Comput. Stat.* **2006**, *21*, 537–555. [[CrossRef](#)]
5. Lagakos, S.W.; Barraj, L.M.; Gruttola, V.D. Nonparametric analysis of truncated survival data, with application to AIDS. *Biometrika* **1988**, *75*, 515–523. [[CrossRef](#)]
6. Dai, H.; Wang, H. *Analysis for Time-to-Event Data Under Censoring and Truncation*; Academic Press: Cambridge, MA, USA, 2016.
7. Lawless, J. Truncated distributions. *Encycl. Actuar. Sci.* **2006**, *3*.
8. Nadarajah, S.; Kotz, S. R programs for computing truncated distributions. *J. Stat. Softw.* **2006**, *16*, 1–8. [[CrossRef](#)]
9. Nadarajah, S. Some truncated distributions. *Acta Appl. Math.* **2009**, *106*, 105. [[CrossRef](#)]
10. Philippe, A. Simulation of right and left truncated gamma distributions by mixtures. *Stat. Comput.* **1997**, *7*, 173–181. [[CrossRef](#)]

11. Tokmachev, M. Modeling of truncated probability distributions. *IOP Conf. Ser. Mater. Sci. Eng.* **2018**, *441*, 012056 [[CrossRef](#)]
12. Lee, L.; Krutchkoff, R.G. Mean and variance of partially-truncated distributions. *Biometrics* **1980**, *36*, 531–536. [[CrossRef](#)]
13. Ho, H.J.; Lin, T.I.; Chen, H.Y.; Wang, W.L. Some results on the truncated multivariate t distribution. *J. Stat. Plan. Inference* **2012**, *142*, 25–40. [[CrossRef](#)]
14. Shonkwiler, J. Variance of the truncated negative binomial distribution. *J. Econom.* **2016**, *195*, 209–210. [[CrossRef](#)]
15. Edwards, A.M.; Phillips, R.A.; Watkins, N.W.; Freeman, M.P.; Murphy, E.J.; Afanasyev, V.; Buldyrev, S.V.; da Luz, M.G.; Raposo, E.P.; Stanley, H.E.; et al. Revisiting Lévy flight search patterns of wandering albatrosses, bumblebees and deer. *Nature* **2007**, *449*, 1044–1048. [[CrossRef](#)]
16. Codling, E.A.; Plank, M.J. Turn designation, sampling rate and the misidentification of power laws in movement path data using maximum likelihood estimates. *Theor. Ecol.* **2011**, *4*, 397–406. [[CrossRef](#)]
17. Auger-Méthé, M.; Derocher, A.E.; Plank, M.J.; Codling, E.A.; Lewis, M.A. Differentiating the Lévy walk from a composite correlated random walk. *Methods Ecol. Evol.* **2015**, *6*, 1179–1189. [[CrossRef](#)]
18. Benguigui, L.; Marinov, M. A classification of natural and social distributions Part one: The descriptions. *arXiv* **2015**, arXiv:1507.03408.
19. Breed, G.A.; Severns, P.M.; Edwards, A.M. Apparent power-law distributions in animal movements can arise from intraspecific interactions. *J. R. Soc. Interface* **2015**, *12*, 20140927. [[CrossRef](#)]
20. Pyke, G.H. Understanding movements of organisms: It's time to abandon the Lévy foraging hypothesis. *Methods Ecol. Evol.* **2015**, *6*, 1–16. [[CrossRef](#)]
21. Broido, A.D.; Clauset, A. Scale-free networks are rare. *Nat. Commun.* **2019**, *10*, 1017. [[CrossRef](#)]
22. Clauset, A.; Shalizi, C.R.; Newman, M.E. Power-law distributions in empirical data. *SIAM Rev.* **2009**, *51*, 661–703. [[CrossRef](#)]
23. Gillespie, C.S. Fitting Heavy Tailed Distributions: The powerLaw Package. *J. Stat. Softw.* **2015**, *64*, 1–16. [[CrossRef](#)]
24. Hanel, R.; Corominas-Murtra, B.; Liu, B.; Thurner, S. Fitting power-laws in empirical data with estimators that work for all exponents. *PLoS ONE* **2017**, *12*, e0170920. [[CrossRef](#)] [[PubMed](#)]
25. Mitzenmacher, M. A brief history of generative models for power law and lognormal distributions. *Internet Math.* **2004**, *1*, 226–251. [[CrossRef](#)]
26. Montebruno, P.; Bennett, R.J.; Van Lieshout, C.; Smith, H. A tale of two tails: Do Power Law and Lognormal models fit firm-size distributions in the mid-Victorian era? *Phys. A Stat. Mech. Its Appl.* **2019**, *523*, 858–875. [[CrossRef](#)]
27. Bartumeus, F.; Catalan, J.; Viswanathan, G.; Raposo, E.; Da Luz, M. The influence of turning angles on the success of non-oriented animal searches. *J. Theor. Biol.* **2008**, *252*, 43–55. [[CrossRef](#)]
28. Codling, E.A.; Bearon, R.N.; Thorn, G.J. Diffusion about the mean drift location in a biased random walk. *Ecology* **2010**, *91*, 3106–3113. [[CrossRef](#)]
29. Clark, T.G.; Bradburn, M.J.; Love, S.B.; Altman, D.G. Survival analysis part I: Basic concepts and first analyses. *Br. J. Cancer* **2003**, *89*, 232–238. [[CrossRef](#)]
30. Lee, S.H.; Urrutia, J.L. Analysis and prediction of insolvency in the property-liability insurance industry: A comparison of logit and hazard models. *J. Risk Insur.* **1996**, *63*, 121–130. [[CrossRef](#)]
31. Block, H.W.; Savits, T.H.; Singh, H. The Reversed Hazard Rate Function. *Probab. Eng. Inform. Sci.* **1998**, *12*, 69–90. [[CrossRef](#)]
32. Chechile, R.A. Properties of reverse hazard functions. *J. Math. Psychol.* **2011**, *55*, 203–222. [[CrossRef](#)]
33. Saichev, A.I.; Malevergne, Y.; Sornette, D. *Theory of Zipf's Law and Beyond*; Springer Science & Business Media: Berlin/Heidelberg, Germany, 2009; Volume 632.
34. R Core Team. *R: A Language and Environment for Statistical Computing*; R Foundation for Statistical Computing: Vienna, Austria, 2021.
35. Goldstein, M.L.; Morris, S.A.; Yen, G.G. Problems with fitting to the power-law distribution. *Eur. Phys. J. B-Condens. Matter Complex Syst.* **2004**, *41*, 255–258. [[CrossRef](#)]
36. Reynolds, A.M.; Smith, A.D.; Reynolds, D.R.; Carreck, N.L.; Osborne, J.L. Honeybees perform optimal scale-free searching flights when attempting to locate a food source. *J. Exp. Biol.* **2007**, *210*, 3763–3770. [[CrossRef](#)]
37. Sims, D.W.; Righton, D.; Pitchford, J.W. Minimizing errors in identifying Lévy flight behaviour of organisms. *J. Anim. Ecol.* **2007**, *76*, 222–229. [[CrossRef](#)]
38. Xiao, X.; White, E.P.; Hooten, M.B.; Durham, S.L. On the use of log-transformation vs. nonlinear regression for analyzing biological power laws. *Ecology* **2011**, *92*, 1887–1894. [[CrossRef](#)]
39. Levine, D. Modeling tail behavior with extreme value theory. *Risk Manag.* **2009**, *17*, 14–18.
40. Newman, M. Power laws, Pareto distributions and Zipf's law. *Contemp. Phys.* **2005**, *46*, 323–351. [[CrossRef](#)]
41. Edwards, A.M.; Freeman, M.P.; Breed, G.A.; Jonsen, I.D. Incorrect likelihood methods were used to infer scaling laws of marine predator search behaviour. *PLoS ONE* **2012**, *7*, e45174. [[CrossRef](#)]

42. Limpert, E.; Stahel, W.A.; Abbt, M. Log-normal Distributions across the Sciences: Keys and Clues: On the charms of statistics, and how mechanical models resembling gambling machines offer a link to a handy way to characterize log-normal distributions, which can provide deeper insight into variability and probability—Normal or log-normal: That is the question. *BioScience* **2001**, *51*, 341–352.
43. Benhamou, S. How many animals really do the Lévy walk? *Ecology* **2007**, *88*, 1962–1969. [[CrossRef](#)]

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.