

A Generalized Benford Framework for Threat Identification in Counter-Intelligence

January 2025

Author

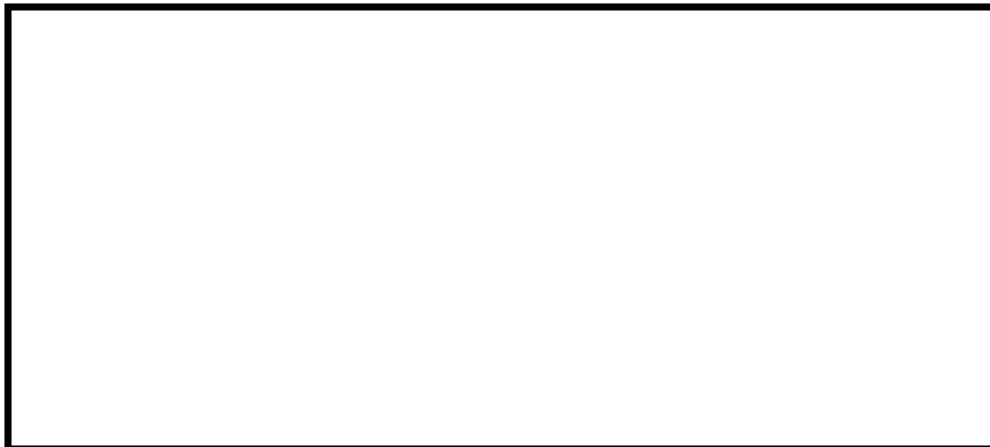
Timothy Tarter, James Madison University
Department of Mathematics
www.tartermathematics.com

Abstract

In this paper, we develop a framework of ‘Benford models’ for counter-intelligence investigations which analyze frequency data of a suspect’s visits to physical locations, online websites, and communication channels. We accomplish this by establishing the Benford measure for continuous & bounded domains, generalizing the accumulated percentage differences between sites in the frequency data with the log-determinant of ‘Benford Matrices,’ employing an estimator to determine a ‘Benford Test Statistic (λ),’ and identifying maximal values of λ across all permutations of included sites in our data. This framework is intended to complement outlier analysis models by finding where hidden Benford patterns ‘break’ in frequency data and telling investigators which sites they should investigate.

1 Overview

Before reading this paper, please take the following space below (or a nearby piece of paper) and draw seven X’s at random.



Without seeing your page in person, the author can tell you that the distance between a given point and each closest point is about the same between all points. To help illustrate this, see figure 1:

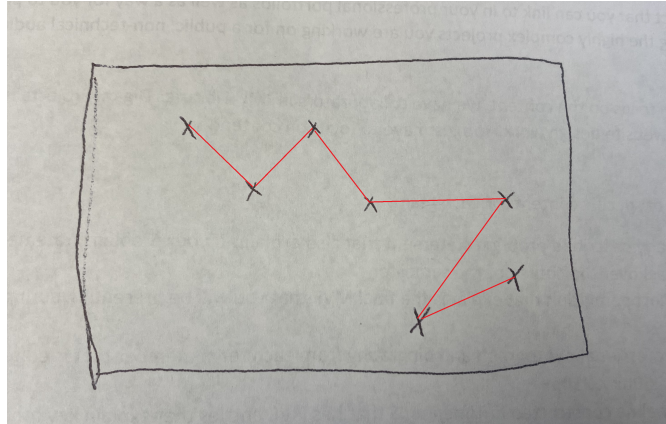


Figure 1: Demonstration

The key word in the prompt provided above was “random.” The human brain is notoriously poor at making random patterns when it is trying to (Małgorzata et. Al., 2008). This is one reflection of Benford’s Law which describes a surprising relationship between data randomly produced by natural processes and the probability distribution of leading digits in base-invariant representation. In the past, Benford’s Law has been used to detect financial fraud because the property is surprising to most individuals attempting to fabricate data.

Law of Anomalous Numbers A very familiar probability model examines an intuitive notion: if we are given nine items and one slot to put one of them in, the probability that a given number ends up in the slot is $\frac{1}{9}$ - i.e. equilikely. If we generalize this principle to nine numbers, $\{1, 2, 3, 4, 5, 6, 7, 8, 9\}$, we would assume that for any given slot, the distribution of numbers would be equilikely. To his surprise, however, in 1881, Simon Newcomb (and later Frank Benford) noticed that the leading digit in table entries of various datasets was more often one than two, two than three, three than four, *et cetera*.

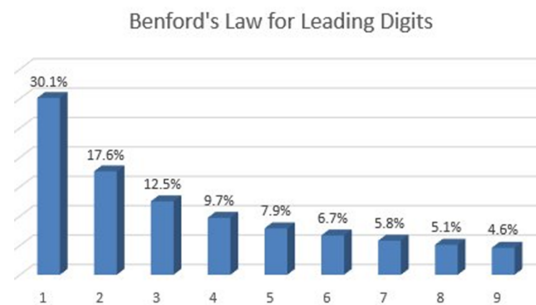


Figure 2: Law of Anomalous Numbers Visualization (Frost, 2022)

Early Applications Slightly after the discovery of base-invariance and implications of goodness of fit of Benford’s distribution across orders of magnitude (Miller, 2017), Benford’s Law was used to unearth Bernie Madoff’s ponzie scheme. Since then, “Benford Tests” have been used in financial crime investigations as a way to catch criminals altering bank statements, tax forms, and other types of money-laundering schemes.

2 Defining the Benford Measure

The Discrete Characterization of Benford Measure The Benford Distribution, denoted \mathbf{B} , is the unique measure on (\mathbb{R}^+, S) , where S is some leading digit (or for our purposes, generalized distance), such that the accumulation of leading digits $[1, t]$ follows a logarithmic distribution (for some base k , which is typically 10):

$$\mathbf{B} = \mu\left(\bigcup_{\ell \in \mathbb{Z}} k^\ell [1, t]\right) \simeq \ln(t), t \in \mathbb{N}_{[k]} \quad (1)$$

Moreover, its base-10 PDF is defined by:

$$P(S \leq t) = \log_{10}(t + 1) - \log_{10}(t), \forall t \in [1, 9] \quad (2)$$

Generalization of Benford Measure to a Continuous Domain This generalization is useful for continuous domains under a distance metric. It is also unique on its domain and carries the same properties as the above measure on \mathbb{R}^+ (Miller, 2017).

We previously established the discrete PDF for leading digits with a log base 10. Using the base invariance property of the distribution established by previous literature, we find that a generalized discrete Benford measure for base k is

$$P(S \leq t) = \log_k\left(1 + \frac{1}{t}\right), \forall t \in [1, k - 1]. \quad (3)$$

We now want to motivate two things:

1. A continuous domain for t
2. Moving \log_k to \log_e

Motivating the log-base e is fairly simple. Since the original measure follows from the accumulation of $\frac{1}{t}$, we can let a be in the neighborhood of some $t \in [1, k]$. Then the image of a through the map:

$$f : X \rightarrow Y, X \in (0, 1], Y \in (-\infty, \ln(k)] \quad (4)$$

will be in the neighborhood of the image of t . Additionally, let d denote a distance metric on both X and Y . Formally, letting $\epsilon > 0$, if $a \in n_\epsilon(t)$ then $\exists \delta > 0$ such that $f(a) \in n_\delta(f(t))$. Moreover, since $\delta \in \mathbb{R}^+$, by the Archimedean property, we find that $\forall a, t \in X, \exists K \in \mathbb{R}^{\geq 0}$ with $d(f(x_1), f(x_2)) \leq K d(x_1, x_2)$. This means that f is Lipschitz continuous. Therefore, as the domain of X expands (moving the base of our distribution towards ∞), our property of Lipschitz continuity allows us to express an approximation of that continuous distribution with log-base e .

Now, we want to justify $t \in \mathbb{R}^+$, not just \mathbb{Z}^+ . For this argument, we consider \mathbf{B} as an approximation of $\ln(t)$. While the measure was defined this way for the aforementioned purpose of analyzing discrete values, we find that moving t to a continuous domain in \mathbb{R} is analogous to the Riemann sum converging to a definite integral. Thus, we introduce the measure

$$\tilde{\mathbf{B}} = \bigcup_{k \in \mathbb{Z}} 10^k [1, t] = \ln(t), t \in \mathbb{R}^+ \quad (5)$$

With a continuity approximation error term of $O(\mathbf{B} - \tilde{\mathbf{B}})$. Our PDF becomes,

$$P(S \leq t) = \ln\left(1 + \frac{1}{t}\right), \forall t \in \left[\frac{1}{e-1}, \infty\right). \quad (6)$$

Again, recall that the justification for the bounds of the distribution comes from our argument of Lipschitz continuity and nearness of fit to continuum-normalized distributions. Making a simple substitution of variables, $x = 1 + \frac{1}{t}$, we obtain a “better behaved” distribution:

$$P(S \leq x) = \ln(x), x \in [1, e] \quad (7)$$

3 Counter-Intelligence Threat Identification

This research started because of a thought experiment about the famous spy, Robert Hanssen. In 1976, Hanssen was hired as an FBI agent and quickly promoted to a position in the FBI’s counter-intelligence department. Between 1979 and 2001, Hanssen played turncoat, providing Soviet and Russian intelligence agencies with secrets vital to the state of American National Security. Arguably the most alarming aspect of Hanssen’s tale is that he managed to do this while under near-constant surveillance. This highlights certain key flaws of otherwise powerful forms of outlier analysis for counter-intelligence investigations: one must assume that outliers will occur before substantial damage is done. Hanssen’s story prompted the question: given similar circumstances (constant surveillance via geo-spatial tagging, monitoring communication, etc.), how can we prevent Hanssen-esque threats in the future? Specifically, can Benford models offer promising solutions by highlighting subtle, hidden patterns in data?

3.1 Frequency Analysis

To specify a use case for our solution, we want to narrow down the types of data that we intend to use to frequency data about the number of times our suspect visits a set of locations over a certain period of time. This is powerful because it easily generalizes from a suspect visiting physical locations to the communication habits, online presence, and spending patterns of a suspect. More specifically, however, we want to compare the percent change between the frequency counts at each site, the same way that we compared distance in [1].

Lemma 1 *Let f_i denote the number of visits to a site i , letting n be the total possible number of sites. Then, $g_{i,j} \simeq \ln(f_i) - \ln(f_j)$ is the percent change between f_i and f_j .*

Proof (Sims - Math Review, 2015):

- We know the relative difference frequency, $g_i = \frac{f_i - f_j}{f_j} = \frac{f_i}{f_j} - 1$

- Then $\frac{f_i}{f_j} = g_{i,j} + 1$
- $\ln(g_i + 1) = \ln(f_i) - \ln(f_j) \simeq g_{i,j}$

□

3.2 Order Invariance of Sites & Generalizing μ

We use this to compare percentage changes between visits at sites using $\ln(\frac{f_i}{f_j})$. That said, in the real world, we aren't just looking at two sites, and there isn't an intrinsically correct ordering (as opposed to time-series data for example). This motivates us to want to generalize the percentage change in each site with respect to every other site via **Benford Matrices**.

Consider a trivial example where there are three sites which yield the 3 x 3 **Benford Matrix**:

$$A = \begin{bmatrix} f_1/f_1 & f_2/f_2 & f_3/f_3 \\ f_1/f_2 & f_2/f_3 & f_3/f_1 \\ f_1/f_3 & f_2/f_1 & f_3/f_2 \end{bmatrix} \quad (8)$$

Aside: for 'n' sites, this will be an n x n matrix.

Then to capture the total percentage of differences between sites, Δf , let $\ln|\det(A)| = \Delta f$ (since the determinant generalizes the accumulated area of a collection of vectors). Additionally, this is justified by the fact that multiples of Benford variables and their reciprocals are also Benford (Pike, 2008). Since the absolute value of a determinant divided by the number of vectors in its matrix tells you the average scaled volume of its column space (i.e. how spread out the vectors are), our formula will tell us how spread out the patterns in our data are. If that spread is similar to a "Benford spread," then the patterns in our suspect's data are Benford patterns. If they aren't similar, then there is an outlier to the expected Benford pattern for natural processes. Note that the expected average scaled volume is the expected value of a Benford variable, which we have previously established to be ~ 2.0973 . This is a particularly nice generalization for the following reason: we can reasonably expect the total percentage of differences between sites to be equal to the expected value of the continuous Benford distribution if and only if the true frequency data is Benford. Thus, for a continuous random variable, $B \sim \text{Benford}$, and for n sites of f ,

Theorem 1 $E[\mathbf{B}] = \frac{\Delta f}{n} = \frac{\ln|\det(A)|}{n}$ iff $\bigcup_i f_i \sim \text{Benford}$

Note: using our "well-behaved" generalization of Benford Measure to a continuous (but bounded) domain in (7), with X a continuous random Benford variable, $f_X(x) = \ln(x)$, we derive:

$$E[X] = \int_1^e x \ln(x) dx \simeq 2.0973 + O(\mathbf{B} - \tilde{\mathbf{B}}) \quad (9)$$

3.3 The Benford Test Statistic

Given the nature of the intended application for this framework, we know that this parity relation will very rarely be held. In fact, because we want to know where to advise investigators to “look” for illicit activity from a suspect, it may even be more useful to study the magnitude of perturbation caused by introducing a site to the dataset than merely noting that the parity relation is broken. Therefore, let us employ an estimator, λ , our **Benford Test Statistic**, to tell us how broken our parity relation is.

$$\lambda = 2.0973 - \frac{\ln|\det(A)|}{n} \quad (10)$$

It is worth noting that if the sample data is Benford, then λ will be in some close neighborhood of zero. Since we assume our sample to be Benford, then the radius of that neighborhood in \mathbb{R}^1 will be the standard deviation of the Benford distribution for one variable.

3.3.1 Higher Moments of the Benford Distribution

We calculate these higher moments using our distribution $\tilde{\mathbf{B}}$ from (7):

$$E[X] = \int_1^e x \ln(x) dx \simeq 2.0973 + O(\mathbf{B} - \tilde{\mathbf{B}}) \quad (11)$$

$$E[X^2] = \int_1^e x^2 \ln(x) dx \simeq 4.5746 + O(\mathbf{B} - \tilde{\mathbf{B}}) \quad (12)$$

$$V[X] = 4.5746 - 2.0973^2 \simeq 0.1759 \quad (13)$$

$$\sigma[X] = \sqrt{0.1759} \simeq 0.4149 \quad (14)$$

3.3.2 Hypothesis Testing Formula

From these higher moments, we can perform hypothesis tests in the following manner:

$$H_0: \lambda = 0; H_A: \lambda \neq 0$$

Test Statistic (Studentized t-test):

$$t = \frac{0 - \lambda}{0.4149} \quad (15)$$

Our p-value can be calculated with df (degrees of freedom) = $n^2 - 1$ (since our Benford Matrix is $n \times n$), and can be reasonably evaluated at the 95% level of significance.

3.4 Maximizing λ

Recalling our motivation for finding λ in [3.3], we really want to determine which entries $a_{i,j} \in A$ maximize λ when added to our Benford Matrix, A . Since our sites are order-invariant, we unfortunately have to select out one column at a time (which of course also removes a row, making

A $(n - 1) \cdot (n - 1)$), which is somewhat computationally expensive. Moreover, it isn't enough to just do this for one removed row / column, we have to do it for every combination of sites which could be removed from the dataset. Thus, for 'n' sites, we will have to run $n \cdot n!$ computations of λ , and then employ an optimization algorithm to find $\max_{i,j} \lambda_{i,j}$. Simplifying this permutation algorithm is classified under active future research.

3.5 Numerical Simulation in Python

To validate our Benford Framework for frequency analysis, we programmed numerical simulations in Python which sample from various distributions and perform the steps detailed above for Benford our tests.

If the reader is so inclined, they can access the code on my website at:

www.tartermathematics.com

Otherwise, here are some screenshots of the example outputs from the code:

```
benfordMtx(8,7,'uniform')
✓ 0.0s
Sites Visited: [7, 1, 3, 5, 6, 3, 6, 7]
-----
Benford Matrix:
-----
```

	1	2	3	4	5	6	7	8
1	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
2	7.0	0.333333	0.6	0.833333	2.0	0.5	0.857143	1.0
3	2.333333	0.2	0.5	1.666667	1.0	0.428571	0.857143	7.0
4	1.4	0.166667	1.0	0.833333	0.857143	0.428571	6.0	2.333333
5	1.166667	0.333333	0.5	0.714286	0.857143	3.0	2.0	1.4
6	2.333333	0.166667	0.428571	0.714286	6.0	1.0	1.2	1.166667
7	1.166667	0.142857	0.428571	5.0	2.0	0.6	1.0	2.333333
8	1.0	0.142857	3.0	1.666667	1.2	0.5	2.0	1.166667

```
-----
Determinant: 35508
-----
Mu(Det):1.3096888215
-----
Lambda:1.601372757839979
-----
Characteristic Equation: 0.0
-----
Test statistic: -1.449440245456686, Degrees of Freedom: 63
P-Value: 0.076087283835627
We fail to reject the null hypothesis and conclude that the sequence is Benford
```

Figure 3: Sample Benford Set

```

benfordMtx(5,7,'uniform')
✓ 0.0s
Sites Visited: [6, 1, 3, 2, 2]
-----
Benford Matrix:
-----

```

	1	2	3	4	5
1	1.0	1.0	1.0	1.0	1.0
2	6.0	0.333333	1.5	1.0	0.333333
3	2.0	0.5	1.5	0.333333	2.0
4	3.0	0.5	0.5	2.0	0.666667
5	3.0	0.166667	3.0	0.666667	1.0

```

-----
Determinant: 22
-----
Mu(Det):0.6138688182
-----
Lambda:3.4165279908172708
-----
Characteristic Equation: 0.0
-----
Test statistic: -5.824362474854834, Degrees of Freedom: 24
P-Value: 2.627898875297457e-06
We reject the null hypothesis and conclude that the sequence is not Benford

```

Figure 4: Sample ‘Not’ Benford Set

4 Conclusion

In this work, we generalize the Benford Measure beyond the Significant σ -Algebra, characterize a continuous distribution, as well as a bijection from it to a bounded continuous distribution, establish an algorithm for determining if a sequence of frequency data is distributed Benford, and provide an algorithm for analyzing system perturbations. These results indicate likely highly material performance yields for future research on Benford models, especially for applications in national security and threat prevention.

5 References

- Frost, J. (2022, October 6). Benford's Law explained with examples. Statistics By Jim.
<https://statisticsbyjim.com/probability/benfords-law/>
- Park, S.-H., Huh, S.-Y., Oh, W., & Han, S. P. (2012). A Social Network-Based Inference Model for Validating Customer Profile Data. *MIS Quarterly*, 36(4), 1217–1237.
<https://doi.org/10.2307/41703505>
- Morzy, Mikolaj & Kajdanowicz, Tomasz & Szymanski, Boleslaw. (2016). Benford's Distribution in Complex Networks. *Scientific Reports*. 6. 34917. 10.1038/srep34917.
- Miller, S. J. (2017). Benford's law: Theory and applications. Princeton University Press.
- (PDF) Benford's law and the $C\beta e$. (n.d.).
https://www.researchgate.net/publication/368304741_Benford's_law_and_the_CbetaE
- Kopczewska, K., & Kopczewski, T. (2022, October 20). Natural spatial pattern-when mutual socio-Geo Distances between cities follow Benford's law. *PLoS one*.
<https://pmc.ncbi.nlm.nih.gov/articles/PMC9584388/>
- Author links open overlay panel Małgorzata Figurska a, a, b, 1, believed, S. is widely, Brown, R. G., Daniels, C., Joppich, G., Persaud, N., Schneider, S., Pollux, P. M. J., & Baddeley, A. D. (2007, September 20). Humans cannot consciously generate random numbers sequences: Polemic study. *Medical Hypotheses*.
<https://www.sciencedirect.com/science/article/abs/pii/S030698770700480X>
- Pike, D. (2008). Testing for the benford property. *SIAM*.
https://www.siam.org/media/r3alxzk0/testing_for_the_benford_property.pdf
- Friar, J. L., Goldman, T., & Perez-Mercader, J. (2016, April 7). Ubiquity of Benford's law and emergence of the Reciprocal Distribution. *Physics Letters A*.
<https://www.sciencedirect.com/science/article/abs/pii/S0375960116300603>
- Sims, E. (n.d.). University of Notre Dame. Eric Sims.
<https://sites.nd.edu/esims/>