

Benford's Law or the law of anomalous digits

In certain natural collections of numbers like sizes of towns, heights of buildings, lengths of rivers, financial transactions the leading digit

feature.

has a initially surprising

Since the lead digit is 1, 9 one would

expect that each digit would occur $\frac{1}{9} \approx 11\%$ of

the time but in actuality 1 is the leading digit about 30% of the time and 9 about

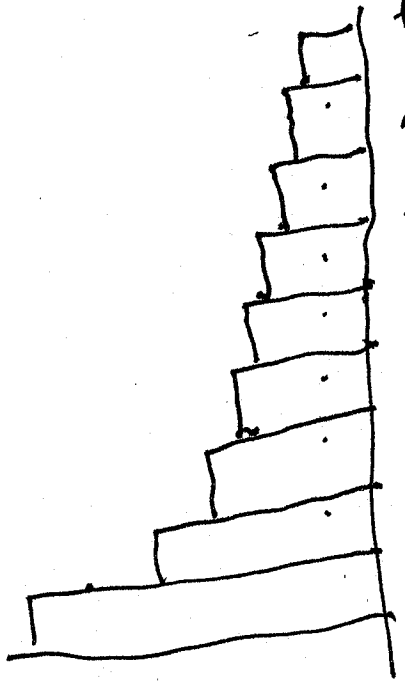
5% of the time.

(2)

DEF: A data set satisfies Benford's law if the frequency of occurrence of d as the leading digit is

$$P(d) = \log_{10} \left(\frac{d+1}{d} \right)$$

approximate plot



has leading

Initial note: If $x \in [d, d+1)$ it

digit d if $x \in [d, d+1)$ or

$$d \leq x < d+1$$

Taking \log_{10}

$$\log_{10} d \leq \log_{10} x < \log_{10} d+1$$

$$\text{or } \log_{10} x \in [\log_{10} d, \log_{10} d+1)$$

which is an interval with width or Lobes or measure

$$\log_{10} d+1 - \log_{10} d = \log_{10} \frac{d+1}{d} = P(d)$$

So if all the data is in $[1, 10)$, being Benford

is assuming that the Lebesgue measure

distributed w.r.t. d, d_1, d_2, \dots

In general, if $x = d_1 d_2 \dots \times 10^k$ for some k so

$$\log_{10} x = \log d_1 d_2 \dots \times 10^k = (\log d_1 d_2 \dots) + k$$

So the general assumption on being Benford is that the log of the data mod 1 is uniformly distributed with respect to Lebesgue.

In certain cases one can prove data is Benford

using Ergodic Theory

Lemma: $R_\alpha: S^1 \rightarrow S^1, x \mapsto x + \alpha$ has the property that for any interval ~~in S^1~~ J in S^1

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=0}^{n-1} \chi_J R_\alpha^i(x) \rightarrow \int \chi_J d\mu = \mu(J)$$

for all $x \in S^1$ where $\mu =$ Lebesgue measure

so $\mu(J)$ is just the length of J .

Proof: This is almost unique ergodicity but χ_J isn't continuous so one needs to approximate it by continuous functions.

Theorem The set $\{2^n\}_{n \in \mathbb{N}} = \{1, 2, 4, 8, 16, \dots\}$

satisfies Benford's law.

As noted above 2^n has leading

Proof digit d if and only if $(2^n) \in [\log_{10} d, \log_{10} d + 1) \pmod{1}$

or $n \log_{10} 2 \in \mathbb{J}_d$

Define $R: \mathbb{S}^1 \rightarrow \mathbb{S}^1$ via $\theta \mapsto \theta + \log_{10} 2$ and note that $\log_{10} 2$ is irrational (1) and

$n \log_{10} 2 = R^n(0)$

3

By the lemma

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{l=0}^{n-1} \chi_{\mathbb{F}_d}(R^n(\omega)) = \mu(\mathbb{F}_d) = \log \frac{p}{1+p}$$

is counting the average number of times (asymptotically) that d occurs as the leading digit of some 2^n .

Similar arguments show that daily balances with compound interest $(R)^n P_0$ are Benford.

deg $(1.03)^n(20,000)$ are usually

- Fraudulent checks, deposits, votes, ... are usually uniformly distributed and so don't satisfy Benford's law, so it can be used as one part of a Fraud detection Suite.

We need one more proof to properly wrap up entropy:

Lemma: If $\rho = \sum A_i, A_n, \dots, m = \{B_1, \dots, B_m\}$

$$\text{Then } H(\rho \vee \eta) \leq H(\rho) + H(\eta)$$

terms

Proof: For simplicity we don't include terms

in the sums when $m(B_j) = 0$ or $m(A_i \wedge B_j) = 0$

$$H(\rho \vee \eta) = - \sum_{i,j} m(A_i \wedge B_j) \log m(A_i \wedge B_j)$$

$$= - \sum_{i,j} m(A_i \wedge B_j) \log \left(\frac{m(A_i \wedge B_j)}{m(B_j)} m(B_j) \right)$$

$$= - \sum_{i,j} \mu(A_i \wedge B_j) \log \left(\frac{\mu(A_i \wedge B_j)}{\mu(B_j)} \right)$$

$$= - \sum_{i,j} \mu(A_i \wedge B_j) \log \mu(B_j)$$

$$= - \sum_{i,j} \frac{\mu(A_i \wedge B_j)}{\mu(B_j)} \mu(B_j) \log \left(\frac{\mu(A_i \wedge B_j)}{\mu(B_j)} \right)$$

$$= - \sum_j \mu(B_j) \log \mu(B_j)$$

since ρ is a partition

$$= \sum_j \mu(A_i \wedge B_j) = \mu(B_j)$$

using the fact that $\sum_i \mu(A_i \wedge B_j) = \mu(B_j)$

$$\text{Thus } H(\rho \vee \eta) = - \sum_{i,j} \frac{\mu(A_i \wedge B_j)}{\mu(B_j)} \mu(B_j) \log \left(\frac{\mu(A_i \wedge B_j)}{\mu(B_j)} \right)$$

$$= H(\eta) \quad (*)$$

Now fix i and let $\gamma_j = \mu(B_j)$ and so $\sum_j \gamma_j = 1$

19

and let $\alpha_j = \frac{\mu(A_L \cap B_j)}{\mu(B_j)}$ so $\alpha_j \in (0, 1)$ and

$$\sum_{j=1}^m \gamma_j \alpha_j = \sum_j \mu(B_j) \frac{\mu(A_L \cap B_j)}{\mu(B_j)} = \sum_j \mu(A_L \cap B_j) = \mu(A_L)$$

Now recall if $\phi(x) = -x \log x$ $x \in (0, 1)$, $\phi(0) = 0$

Then we showed $\sum_{j=1}^m \gamma_j \phi(\gamma_j) \leq \phi\left(\sum_{j=1}^m \gamma_j \alpha_j\right)$

$$\text{Thus } -\sum_{j=1}^m \mu(B_j) \frac{\mu(A_L \cap B_j)}{\mu(B_j)} \log \frac{\mu(A_L \cap B_j)}{\mu(B_j)} \leq \sum_j \gamma_j \phi(\gamma_j)$$

$$\leq \phi\left(\sum_j \gamma_j \alpha_j\right) = \phi(\mu(A_L)) = -\mu(A_L) \log \mu(A_L)$$

For each i .

Summing over i yields

$$- \sum_{i,j} \mu(B_j) \frac{\mu(A_i | B_j)}{\mu(B_j)} \log \frac{\mu(A_i | B_j)}{\mu(B_j)}$$

$$\leq H(p)$$

back to equation (*) we have

Then going

$$H(p, v_m) \leq H(p) + H(m)$$