Introduction	Preliminaries	Results	Conclusion

# Benford's Law and Dependent Random Variables

Thealexa Becker<sup>1</sup>, Alexander Greaves-Tunnell<sup>2</sup>, Steven J. Miller<sup>2</sup>, and Ryan Ronan<sup>3</sup>

<sup>1</sup>Smith College <sup>2</sup>Williams College <sup>3</sup>Cooper Union

Joint Meetings of the AMS / MAA, Boston January 5, 2011

Introduction ●○○	Preliminaries	Results 000	Conclusion
Ronford's Low			

### Benford's Law: Newcomb (1881), Benford (1938)

For many data sets, probability of observing a first digit of *d* base *B* is  $\log_B(\frac{d+1}{d})$ .

Introduction	Preliminaries	Results	Conclusion
●○○	000	000	
Benford's Law			

# Benford's Law: Newcomb (1881), Benford (1938)

For many data sets, probability of observing a first digit of *d* base *B* is  $\log_B(\frac{d+1}{d})$ .

# Examples:

- iterates of 3x+1, along with power, exponential and rational maps
- products of independent random variables
- particle decay, hydrology and financial data



Introduction ●○○	Preliminaries	Results 000	Conclusion
Benford's Law			

# Benford's Law: Newcomb (1881), Benford (1938)

For many data sets, probability of observing a first digit of *d* base *B* is  $\log_B(\frac{d+1}{d})$ .

### Examples:

- iterates of 3x+1, along with power, exponential and rational maps
- products of independent random variables
- particle decay, hydrology and financial data

Question: Which systems lead to Benford behavior?

Introduction	Preliminaries	Results	Conclusion
○●○	000	000	
Dependence Consid	lerations		

# Most efforts to identify Benford processes assume random variables are *independent*.

Introduction	Preliminaries	Results	Conclusion
○●○	000	000	
Dependence Con	siderations		

# Most efforts to identify Benford processes assume random variables are *independent*.

Motivated by "real world" processes shown to be Benford, we want to explore dependent systems.



Introduction	Preliminaries	Results	Conclusion
○●○	000	000	oo
Dependence Con	siderations		

# Most efforts to identify Benford processes assume random variables are *independent*.

Motivated by "real world" processes shown to be Benford, we want to explore dependent systems.

**Our motivating example (Lemons 1986):** Decomposition of a conserved quantity as a model for particle decay.

Introduction	Preliminaries	Results	Conclusion
000			

A discrete-time, continuous-division model of decay:

• "Stick" of length *L* divided in two at each stage.



Introduction	Preliminaries	Results	Conclusion
000			

A discrete-time, continuous-division model of decay:

- "Stick" of length *L* divided in two at each stage.
- Cuts from continuous random variable K on (0, 1).

Introduction	Preliminaries	Results	Conclusion
000			

A discrete-time, continuous-division model of decay:

- "Stick" of length *L* divided in two at each stage.
- Cuts from continuous random variable K on (0, 1).
- Pieces  $\{X_i\}$  not independent, as must sum to *L*.

Introduction	Preliminaries	Results	Conclusion
000			

A discrete-time, continuous-division model of decay:

- "Stick" of length *L* divided in two at each stage.
- Cuts from continuous random variable K on (0, 1).
- Pieces  $\{X_i\}$  not independent, as must sum to *L*.



Introduction	Preliminaries	Results	Conclusion
000	000	000	00

#### Equidistribution and Benford's Law

**Significand:** for positive x, write  $x = S_{10}(x) \cdot 10^{k(x)}$ .

Introduction	Preliminaries	Results	Conclusion
	<b>•</b> 00		

#### Equidistribution and Benford's Law

**Significand:** for positive x, write  $x = S_{10}(x) \cdot 10^{k(x)}$ .

**Equidistribution:**  $\{y_i\}$  equidistributed mod 1 if for any  $[a, b] \subset [0, 1]$  have  $\lim_{N \to \infty} \frac{\#\{n \le N: y_n \in [a, b]\}}{N} = b - a$ .

For  $s \in [1, 10)$ , let

 $\varphi_{s}(u) := \begin{cases} 1 & \text{if } u \text{'s significand is at most } s \\ 0 & \text{otherwise.} \end{cases}$ 

Introduction	Preliminaries	Results	Conclusion
	<b>0</b> 00		

#### Equidistribution and Benford's Law

**Significand:** for positive x, write  $x = S_{10}(x) \cdot 10^{k(x)}$ .

**Equidistribution:**  $\{y_i\}$  equidistributed mod 1 if for any  $[a, b] \subset [0, 1]$  have  $\lim_{N \to \infty} \frac{\#\{n \le N: y_n \in [a, b]\}}{N} = b - a$ .

For  $s \in [1, 10)$ , let

 $\varphi_{s}(u) := \begin{cases} 1 & \text{if } u \text{'s significand is at most } s \\ 0 & \text{otherwise.} \end{cases}$ 

#### **Fundamental Equivalence**

Data set  $\{x_i\}$  is Benford base B if  $\{y_i\}$  is equidistributed modulo 1, where  $y_i = \log_B x_i$ .

Introduction	Preliminaries	Results	Conclusion
000	○●○	000	oo
The Mellin Tra	nsform		

Let f(x) be a continuous real-valued function on  $[0, \infty)$ . Define its Mellin transform,  $(\mathcal{M}f)(s)$ , by

$$(\mathcal{M}f)(s) := \int_0^\infty f(x) x^s \frac{dx}{x}.$$



Let f(x) be a continuous real-valued function on  $[0, \infty)$ . Define its Mellin transform,  $(\mathcal{M}f)(s)$ , by

$$(\mathcal{M}f)(s) := \int_0^\infty f(x) x^s \frac{dx}{x}.$$

Note that (*Mf*)(*s*) = E[*x*<sup>*s*-1</sup>], thus results concerning expected values translate to results on Mellin transforms.



Let f(x) be a continuous real-valued function on  $[0, \infty)$ . Define its Mellin transform,  $(\mathcal{M}f)(s)$ , by

$$(\mathcal{M}f)(s) := \int_0^\infty f(x) x^s \frac{dx}{x}.$$

- Note that (*Mf*)(*s*) = E[*x*<sup>*s*-1</sup>], thus results concerning expected values translate to results on Mellin transforms.
- With a logarithmic change of variables, we can translate between Mellin and Fourier transforms.

Introduction	Preliminaries	Results	Conclusion
	00•		

#### **Products of Independent Random Variables**

# Jang, Kang, Kruckman, Kudo and Miller (2009)

 $\Xi_1, \ldots, \Xi_N$  independent variables with densities  $f_{\Xi_m}$  and

$$\lim_{N\to\infty}\sum_{\ell=-\infty\atop\ell\neq 0}^{\infty}\prod_{m=1}^{N}(\mathcal{M}f_{\Xi_m})\left(1-\frac{2\pi i\ell}{\log B}\right) = 0 \quad (C1)$$

Then as  $\Xi_1 \cdots \Xi_N$  converges to Benford's law.

Introduction	Preliminaries	Results	Conclusion
	000		

#### **Products of Independent Random Variables**

# Jang, Kang, Kruckman, Kudo and Miller (2009)

 $\Xi_1, \ldots, \Xi_N$  independent variables with densities  $f_{\Xi_m}$  and

$$\lim_{N\to\infty}\sum_{\ell=-\infty\atop\ell\neq 0}^{\infty}\prod_{m=1}^{N}(\mathcal{M}f_{\Xi_m})\left(1-\frac{2\pi i\ell}{\log B}\right) = 0 \quad (C1)$$

Then as  $\Xi_1 \cdots \Xi_N$  converges to Benford's law.

• (C1) is quite weak, met by most distributions.

Introduction	Preliminaries	Results	Conclusion
	00•		

#### **Products of Independent Random Variables**

# Jang, Kang, Kruckman, Kudo and Miller (2009)

 $\Xi_1, \ldots, \Xi_N$  independent variables with densities  $f_{\Xi_m}$  and

$$\lim_{N\to\infty}\sum_{\substack{\ell=-\infty\\\ell\neq 0}}^{\infty}\prod_{m=1}^{N}(\mathcal{M}f_{\Xi_m})\left(1-\frac{2\pi i\ell}{\log B}\right) = 0 \quad (C1)$$

Then as  $\Xi_1 \cdots \Xi_N$  converges to Benford's law.

- (C1) is quite weak, met by most distributions.
- **Corollary:** For  $\Xi_1, \ldots, \Xi_N$  uniformly distributed on (0, 1) and  $N \ge 4$ ,  $|\varphi_s(u) - \log_{10} s| \le \left(\frac{1}{2.9^N} + \frac{\zeta(N) - 1}{2.7^N}\right) 2 \log_{10} s.$

Introduction	Preliminaries	Results	Conclusion
		000	

#### **Limiting Behavior of Decompostions**

Theorem (Decomposition model as above)

Fix a continuous density f on [0, 1] st f(x), f(1 - x) satisfy (C1), set

$$P_N(s) := rac{\sum_{i=1}^{2^N} \varphi_s(X_i)}{2^N},$$

the fraction of pieces with significand less than or equal to *s*. Then

$$Iim_{N\to\infty} \mathbb{E}[P_N(s)] = \log_{10} s.$$

$$im_{N\to\infty}\operatorname{Var}(P_N(s))=0.$$

Introduction	Preliminaries	Results	Conclusion
		000	

#### **Limiting Behavior of Decompostions**

Theorem (Decomposition model as above)

Fix a continuous density f on [0, 1] st f(x), f(1 - x) satisfy (C1), set

$$P_N(s) := \frac{\sum_{i=1}^{2^N} \varphi_s(X_i)}{2^N},$$

the fraction of pieces with significand less than or equal to *s*. Then

$$Iim_{N\to\infty} \mathbb{E}[P_N(s)] = \log_{10} s.$$

$$im_{N\to\infty}\operatorname{Var}(P_N(s))=0.$$

Amalgamation of processes converges to Benford.

Introduction	Preliminaries	Results	Conclusion
		00	

#### Limiting Behavior of Decompostions

Theorem (Decomposition model as above)

Fix a continuous density f on [0, 1] st f(x), f(1 - x) satisfy (C1), set

$$P_N(s) := \frac{\sum_{i=1}^{2^N} \varphi_s(X_i)}{2^N},$$

the fraction of pieces with significand less than or equal to *s*. Then

$$Iim_{N\to\infty} \mathbb{E}[P_N(s)] = \log_{10} s.$$

$$\lim_{N\to\infty}\operatorname{Var}\left(P_N(s)\right)=0.$$

- Amalgamation of processes converges to Benford.
- May consider a single process (if number stages tend to infinity).

Introduction	Preliminaries	Results	Conclusion
		000	

Convergence of Amalgamated Lengths to Benford

$$\mathbb{E}[\boldsymbol{P}_N(\boldsymbol{s})] = \mathbb{E}\left[\frac{\sum_{i=1}^{2^N}\varphi_s(X_i)}{2^N}\right] = \frac{1}{2^N}\sum_{i=1}^{2^N}\mathbb{E}[\varphi_s(X_i)].$$

Introduction	Preliminaries	Results	Conclusion
		000	

**Convergence of Amalgamated Lengths to Benford** 

$$\mathbb{E}[P_N(s)] = \mathbb{E}\left[\frac{\sum_{i=1}^{2^N}\varphi_s(X_i)}{2^N}\right] = \frac{1}{2^N}\sum_{i=1}^{2^N}\mathbb{E}[\varphi_s(X_i)].$$

**Note:** Dependencies exist in the pieces  $\{X_i\}$ , not the random variables  $K_i$ .

Introduction	Preliminaries	Results	Conclusion
		000	

**Convergence of Amalgamated Lengths to Benford** 

$$\mathbb{E}[\boldsymbol{P}_{N}(\boldsymbol{s})] = \mathbb{E}\left[\frac{\sum_{i=1}^{2^{N}}\varphi_{s}(X_{i})}{2^{N}}\right] = \frac{1}{2^{N}}\sum_{i=1}^{2^{N}}\mathbb{E}[\varphi_{s}(X_{i})].$$

**Note:** Dependencies exist in the pieces  $\{X_i\}$ , not the random variables  $K_i$ .

We can apply JKKKM's theorem, as the Mellin transform at  $1 - \frac{2\pi i \ell}{\log 10}$  is strictly less than 1.

For specific choices of *f*, can obtain precise bounds on the error. Ex:  $\mathbb{E}[\varphi_s(X_i)] - \log_{10} s \ll \frac{1}{2.9^N}$ .

Introduction	Preliminaries	Results	Conclusion
000	000	○○●	

#### Analysis of a Single Process

Dependence greatly complicates analysis here.

**Problem:** Evaluating cross terms  $\mathbb{E}[\varphi_s(X_i)\varphi_s(X_j)]$  for  $i \neq j$ .

Introduction 000	Preliminaries 000	Results ○○●	Conclusion
	Durana		

Analysis of a Single Process

Dependence greatly complicates analysis here.

**Problem:** Evaluating cross terms  $\mathbb{E}[\varphi_s(X_i)\varphi_s(X_j)]$  for  $i \neq j$ .

# Solution:

• Pair  $(X_i, X_j)$  shares *M* factors, split at M + 1. Remaining elements in the product are independent, so we can simplify our integral.

Introduction	Preliminaries	Results	Conclusion
000	000	○○●	00
Analysis of a	Single Process		

Dependence greatly complicates analysis here.

**Problem:** Evaluating cross terms  $\mathbb{E}[\varphi_s(X_i)\varphi_s(X_j)]$  for  $i \neq j$ .

# Solution:

- Pair (X<sub>i</sub>, X<sub>j</sub>) shares *M* factors, split at *M* + 1.
  Remaining elements in the product are independent, so we can simplify our integral.
- Study integral as a function of the significand of first M + 1 variables, which we can bound by JKKKM.

Introduction	Preliminaries	Results	Conclusion
000	000	○○●	

Analysis of a Single Process

Dependence greatly complicates analysis here.

**Problem:** Evaluating cross terms  $\mathbb{E}[\varphi_s(X_i)\varphi_s(X_j)]$  for  $i \neq j$ .

# Solution:

- Pair (X<sub>i</sub>, X<sub>j</sub>) shares *M* factors, split at *M* + 1.
  Remaining elements in the product are independent, so we can simplify our integral.
- Study integral as a function of the significand of first M + 1 variables, which we can bound by JKKKM.
- Remaining task is to count for each of the 2<sup>N</sup> choices of *i* and for 1 ≤ n ≤ N, how many choices of X<sub>j</sub> have n factors not in common with X<sub>j</sub>.

Introduction	Preliminaries	Results	Conclusion
000		○○●	00

#### Analysis of a Single Process

Dependence greatly complicates analysis here.

**Problem:** Evaluating cross terms  $\mathbb{E}[\varphi_s(X_i)\varphi_s(X_j)]$  for  $i \neq j$ .

# Solution:

- Pair (X<sub>i</sub>, X<sub>j</sub>) shares *M* factors, split at *M* + 1.
  Remaining elements in the product are independent, so we can simplify our integral.
- Study integral as a function of the significand of first M + 1 variables, which we can bound by JKKKM.
- Remaining task is to count for each of the 2<sup>N</sup> choices of *i* and for 1 ≤ n ≤ N, how many choices of X<sub>j</sub> have n factors not in common with X<sub>j</sub>.
- Resulting sum bounds variance above and goes to 0, thus lim<sub>N→∞</sub> Var (P<sub>N</sub>(s)) = 0.

Introduction	Preliminaries	Results	Conclusion
000	000	000	●○
Summary of Re	sults		

- Decay model: Stick decomposes in discrete stages, cut determined by continuous density function.
  - Pieces must sum to original stick length, thus they are dependent.
  - Allowable densities obey weak condition C1.

Introduction	Preliminaries	Results	Conclusion
000		000	●○
Summary of Re	esults		

- Decay model: Stick decomposes in discrete stages, cut determined by continuous density function.
  - Pieces must sum to original stick length, thus they are dependent.
  - Allowable densities obey weak condition C1.
- Expected lengths of pieces from amalgamation of processes converges to Benford distribution.
  - Key observation: dependencies exist among piece lengths, not densities.

Introduction	Preliminaries	Results	Conclusion
000	000	000	●○
Summary of Re	sults		

- Decay model: Stick decomposes in discrete stages, cut determined by continuous density function.
  - Pieces must sum to original stick length, thus they are dependent.
  - Allowable densities obey weak condition C1.
- Expected lengths of pieces from amalgamation of processes converges to Benford distribution.
  - Key observation: dependencies exist among piece lengths, not densities.
- Variance in piece length distribution goes to zero for a single process. Calculation complicated by dependencies:
  - Study expectation cross terms as integrals over their independent factors.

Introduction	Preliminaries	Results	Conclusion
000	000	000	○●
Thanks to:			

- AMS / MAA
- Williams College SMALL 2011
- National Science Foundation
- Thealexa Becker, Professor Steven J. Miller, Ryan Ronan, Professor Frederick W. Strauch