Introduction
oooooo

General Theory
ooooooo

Applications
oooo

3x + 1
oooo

Copulas
oooo

Conclusions

## Benford's law, or: Why the IRS cares about number theory!

Steven J Miller (Smith / Mount Holyoke / Williams Colleges)

sjm1@williams.edu
http://www.williams.edu/Mathematics/sjmiller/

Smith College, November 15, 2011

## Introduction

### Interesting Question

For a nice data set, such as the Fibonacci numbers, stock prices, street addresses of Smith professors, ..., what percent of the leading digits are 1?

Plausible answers:

## Introduction

### Interesting Question

For a nice data set, such as the Fibonacci numbers, stock prices, street addresses of Smith professors, ..., what percent of the leading digits are 1?

Plausible answers: 10%

## Introduction

### Interesting Question

For a nice data set, such as the Fibonacci numbers, stock prices, street addresses of Smith professors, ..., what percent of the leading digits are 1?
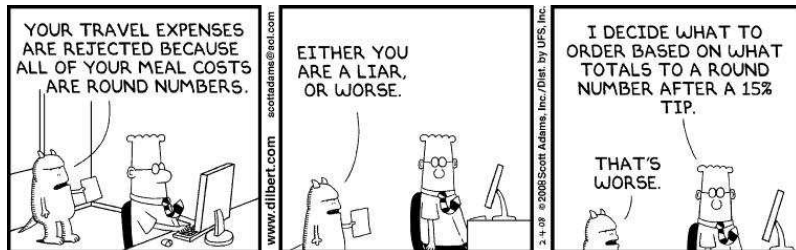
Plausible answers: 10%, 11%

## Introduction

### Interesting Question

For a nice data set, such as the Fibonacci numbers, stock prices, street addresses of Smith professors, ..., what percent of the leading digits are 1?

Plausible answers: 10%, 11%, about 30%.

## Summary

- State Benford's Law.

- Discuss examples and applications.

- Sketch proofs.

- Describe open problems.

**Caveats!**

- A math test indicating fraud is *not* proof of fraud: unlikely events, alternate reasons.

**Caveats!**

- A math test indicating fraud is *not* proof of fraud: unlikely events, alternate reasons.

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

### Statement

For many data sets, probability of observing a first digit of $d$ base $B$ is $\log_B\left(\frac{d+1}{d}\right)$; base 10 about 30% are 1s.

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

### Statement

For many data sets, probability of observing a first digit of $d$ base $B$ is $\log_B\left(\frac{d+1}{d}\right)$; base 10 about 30% are 1s.

- Not all data sets satisfy Benford's Law.

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

### Statement

For many data sets, probability of observing a first digit of $d$ base $B$ is $\log_B\left(\frac{d+1}{d}\right)$; base 10 about 30% are 1s.

- Not all data sets satisfy Benford's Law.
  ◇ Long street $[1, L]$: $L = 199$ versus $L = 999$.

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

### Statement

For many data sets, probability of observing a first digit of $d$ base $B$ is $\log_B\left(\frac{d+1}{d}\right)$; base 10 about 30% are 1s.

- Not all data sets satisfy Benford's Law.
  - ◇ Long street $[1, L]$: $L = 199$ versus $L = 999$.
  - ◇ Oscillates between $1/9$ and $5/9$ with first digit 1.

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

### Statement

For many data sets, probability of observing a first digit of $d$ base $B$ is $\log_B\left(\frac{d+1}{d}\right)$; base 10 about 30% are 1s.

- Not all data sets satisfy Benford's Law.
  ◇ Long street $[1, L]$: $L = 199$ versus $L = 999$.
  ◇ Oscillates between $1/9$ and $5/9$ with first digit 1.
  ◇ Many streets of different sizes: close to Benford.

**Examples**

- recurrence relations
- special functions (such as $n!$)
- iterates of power, exponential, rational maps
- products of random variables
- $L$-functions, characteristic polynomials
- iterates of the $3x + 1$ map
- differences of order statistics
- hydrology and financial data
- many hierarchical Bayesian models

**Applications**

- analyzing round-off errors

- determining the optimal way to store numbers

- detecting tax and image fraud, and data integrity

# General Theory

**Mantissas (or Significands)**

$x \bmod 1$ means the fractional part of $x$: $x - \lfloor x \rfloor$.

**Mantissas (or Significands)**

$x$ mod 1 means the fractional part of $x$: $x - \lfloor x \rfloor$.

Mantissa: $x = M_{10}(x) \cdot 10^k$, $k$ integer.

**Mantissas (or Significands)**

$x \bmod 1$ means the fractional part of $x$: $x - \lfloor x \rfloor$.

Mantissa: $x = M_{10}(x) \cdot 10^k$, $k$ integer.

$M_{10}(x) = M_{10}(\widetilde{x})$ if and only if $x$ and $\widetilde{x}$ have the same leading digits.

**Mantissas (or Significands)**

$x$ mod 1 means the fractional part of $x$: $x - \lfloor x \rfloor$.

Mantissa: $x = M_{10}(x) \cdot 10^k$, $k$ integer.

$M_{10}(x) = M_{10}(\widetilde{x})$ if and only if $x$ and $\widetilde{x}$ have the same leading digits.

Key observation: $\log_{10}(x) = \log_{10}(\widetilde{x})$ mod 1 if and only if $x$ and $\widetilde{x}$ have the same leading digits. Thus often study $y = \log_{10} x$.

**Equidistribution and Benford's Law**

## Equidistribution

$\{y_n\}_{n=1}^{\infty}$ is equidistributed modulo 1 if probability $y_n \bmod 1 \in [a, b]$ tends to $b - a$:

$$\frac{\#\{n \leq N : y_n \bmod 1 \in [a, b]\}}{N} \to b - a.$$

**Equidistribution and Benford's Law**

## Equidistribution

$\{y_n\}_{n=1}^{\infty}$ is equidistributed modulo 1 if probability $y_n \bmod 1 \in [a, b]$ tends to $b - a$:

$$\frac{\#\{n \leq N : y_n \bmod 1 \in [a, b]\}}{N} \to b - a.$$

- Thm: $\beta \notin \mathbb{Q}$, $n\beta$ is equidistributed mod 1.

**Equidistribution and Benford's Law**

### Equidistribution

$\{y_n\}_{n=1}^{\infty}$ is equidistributed modulo 1 if probability $y_n \bmod 1 \in [a, b]$ tends to $b - a$:

$$\frac{\#\{n \leq N : y_n \bmod 1 \in [a, b]\}}{N} \to b - a.$$

- Thm: $\beta \notin \mathbb{Q}$, $n\beta$ is equidistributed mod 1.

- Examples: $\log_{10} 2, \log_{10}\left(\frac{1+\sqrt{5}}{2}\right) \notin \mathbb{Q}$.

**Equidistribution and Benford's Law**

## Equidistribution

$\{y_n\}_{n=1}^{\infty}$ is equidistributed modulo 1 if probability $y_n \bmod 1 \in [a, b]$ tends to $b - a$:

$$\frac{\#\{n \le N : y_n \bmod 1 \in [a, b]\}}{N} \to b - a.$$

- Thm: $\beta \notin \mathbb{Q}$, $n\beta$ is equidistributed mod 1.

- Examples: $\log_{10} 2, \log_{10}\left(\frac{1+\sqrt{5}}{2}\right) \notin \mathbb{Q}$.
  *Proof:* if rational: $2 = 10^{p/q}$.

**Equidistribution and Benford's Law**

### Equidistribution

$\{y_n\}_{n=1}^{\infty}$ is equidistributed modulo 1 if probability $y_n \bmod 1 \in [a, b]$ tends to $b - a$:

$$\frac{\#\{n \leq N : y_n \bmod 1 \in [a, b]\}}{N} \to b - a.$$

- Thm: $\beta \notin \mathbb{Q}$, $n\beta$ is equidistributed mod 1.

- Examples: $\log_{10} 2, \log_{10}\left(\frac{1+\sqrt{5}}{2}\right) \notin \mathbb{Q}$.
  *Proof:* if rational: $2 = 10^{p/q}$.
  Thus $2^q = 10^p$ or $2^{q-p} = 5^p$, impossible.

## Example of Equidistribution: $n\sqrt{\pi}$ mod 1



$n\sqrt{\pi}$ mod 1 for $n \leq 10$

## Example of Equidistribution: $n\sqrt{\pi}$ mod 1



$n\sqrt{\pi}$ mod 1 for $n \le 100$

## Example of Equidistribution: $n\sqrt{\pi} \bmod 1$



$n\sqrt{\pi} \bmod 1$ for $n \leq 1000$

## Example of Equidistribution: $n\sqrt{\pi}$ mod 1



$n\sqrt{\pi}$ mod 1 for $n \leq 10,000$

**Logarithms and Benford's Law**

### Fundamental Equivalence

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

**Logarithms and Benford's Law**

### Fundamental Equivalence

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

**Logarithms and Benford's Law**

### Fundamental Equivalence

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

Introduction
oooooo

General Theory
oooo●oo

Applications
oooo

$3x + 1$
oooo

Copulas
oooo

Conclusions

## Examples

- $2^n$ is Benford base 10 as $\log_{10} 2 \notin \mathbb{Q}$.

## Examples

- $2^n$ is Benford base 10 as $\log_{10} 2 \notin \mathbb{Q}$.

- Fibonacci numbers are Benford base 10.

**Examples**

- $2^n$ is Benford base 10 as $\log_{10} 2 \notin \mathbb{Q}$.

- Fibonacci numbers are Benford base 10.
  $a_{n+1} = a_n + a_{n-1}$.

## Examples

- $2^n$ is Benford base 10 as $\log_{10} 2 \notin \mathbb{Q}$.

- Fibonacci numbers are Benford base 10.
  $a_{n+1} = a_n + a_{n-1}$.
  Guess $a_n = r^n$: $r^{n+1} = r^n + r^{n-1}$ or $r^2 = r + 1$.

**Examples**

- $2^n$ is Benford base 10 as $\log_{10} 2 \notin \mathbb{Q}$.

- Fibonacci numbers are Benford base 10.

  $a_{n+1} = a_n + a_{n-1}$.

  Guess $a_n = r^n$: $r^{n+1} = r^n + r^{n-1}$ or $r^2 = r + 1$.

  Roots $r = (1 \pm \sqrt{5})/2$.

**Examples**

- $2^n$ is Benford base 10 as $\log_{10} 2 \notin \mathbb{Q}$.

- Fibonacci numbers are Benford base 10.

  $a_{n+1} = a_n + a_{n-1}$.
  Guess $a_n = r^n$: $r^{n+1} = r^n + r^{n-1}$ or $r^2 = r + 1$.
  Roots $r = (1 \pm \sqrt{5})/2$.
  General solution: $a_n = c_1 r_1^n + c_2 r_2^n$.

## Examples

- $2^n$ is Benford base 10 as $\log_{10} 2 \notin \mathbb{Q}$.

- Fibonacci numbers are Benford base 10.

  $a_{n+1} = a_n + a_{n-1}$.

  Guess $a_n = r^n$: $r^{n+1} = r^n + r^{n-1}$ or $r^2 = r + 1$.

  Roots $r = (1 \pm \sqrt{5})/2$.

  General solution: $a_n = c_1 r_1^n + c_2 r_2^n$.

  Binet: $a_n = \frac{1}{\sqrt{5}} \left( \frac{1+\sqrt{5}}{2} \right)^n - \frac{1}{\sqrt{5}} \left( \frac{1-\sqrt{5}}{2} \right)^n$.

**Examples**

- $2^n$ is Benford base 10 as $\log_{10} 2 \notin \mathbb{Q}$.

- Fibonacci numbers are Benford base 10.

  $a_{n+1} = a_n + a_{n-1}$.

  Guess $a_n = r^n$: $r^{n+1} = r^n + r^{n-1}$ or $r^2 = r + 1$.

  Roots $r = (1 \pm \sqrt{5})/2$.

  General solution: $a_n = c_1 r_1^n + c_2 r_2^n$.

  Binet: $a_n = \frac{1}{\sqrt{5}} \left( \frac{1+\sqrt{5}}{2} \right)^n - \frac{1}{\sqrt{5}} \left( \frac{1-\sqrt{5}}{2} \right)^n$.

- Most linear recurrence relations Benford.

**Digits of** $2^n$

First 60 values of $2^n$ (only displaying 30)

| | | | digit | # | Obs Prob | Benf Prob |
|---|---|---|---|---|---|---|
| 1 | 1024 | 1048576 | 1 | 18 | .300 | .301 |
| 2 | 2048 | 2097152 | 2 | 12 | .200 | .176 |
| 4 | 4096 | 4194304 | 3 | 6 | .100 | .125 |
| 8 | 8192 | 8388608 | 4 | 6 | .100 | .097 |
| 16 | 16384 | 16777216 | 5 | 6 | .100 | .079 |
| 32 | 32768 | 33554432 | 6 | 4 | .067 | .067 |
| 64 | 65536 | 67108864 | 7 | 2 | .033 | .058 |
| 128 | 131072 | 134217728 | 8 | 5 | .083 | .051 |
| 256 | 262144 | 268435456 | 9 | 1 | .017 | .046 |
| 512 | 524288 | 536870912 | | | | |

**Digits of** $2^n$

First 60 values of $2^n$ (only displaying 30)

| 1 | 1024 | 1048576 | digit | # | Obs Prob | Benf Prob |
|---|------|---------|-------|---|----------|-----------|
| 2 | 2048 | 2097152 | 1 | 18 | .300 | .301 |
| 4 | 4096 | 4194304 | 2 | 12 | .200 | .176 |
| 8 | 8192 | 8388608 | 3 | 6 | .100 | .125 |
| 16 | 16384 | 16777216 | 4 | 6 | .100 | .097 |
| 32 | 32768 | 33554432 | 5 | 6 | .100 | .079 |
| 64 | 65536 | 67108864 | 6 | 4 | .067 | .067 |
| 128 | 131072 | 134217728 | 7 | 2 | .033 | .058 |
| 256 | 262144 | 268435456 | 8 | 5 | .083 | .051 |
| 512 | 524288 | 536870912 | 9 | 1 | .017 | .046 |

**Digits of $2^n$**

First 60 values of $2^n$ (only displaying 30): $2^{10} = 1024 \approx 10^3$.

| 1 | 1024 | 1048576 | digit | # | Obs Prob | Benf Prob |
|---|------|---------|-------|---|----------|-----------|
| 2 | 2048 | 2097152 | 1 | 18 | .300 | .301 |
| 4 | 4096 | 4194304 | 2 | 12 | .200 | .176 |
| 8 | 8192 | 8388608 | 3 | 6 | .100 | .125 |
| 16 | 16384 | 16777216 | 4 | 6 | .100 | .097 |
| 32 | 32768 | 33554432 | 5 | 6 | .100 | .079 |
| 64 | 65536 | 67108864 | 6 | 4 | .067 | .067 |
| 128 | 131072 | 134217728 | 7 | 2 | .033 | .058 |
| 256 | 262144 | 268435456 | 8 | 5 | .083 | .051 |
| 512 | 524288 | 536870912 | 9 | 1 | .017 | .046 |

**Logarithms and Benford's Law**

$\chi^2$ values for $\alpha^n$, $1 \leq n \leq N$ (5% 15.5).

| $N$ | $\chi^2(\gamma)$ | $\chi^2(e)$ | $\chi^2(\pi)$ |
|---|---|---|---|
| 100 | 0.72 | 0.30 | 46.65 |
| 200 | 0.24 | 0.30 | 8.58 |
| 400 | 0.14 | 0.10 | 10.55 |
| 500 | 0.08 | 0.07 | 2.69 |
| 700 | 0.19 | 0.04 | 0.05 |
| 800 | 0.04 | 0.03 | 6.19 |
| 900 | 0.09 | 0.09 | 1.71 |
| 1000 | 0.02 | 0.06 | 2.90 |

**Logarithms and Benford's Law: Base 10**

$\log_{10}(\chi^2)$ vs $N$ for $\pi^n$ (red) and $e^n$ (blue), $n \in \{1, \ldots, N\}$. Note $\pi^{175} \approx 1.0028 \cdot 10^{87}$, (5% and 8 d.f., $\log_{10}(\chi^2) \approx .44$).

# Applications

## Applications for the IRS: Detecting Fraud

## Applications for the IRS: Detecting Fraud

**Detecting Fraud**

## Bank Fraud

- Audit of a bank revealed huge spike of numbers

**Detecting Fraud**

## Bank Fraud

- Audit of a bank revealed huge spike of numbers starting with 4

**Detecting Fraud**

## Bank Fraud

- Audit of a bank revealed huge spike of numbers starting with 48 and 49

**Detecting Fraud**

## Bank Fraud

- Audit of a bank revealed huge spike of numbers starting with 48 and 49, most due to one person.

**Detecting Fraud**

## Bank Fraud

- Audit of a bank revealed huge spike of numbers starting with 48 and 49, most due to one person.

- Write-off limit of \$5,000. Officer had friends applying for credit cards, ran up balances just under \$5,000 then he would write the debts off.

**Data Integrity: Stream Flow Statistics: 130 years, 457,440 records**

**Election Fraud: Iran 2009**

Numerous protests/complaints over Iran's 2009 elections.

Lot of analysis; data moderately suspicious:

- First and second leading digits;
- Last two digits (should almost be uniform);
- Last two digits differing by at least 2.

Warning: enough tests, even if nothing wrong will find a suspicious result (but when all tests are on the boundary...).

The $3x + 1$ Problem
and
Benford's Law

## $3x + 1$ **Problem**

- Kakutani (conspiracy), Erdös (not ready).

- $x$ odd, $T(x) = \frac{3x+1}{2^k}$, $2^k || 3x + 1$.

- Conjecture: for some $n = n(x)$, $T^n(x) = 1$.

- $7 \rightarrow_1 11 \rightarrow_1 17 \rightarrow_2 13 \rightarrow_3 5 \rightarrow_4 1 \rightarrow_2 1$

## 3x + 1 **and Benford**

**Theorem (Kontorovich and M–, 2005)**

As $m \to \infty$, $x_m/(3/4)^m x_0$ is Benford.

**Theorem (Lagarias-Soundararajan 2006)**

$X \geq 2^N$, for all but at most $c(B)N^{-1/36}X$ initial seeds the distribution of the first $N$ iterates of the $3x + 1$ map are within $2N^{-1/36}$ of the Benford probabilities.

**Sketch of the proof**

- Failed Proof: lattices, bad errors.

- CLT: $(S_m - 2m)/\sqrt{2m} \to N(0, 1)$:

$$\mathbb{P}\left(S_m - 2m = k\right) = \frac{\eta(k/\sqrt{m})}{\sqrt{m}} + O\left(\frac{1}{g(m)\sqrt{m}}\right).$$

- Quantified Equidistribution: $I_\ell = \{\ell M, \ldots, (\ell + 1)M - 1\}$,
  $M = m^c$, $c < 1/2$
  $k_1, k_2 \in I_\ell$: $\left|\eta\left(\frac{k_1}{\sqrt{m}}\right) - \eta\left(\frac{k_2}{\sqrt{m}}\right)\right|$ small
  $C = \log_B 2$ of irrationality type $\kappa < \infty$:

  $$\#\{k \in I_\ell : \overline{kC} \in [a, b]\} = M(b - a) + O(M^{1+\epsilon-1/\kappa}).$$

## Sketch of the proof

- Failed Proof: lattices, bad errors.

- CLT: $(S_m - 2m)/\sqrt{2m} \to N(0, 1)$:

$$\mathbb{P}\left(S_m - 2m = k\right) = \frac{\eta(k/\sqrt{m})}{\sqrt{m}} + O\left(\frac{1}{g(m)\sqrt{m}}\right).$$

- Quantified Equidistribution: $I_\ell = \{\ell M, \ldots, (\ell + 1)M - 1\}$,
  $M = m^c$, $c < 1/2$
  $k_1, k_2 \in I_\ell$: $\left|\eta\left(\frac{k_1}{\sqrt{m}}\right) - \eta\left(\frac{k_2}{\sqrt{m}}\right)\right|$ small
  $C = \log_B 2$ of irrationality type $\kappa < 1.2 \cdot 10^{602} < \infty$:

$$\#\{k \in I_\ell : \overline{kC} \in [a, b]\} = M(b - a) + O(M^{1+\epsilon-1/\kappa}).$$

## $3x + 1$ **Data: random 10,000 digit number,** $2^k||3x + 1$

80,514 iterations $((4/3)^n = a_0$ predicts 80,319);
$\chi^2 = 13.5$ (5% 15.5).

| Digit | Number | Observed | Benford |
|-------|--------|----------|---------|
| 1 | 24251 | 0.301 | 0.301 |
| 2 | 14156 | 0.176 | 0.176 |
| 3 | 10227 | 0.127 | 0.125 |
| 4 | 7931 | 0.099 | 0.097 |
| 5 | 6359 | 0.079 | 0.079 |
| 6 | 5372 | 0.067 | 0.067 |
| 7 | 4476 | 0.056 | 0.058 |
| 8 | 4092 | 0.051 | 0.051 |
| 9 | 3650 | 0.045 | 0.046 |

3x + 1 **Data: random 10,000 digit number,** $2|3x + 1$

241,344 iterations, $\chi^2 = 11.4$ (5% 15.5).

| Digit | Number | Observed | Benford |
|-------|--------|----------|---------|
| 1 | 72924 | 0.302 | 0.301 |
| 2 | 42357 | 0.176 | 0.176 |
| 3 | 30201 | 0.125 | 0.125 |
| 4 | 23507 | 0.097 | 0.097 |
| 5 | 18928 | 0.078 | 0.079 |
| 6 | 16296 | 0.068 | 0.067 |
| 7 | 13702 | 0.057 | 0.058 |
| 8 | 12356 | 0.051 | 0.051 |
| 9 | 11073 | 0.046 | 0.046 |

Introduction
oooooo

General Theory
ooooooo

Applications
oooo

$3x + 1$
oooo

Copulas
oooo

Conclusions

Copulas and Benford's Law
(joint with Thealexa Becker '13)

### Definition of Copulas

Copula: A form of joint CDF between multiple variables with given uniform marginals on the d-dimensional unit cube.

### Sklar's Theorem

Let X and Y be random variables with joint distribution function H and marginal distribution fucntions F and G respectively. There exists a copula, C, such that

$$\text{for all } x, y \in \mathbb{R}, \ \ H(x, y) = C(F(x), G(y)).$$

Introduction
000000

General Theory
0000000

Applications
0000

$3x + 1$
0000

**Copulas**
0●00

Conclusions

### Archimedean Copulas

A commonly used / studied family of copulas is of the form

$$C(x, y) = \phi^{-1}(\phi(x) + \phi(y))$$

where $\phi$ is the generator and $\phi^{-1}$ is the inverse generator of the copula.

Investigating the Benfordness of the product of random variables arising from copulas.

Clayton Copula: $C(x, y) = (x^{-\theta} + y^{-\theta} - 1)^{-1/\theta}$.

PDF (bivariate): $\theta(\theta^{-1} + 1)(xy)^{-\theta-1}(x^{-\theta} + y^{-\theta} - 1)^{-2-1/\theta}$.

PDF (general case):
$\theta^{n-1}\frac{\Gamma(n+\theta^{-1})}{\Gamma(1+\theta^{-1})}(x_1 \cdots x_n)^{-\theta-1}(x_1^{-\theta} + \cdots + x_n^{-\theta} - 1)^{-n-1/\theta}$.

**Results**

- Early data and chi-square tests of multivariate copulas suggest Benford behavior of the products of copulas.
- Proof strategy includes the integration of the PDF over the region in which the product has first digit $d$ using Poisson summation:

$$\int_0^1 \cdots \int_0^1 \sum_k \widehat{\phi}_{\log_{10}(x_1 \cdots x_n)}(k) p(x_1, \ldots, x_n) dx_1 \cdots dx_n,$$

where

$$\phi_a(u) = \chi_{[1,2)}(10^{u+a}) = \begin{cases} 1 & \text{if } 10^{u+a} \in [1,2) \\ 0 & \text{otherwise.} \end{cases}$$

# Conclusions

**Conclusions and Future Investigations**

- Many different systems are Benford.

- Ingredients of proofs (logarithms, equidistribution).

- Applications to fraud detection / data integrity.

- Future work:
  ◇ Study digits of other systems.
  ◇ Develop more sophisticated tests for fraud.