

# Anomaly Detection in Network Flows

## Benford's Law

**Bianca I. Colón Rosado**

`bianca.colon1@upr.edu`

Computer Science Department

University of Puerto Rico - Río Piedras

Advisor:

**Humberto Ortiz Zuazaga**

`humberto@hpcf.upr.edu`

Computer Science Department

University of Puerto Rico - Río Piedras

May 2016

### Abstract

The Benford's Law, is a phenomenological law about the frequency distribution of leading digits in many real-life sets of numerical data. Our assumption is that Benford's law applies to the TCP flow inter-arrival times as well, and therefore, simpler approaches for computer network traffic analysis can be applied, specifically with the aim of fault and intrusion detection. To prove it we are using the Numenta Anomaly Benchmark, that contains labeled data with anomalies. We are classifying as flow anomalies the peaks that deviate from the baseline when we apply the Benford's Law. An important advantage of this method is that malware cannot easily adapt their communication pattern to conform to the logarithmic distribution of first digits. If the Benford's law works in TCP flows, we will be close to find a general method to detect anomalies in network traffic.

# 1 Description

Last semester we implemented the Benford’s law [3]. The Benford’s law also called the First-Digit Law, is a phenomenological law about the frequency distribution of leading digits in many real-life sets of numerical data. That law states that in many naturally occurring collections of numbers the small digits occur disproportionately often as leading significant digits.

We started looking for a general technique for anomaly detection in Network Flows. Our assumption was that Benford’s law applies to the TCP flow inter-arrival times as well, and therefore, simpler approaches for computer network traffic analysis can be applied, specifically with the aim of fault and intrusion detection. We expect that intentional attacks alter the first digit distribution of the inter-arrival times can simply be detected without the need of packet header inspection.

This semester, we implemented the Evaluating Real-time Anomaly Detection Algorithms – the Numenta Anomaly Benchmark (NAB) [9], in order to compare and evaluate different algorithms for detecting anomalies in streaming data. Using NAB we identify two phases where we want to work.

1. Use NAB with our SiLK flows, to identify anomalies in our data.
2. Run our Benford’s Algorithm with their flows data, that is already labeled with real anomalies. With this approach we can prove if our algorithm work with TCP flow inter-arrival times.

## 2 Background

### 2.1 Benford’s Law

According to Benford’s law of anomalous numbers the frequency of the digit  $d$ , appearing as the first significant digit in a collection of numbers, is no uniform as expected intuitively, rather it follows closely the logarithmic relation:

$$P_d = \log_{10} \frac{d+1}{d}, \sum_{d=1}^9 P(d)$$

Sets which obey the law the number 1 would appear as the most significant digit about 30% of the time while larger digits would occur in that position less frequently: 9 would appear less than 5% of the time. If all digits were

distributed uniformly, they would each occur about 11.1% of the time. See Figure 1.

### 3 Methodology

In order to use NAB with our SiLK flows, to identify anomalies in our data.

1. We start the installation process of NAB in HULK. After numerous attempts of trying, we cannot install it in HULK.
2. Then we start installing NAB locally in my computer. We install it, but only work's once before it crash. After numerous attempts to make it work we cannot make it work.

We documented the followed steps in GitHub [10].

Then we decide to move to the second phase that is run our Bendford's Algorithm with the NAB flows data, that is already labeled with real anomalies.

1. We adapted our Benford's program to receive the format provided by NAB. We include the program in the Appendix.
2. We graphed the results in Plotly [8].

### 4 Results

We start comparing the expected curve provided by Benford's Law (Figure 1) with the results obtained from the data labeled as no anomalies (Figure 2).

We notice some differences that should not be there assuming the data don't have anomalies. To find more information about this differences we search the deviations from the Benford theoretical curve.

In the graph of the Deviations from the Benford Theoretical Curve, (Figure 3), there's a gap between the days from September 5 and some hours of September 9 in the provided data. We didn't notice that until we make this graph.

Also, looking the graph we need to find information to define the real anomalies. We want to prognosticate how close or far the deviation need to be from zero to be counted as a TP, FP, FN, TN.

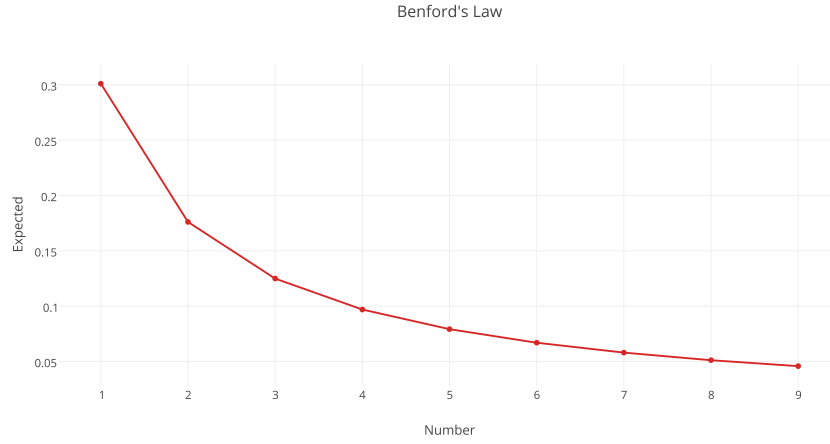


Figure 1: Expected curve

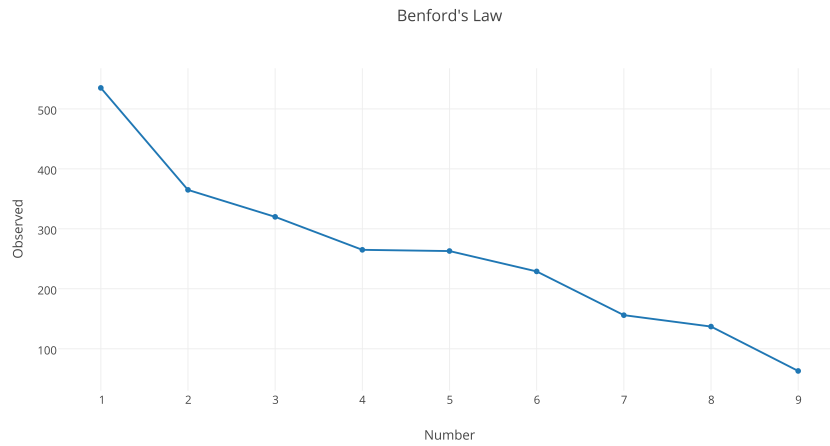


Figure 2: Observed curve

## 5 Conclusion

In order to prove that the Benford's Law can detect anomalies in Network Flows we need to continue validating the algorithm with labeled data. NAB provide a lot of files with data, but we don't get enough time to inspect all the files. Also we need to define when we classify the results of the Deviations

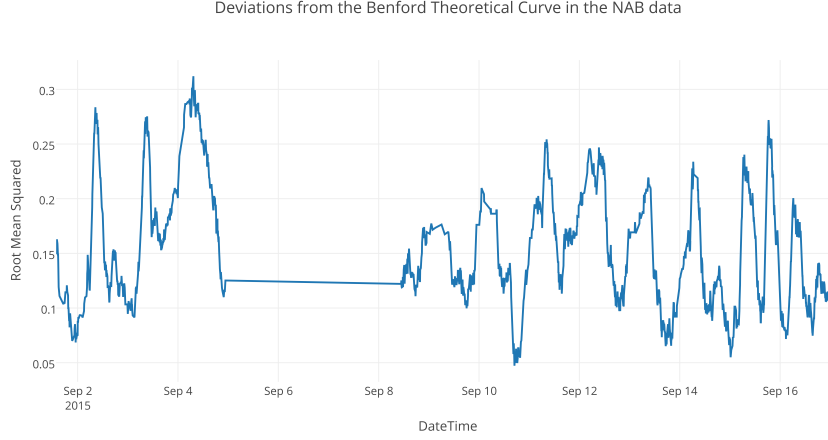


Figure 3: Deviations from the Benford Theoretical Curve in the NAB data

from the Benford theoretical curve as a real anomaly.

An important advantage of this method is that malware cannot easily adapt their communication pattern to conform to the logarithmic distribution of first digits.

## 6 Future Work

We want to analyze more the utility of this law, we want to compare the results of the Deviations from the Benford theoretical curve with the octets in the same flows.

Also, we will explore new approaches to find new techniques. Implement these techniques for anomaly detection to our collection of flows from UPR's network, and compare results with the results of current techniques. If those techniques are effective we can use it in real time flow collection and build an alerting system to notify the anomalies as soon as they are detected.

## Acknowledgement

This work is supported by the scholarship Academics and Training for the Advancement of Cybersecurity Knowledge in Puerto Rico (ATAACK-PR) supported by the National Science Foundation under Grant No. DUE-1438838.

## References

- [1] Bianca Colón-Rosado, Humberto Ortiz-Zuazaga. Techniques for Anomaly Detection in Network Flows. 2014.  
[http://figshare.com/articles/Techniques\\_for\\_Anomaly\\_Detection\\_in\\_Network\\_Flows/1424475](http://figshare.com/articles/Techniques_for_Anomaly_Detection_in_Network_Flows/1424475)     <http://ccom.uprrp.edu/~humberto/megaprobe/tag/flows.html>
- [2] Iván García, Humberto Ortiz-Zuazaga. Techniques for Anomaly Detection in Network Flows. 2013. <http://ccom.uprrp.edu/~humberto/research/anomaly-detection.pdf>
- [3] Bianca Colón-Rosado, Humberto Ortiz-Zuazaga. Techniques for Anomaly Detection in Network Flows: Benford's Law 2015.
- [4] Bingdong Li, Jeff Springer, George Bebis and Mehmet Hadi Gunes. A survey of network flow applications. *Journal of Network and Computer Applications*, 2013. doi:10.1016/j.jnca.2012.12.020
- [5] Lakhina Anukool, Mark Crovella, and Christophe Diot. Characterization of network-wide anomalies in traffic flows. doi=10.1.1.92.7738, 2004.
- [6] Laleh Arshadi and Amir Hossein Jahangir. Benford's law behavior of internet traffic. *Journal of Network and Computer Applications*, 2013. doi:10.1016/j.jnca.2013.09.007.
- [7] Ron Bandes, Timothy Shimeall, Matt Heckathorn, Sidney Faber (2014) Using SiLK for Network Traffic Analysis. CERT Coordination Center. 2014.  
<https://tools.netsa.cert.org/silk/>  
<http://tools.netsa.cert.org/silk/analysis-handbook.pdf>
- [8] Plotly. <https://plot.ly/>
- [9] NAB. Numenta Anomaly Benchmark  
<http://arxiv.org/pdf/1510.03336v4.pdf>  
<https://github.com/numenta/NAB>
- [10] Documented Steps. [bit.ly/stepsNAB](http://bit.ly/stepsNAB)

## 7 Appendix

```
1 #!/usr/bin/env python
2
3 # Research Investigation : Spring 2015–2016
4 # Techniques for Anomaly Detection using Benford's Law
5 #
6 # Using the data provided by NAB (https://github.com/numenta/NAB)
7 # we will analyze the results of the Benford's Law.
8 #
9 import math
10 from sys import argv
11
12 def significantNumber(num): #recieve 0.0024
13     number = num.split('.') # '0' '0024'
14     aDot = number[1] # '0024'
15     for d in range(0, len(aDot)):
16         if aDot[d] != '0':
17             return aDot[d]
18     else:
19         pass
20
21 def howManyNumbers(listCounter):
22     allNum = [] # List to put all the numbers
23     diNum = [] * 10 # List to put counts of numbers
24
25     for i in range(0, len(listCounter)):
26         allNum.append(int(listCounter[i])) # Change strings
27         for numbers
28
29     totalNum = len(allNum) # Total of numbers
30     for c in range(1,10):
31         diNum.insert(c, 1.0 * allNum.count(c)/totalNum)
32
33     return diNum
```

```

34 def benford(listCounter):
35     allNum = [] # List to put all the numbers
36     diNum = [] * 10 # List to put counts of numbers
37     for i in range(0, len(listCounter)):
38         allNum.append(int(listCounter[i])) # Change strings
           for numbers
39     for c in range(1,10):
40         diNum.insert(c, allNum.count(c))
41     return diNum
42
43 def rms(freq, expected):
44     sum = 0
45     for i in range(len(freq)):
46         sum += (freq[i] - expected[i])**2
47     return math.sqrt(sum)
48
49 ## main
50 if (len(argv) == 4):
51     filename = argv[1]
52     columnNumber = int(argv[2]) - 1
53     windowSize = int(argv[3])
54
55     with open(filename) as f:
56
57         content = f.read().splitlines()
58
59     f.close()
60
61     value = [] # For the values
62     fDigit = [] # For the first significant digit
63
64     # Take the only values of the line
65     # 1.Starts in line 1 because line 0 contain column
       title. (range(1 ... )
66     ## timestamp, value, anomaly_score, raw_score, label, S(t)
       _reward_low_FP_rate, S(t)_reward_low_FN_rate, S(t)
       _standard
67

```



```

68 # 2. The value is positioned in the second column. (.
    split(',') [1] ) columnNumber should be 1.
69 ### 2015-09-01
    13:45:00,3.06,0.0301029996659,1.0,0,0.0,0.0,0.0
70 for i in range(1,len(content)): # 1.
71     if float(content[i].split(',')[columnNumber]) !=
        int(0):
72         value.append(content[i].split(',')[columnNumber])
        # 2.
73     else:
74         pass
75
76 # To that values, just select the first significant
    digit
77 for j in range(0,len(value)):
78     if value[j].isdigit() == True: # For int values
79         if value[j][0].isdigit() == True:
80             fDigit.append(value[j][0])
81         else:
82             print "Alert Anomaly: value = %d" % (value[j])
83         else: # For float values
84             if value[j].split('.')[0] != '0' : # Save the
                left dot number for a number like 12.003
85                 fDigit.append(value[j].split('.')[0][0]) # From
                12.003 only save the 1
86             else: # Save the right dot number
                for a number like 0.0034
87                 fDigit.append(significantNumber(value[j]))
88
89
90 expected = [ math.log(1.0/d+1.0,10) for d in range
    (1,10)]
91
92 #### ONLY BENFORD GRAPH
93 #bVal = benford(fDigit)
94 #for n in range(0,len(bVal)):
95 # print n+1, bVal[n], expected[n]
96 ###

```

```

97
98
99     for i in range(len(fDigit) - windowSize):
100         valuesFD = howManyNumbers(fDigit[i:i+windowSize])
101         print content[i].split(',')[0], rms(valuesFD,
102             expected)
103         # The timestamp is in the first column, that's why
104         content[i].split(',')[0]
105
106 else:
107     print ""
108     print "Need to provide the name of the file that have
109         the data"
110     print "and in what column is the value data"
111     print "Example:"
112     print "timestamp,value,anomaly_score,raw_score,label,
113         S(t)_reward_low_FP_rate"
114     print "The value column is the second"
115     print ""
116     print "Format to write the missing information"
117     print "$ python benford-NAB.py fileName columnNumber
118         windowSize"
119     print "$ python benford-NAB.py numenta_occupancy_6005
120         .csv 2 100"
121     print ""

```