

THE UNIVERSITY OF FRANCHE-COMTÉ

MASTER COMPU<sub>PHYS</sub>

---

## Statistical physics home assignment : Benford's law

---



Chiari EVEN  
M1 Compuphys 2022-2023

*Supervisor :* Jose LAGES

January 15, 2023

**Contents**

<b>I</b>	<b>Introduction</b>	<b>2</b>
<b>II</b>	<b>Question 5</b>	<b>2</b>
<b>III</b>	<b>Question 6</b>	<b>5</b>
<b>IV</b>	<b>Conclusion</b>	<b>8</b>

## I Introduction

In this home assignment we will study the Benford's law define like :

$$P(i) = \log(i+1) - \log(i) = \log\left(\frac{i+1}{i}\right) = \log\left(1 + \frac{1}{i}\right)$$

with  $i \in \{1, \dots, 9\}$

obviously the total probability should be equal to one. In the other basis the law is define like :

$$P(i) = \log_b(i+1) - \log_b(i) = \log_b\left(\frac{i+1}{i}\right) = \log_b\left(1 + \frac{1}{i}\right)$$

with  $i \in$  various digit in other basis

we can easily compute this formula is base 10 with the change of base formula of log :

$$\log_b(M) = \frac{\log_x(M)}{\log_x(b)}$$

## II Question 5

Now we want to check if the Fibonacci numbers follow the Benford's law, in order to do that we have to to the distribution of the first digit for Fibonacci numbers in base 10 and 16:

First i will code the a Fibonacci function to generate a list of them :

```
1 def fibo(n):
2     """Function for fibonnaci numbers
3
4     Args:
5         n (integer): n fibonnaci numbers
6
7     Returns:
8         fibo (list) : list of n fibonnaci numbers
9     """
10    list=[0,1]
11    for i in range(n):
12        list.append(list[-2] +list[-1])
13    return list
```

Now in order to understand hexadecimal number i will code two function one to convert decimal into hexadecimal and the second to convert hexadecimal into decimal :

```
1 def CDH(n):
2     """Function to convert n in base 10 to n in base 16
3         CDH for Convert Decimal to Hexadecimal
4
5     Args:
6         n (integer): number that we want to convert
7
8     Returns:
9         r (str): return n in hexadecimal so in str
10    """
11    r=""
12    d={0:"0",1:"1",2:"2",3:"3",4:"4",5:"5",6:"6",7:"7",8:"8",9:"9",
13        10:"A",11:"B",12:"C",13:"D",14:"E",15:"F"}
14    while True :
15        r += d[n%16]
16        n = n // 16
17        if n == 0 :
18            return r[::-1]
19
20
21 def CHD(n):
22     """Function to convert n in base 16 to n in base 10
23         CHD for Convert Hexadecimal to Decimal
24
```

```

25     Args:
26         n (str): number that we want to convert
27
28     Returns:
29         r (integer): return n in decimal so in integer
30     """
31     d={0:"0",1:"1",2:"2",3:"3",4:"4",5:"5",6:"6",7:"7",8:"8",9:"9",
32        10:"A",11:"B",12:"C",13:"D",14:"E",15:"F"}
33     d= {v: k for k, v in d.items()}
34     r, it = 0, 0
35     n=n[::-1]
36     for i in n:
37         r+= d[i] * 16**it
38         it +=1
39     return r

```

In python hexadecimal number are like  $0x$  — — — but my function generate these number like — — — like that its easier to get the first digit of a hexadecimal number.

To see if the Fibonacci numbers follow the Benford's law i will plot in the same figure the frequency of each digit for Fibonacci numbers and Benford's law.

The distribution of each digit for benford's law are stored with a dictionary and compute like that in base 10 and in base 16 (with the change of base formula) :

```

1 #compute the Benford's law in base 10
2 #create dictionary to have "i_digit" : frequency
3 it=1 #start at one because we don't take in account 0
4 ben_d={"1":0,"2":0,"3":0,"4":0,"5":0,"6":0,"7":0,"8":0,"9":0} #ben_d for Benford's_law_Decimal
5 for i in ben_d.keys():
6     ben_d[i] = np.log10(1+1/it) #add np.log10(1+1/it) for Benford's law in base 10 with it the
7     i_digit
8     it += 1
9
10 #compute the Benford's law in base 16
11 #create dictionary to have "i_digit" : frequency
12 it=1 #start at one because we don't take in account 0
13 ben_h={"1":0,"2":0,"3":0,"4":0,"5":0,"6":0,"7":0,"8":0,"9":0,
14        "A":0,"B":0,"C":0,"D":0,"E":0,"F":0} #ben_h for Benford's_law_Hexadecimal
15 for i in ben_h.keys():
16     ben_h[i] = np.log10(1+1/it)/np.log10(16) #add np.log10(1+1/it)/np.log10(16)
17     i_digit
18     it += 1

```

and for the fibonacci number :

```

1 data=fibo(N)
2 for i in range(1,N):#loop over the number of elements
3     nd=str(data[i]) #convert the number into str to get a list and have the first digit by list[0]
4     in base 10
5     nh=CDH(data[i]) #convert the number into str to get a list and have the first digit by list[0]
6     in base 16
7     fd_fib_d[nd[0]] = fd_fib_d.get(nd[0]) + 1/N #update the values for i digit by 1/N (already
8     normalised) in base 10
9     fd_fib_h[nh[0]] = fd_fib_h.get(nh[0]) + 1/N #update the values for i digit by 1/N (already
10    normalised) in base 16

```

Now we can plot the histogram for the base 10 :

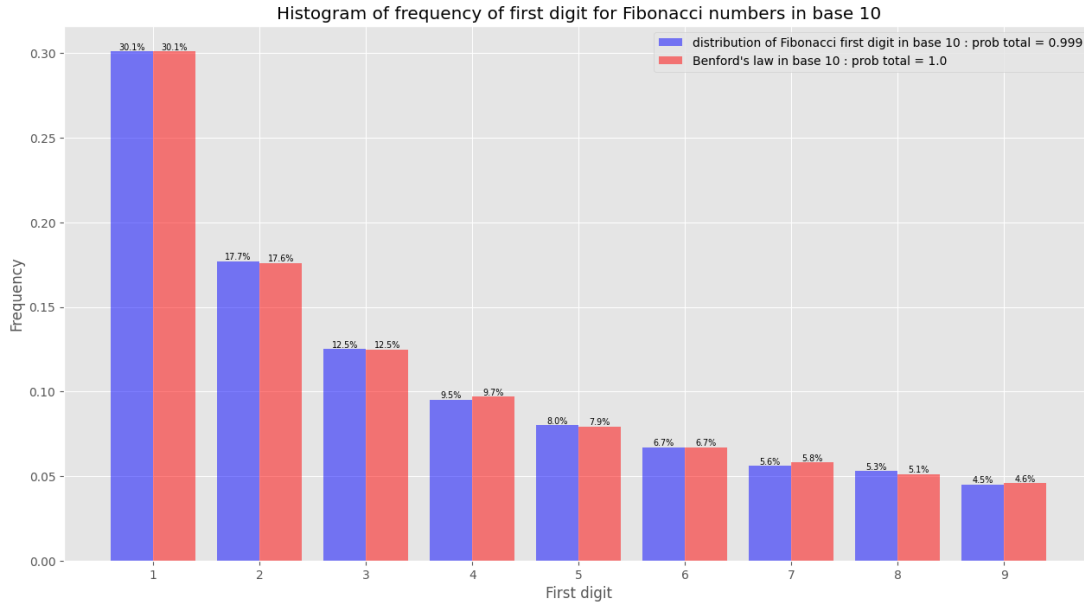


Figure 1: Histogram of distribution of the first digit for Fibonacci numbers in base 10

On this graph we can see clearly the Fibonacci numbers are following the Benford's law in base 10, the total probability should be one for any data set but for the Fibonacci numbers i have 0.999 this come from the precision of a programming language and float number because when i compute the distribution I normalise by the length of the Fibonacci list, here the number I took 1000 numbers to have enough data obviously with a number  $N$  of Fibonacci numbers too few we will not see the correlation between the two law.

We can check the correlation in base 16 :

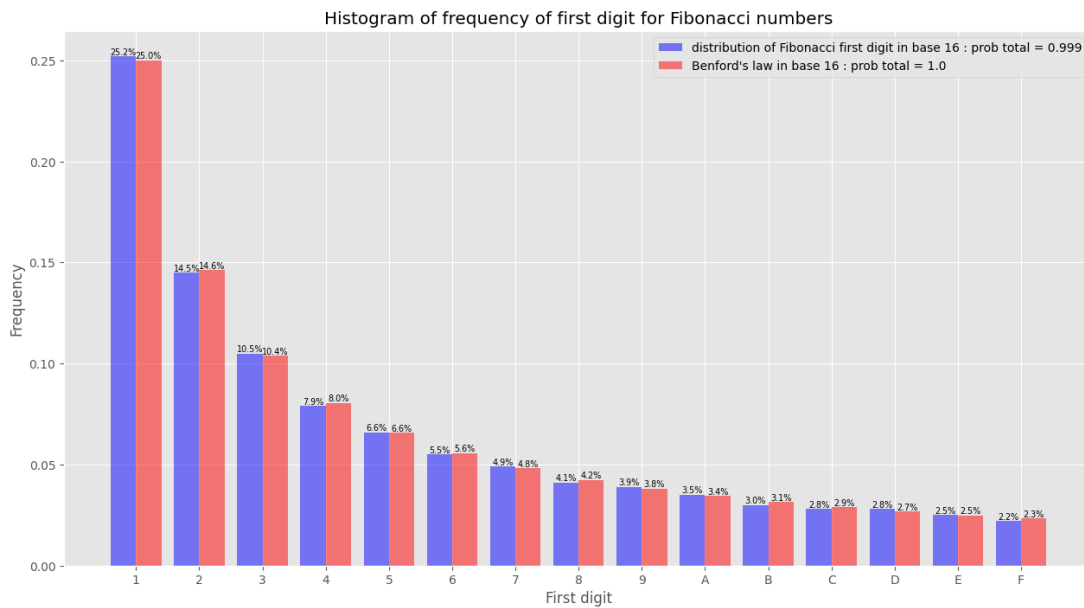


Figure 2: Histogram of distribution of the first digit for Fibonacci numbers in base 16

Even is base 16 we have a good behavior of the distribution of the first digit for Fibonacci numbers, we have the same problem of probability for the Fibonacci numbers but the probability is very close to one so we can say the correlation is correct.

### III Question 6

Now we will check a correlation between two data set and the Benford's law. I am very surprising by this law so i wanted to find data set very tricky so i took the numbers of bar per country and the numbers of bar per city in US. I have found the two data set here :

- bar per city : [click here](#)
- bar per country : [click here](#)

To obtains the number of bar per country I have code that :

```
1 #unpack the data from data_beer.csv
2 path=os.path.join(os.path.dirname(__file__),"data_beer.csv") #path of csv
3 df=pd.read_csv(path) #unpack with pandas
4 nb_bar=df['country'].value_counts() #take the colomun that we want
5 nb_bar=nb_bar.to_dict() #transform dataframe into dictionary
```

To obtains the distribution of the first digit of the number of bar I have proceed like before :

```
1 #compute the distribution of first digit for bar per country in base 10 and 16
2 #create dictionary to have "i_digit" : frequency
3 fd_bar_d={"1":0,"2":0,"3":0,"4":0,"5":0,"6":0,"7":0,"8":0,"9":0} #fd_bar_d for
   First_Digit_Bar_Decimal
4 fd_bar_h={"1":0,"2":0,"3":0,"4":0,"5":0,"6":0,"7":0,"8":0,"9":0,
   "A":0,"B":0,"C":0,"D":0,"E":0,"F":0} #fd_bar_h for First_Digit_Bar_Hexadecimal
5
```

Now we can plot the histogram in base 10 :

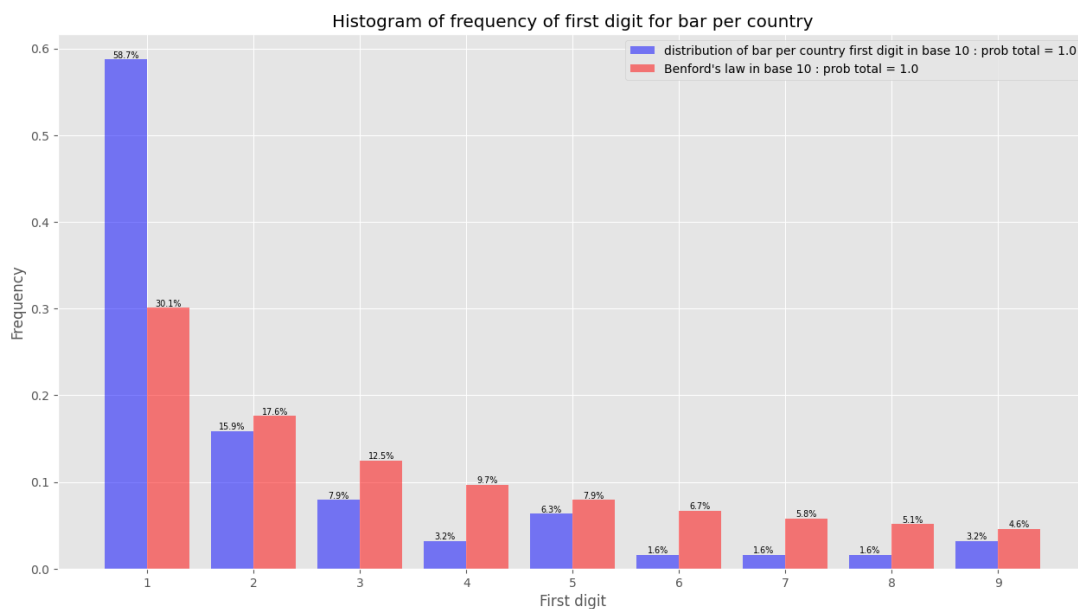


Figure 3: Histogram of distribution of the first digit for bar per country in base 10

I base 10 we can see the number of bar per country has the same tendency of benford's law but the value of distribution are not the same, they are quite close if we forgot the number one.

The histogram in base 16 :

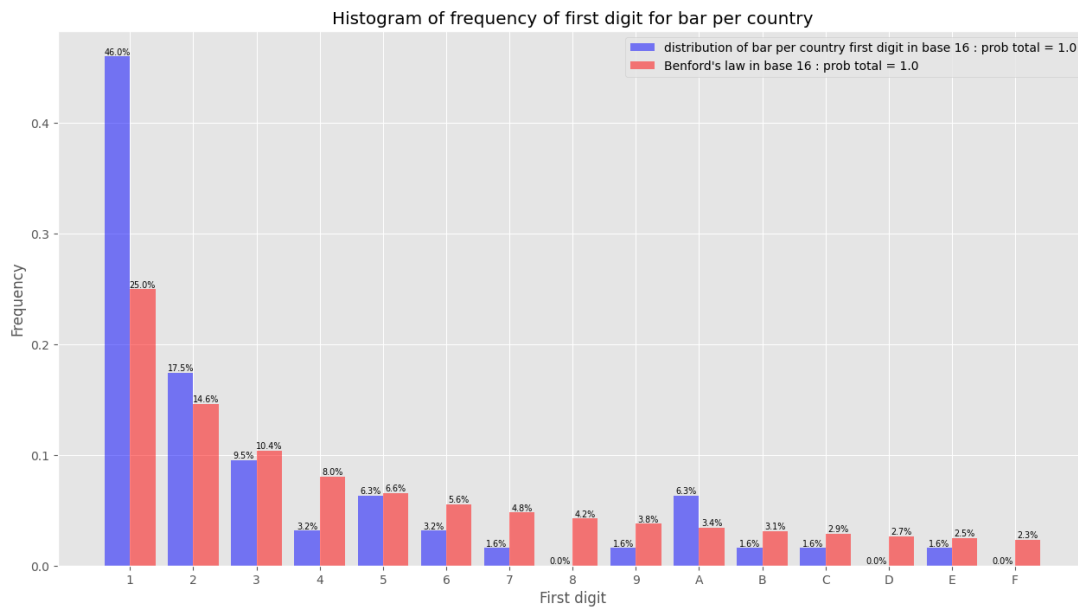


Figure 4: Histogram of distribution of the first digit for bar per country in base 16

In base 16 the tendency between the number of bar per country and the Benford's law is visible but we have like in the previous graph the values quite close for all number but not for 1 and 16.

For the second data set i have choose the city in US, to count i did like for the first data set :

```

1 #compute the distribution of first digit for bar per city in base 10 and 16
2 #create dictionary to have "i_digit" : frequency
3 fd_bar_d={"1":0,"2":0,"3":0,"4":0,"5":0,"6":0,"7":0,"8":0,"9":0} #fd_bar_d for
   First_Digit_Bar_Decimal
4 fd_bar_h={"1":0,"2":0,"3":0,"4":0,"5":0,"6":0,"7":0,"8":0,"9":0,
5           "A":0,"B":0,"C":0,"D":0,"E":0,"F":0} #fd_bar_h for First_Digit_Bar_Hexadecimal
6
7 for i in nb_bar.keys(): #loop over the number of elements
8     nd=str(nb_bar[i]) #convert the number into str to get a list and have the first digit by list[0]
   in base 10
9     nh=CDH(nb_bar[i]) #convert the number into str to get a list and have the first digit by list[0]
   in base 16
10    fd_bar_d[nd[0]] = fd_bar_d.get(nd[0]) + 1/len(nb_bar) #update the values for i digit by 1/N (
   already normalised) in base 10
11    fd_bar_h[nh[0]] = fd_bar_h.get(nh[0]) + 1/len(nb_bar) #update the values for i digit by 1/N (
   already normalised) in base 16

```

Now we can plot the histogram in base 10 :

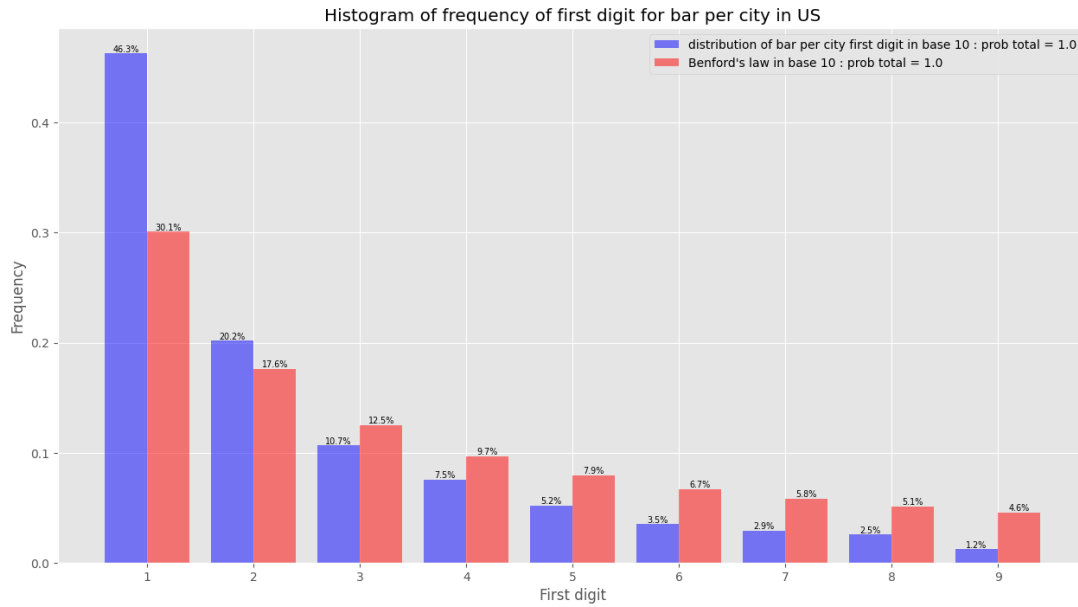


Figure 5: Histogram of distribution of the first digit for bar per city in base 10

In base 10 the number of bar per city has the same tendency of Benford's law and the same probability if we forgot the number one. In this case like in the Benford's law the probability of 1 is greater than the probability of 2 and this probability of 2 is greater than the probability in 3 and so on until 9.

Now we can plot the histogram in base 16 :

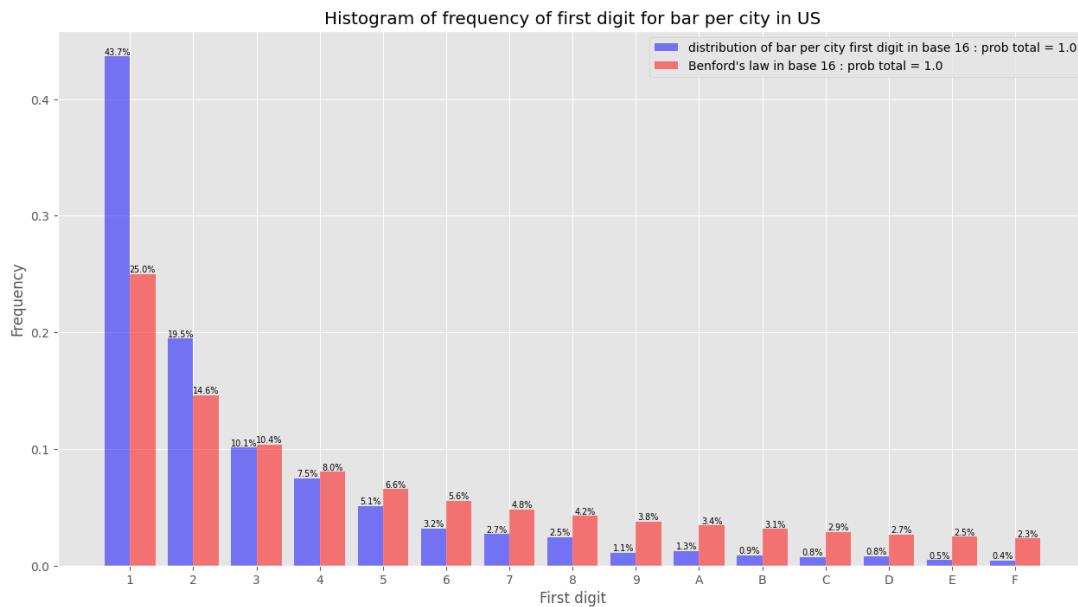


Figure 6: Histogram of distribution of the first digit for bar per city in base 16

Even in base 16 we have the same behavior of number of bar per city like in the previous graph and the same quite close probability but for the number A, B, C, D, E, F their probability are less than the other.



## IV Conclusion

The Benford's law is surprising, we have checked the correlation of it and Fibonacci numbers in base 10 and 16 and we have found very close probability so we can say the Fibonacci numbers follow the Benford's law. For the data set the correlation is more nuanced we have a good tendency but the probability is not exactly the same, we can say the second data set (number of bar per city in US) follows more the Benford's law than the first (number of bar per country) and this can be the consequence of the length of them because the first data set is much more larger (around 7000 bar) than the second (around 1400 bar). The probability of one is much more important in the data set because the probability to have one bar is more important especially in small towns.

With this conclusion we have more information to find a data set that will follow the Benford's law :

- data set with a big quantity of values better
- the data set must not have artificial minimum and maximum or have a big standard deviation
- if the data set follows the Benford's law it should follow in base 16 (maybe in all other bases ? (obviously in base 1 and 2 because all numbers start with one))