

# Most Frequent Letters

Eduardo Santos, n°mec 93107, eduardosantoshf@ua.pt

**Abstract** – The objective of this assignment was to identify the most frequent letters in text files using different methods and to evaluate the quality of estimates regarding the exact counts. Three types of counters were implemented: an Exact Counter, a Decreasing Probability Counter, and a Frequent Counter.

**Keywords** – Counting Algorithms, Data Stream Algorithms, Probabilistic Counters, Exact Counter, Decreasing Probability Counter, Frequent Counter

## I. INTRODUCTION

### A. Decreasing Probability Counter

The aim of a probabilistic counter is to count very large numbers using only a little space to store the counter. Counting a very large number of events using an exact counter will result in a large memory usage, which is something that should be avoided, as memory is expensive, and makes a program less efficient. To mitigate this problem, probabilistic counters were created. So, for each call to an increment method on the counter, its actual value is updated with probability  $p$ . Using this method, we are trading accuracy for the ability to count up to very large numbers with little storage space.

Even though nowadays memory is no longer scarce, this approach is still useful when treating/counting massive data volumes, when there is a need for quick and memory-efficient processing.

### B. Frequent Counter

A streaming algorithm is an algorithm for processing data streams in which the input is presented as a sequence of items and can be examined in only a few passes, typically just one. The frequent problem happens when, given a sequence of items, we want to identify those which occur most frequently. This can also be expressed as finding all items whose frequency exceeds a specified fraction of the total number of items.

## II. PROBLEM DESCRIPTION

The goal of this assignment was to identify the most frequent letters in literary works from Project Gutenberg [1] using different methods and to evaluate the quality of estimates regarding the exact counts.

In order to accomplish that, three different approaches were developed and tested:

- exact counter

- approximate counter - decreasing probability counter with probability  $\frac{1}{\sqrt{2}^k}$
- frequent counter - Misra & Gries algorithm to identify frequent items in data streams

The testing involved an analysis of the computational efficiency and limitations of the developed approaches was carried out, in terms of absolute and relative errors, computing the mean, minimum, and maximum of each error, as well as computing the standard deviation and variance. Finally, for each method, the most frequent letters were identified, checking if they were in the same relative order.

## III. IMPLEMENTATION DESCRIPTION

### A. Main

Running the **main.py**, with the **-help** flag, a few running options are presented.

```
+ python3 main.py --help
usage: main.py [-h] [-t TEXT] [-s TEXT] {exact,decreasing,frequent} ...

Most Frequent Letters

optional arguments:
  -h, --help            show this help message and exit
  -t TEXT, --text TEXT  Load literary work from text file
  -s TEXT, --stopwords TEXT
                        Load stop-words from text file

Counter type:
{exact,decreasing,frequent}
    Counter type/approach
    exact                Run with an exact counter
    decreasing           Run with a counter based on a decreasing probability counter
    frequent             Run with a counter based on frequent-count
```

Fig. 1: Help Menu of the Main Program

The **-t** flag represents the filename of the literary work to be read and processed by the program, which can be found inside the **/texts/** directory. The **-s** flag represents the filename of the stop-words to be ignored when processing the text file. The stop-words files can be found inside the **/stop-words/** directory. After the previous parameters, we have to specify the counter type to be computed, using one of the following:

- **exact** - count using the exact counter
- **decreasing** - count using the decreasing probability counter (with probability  $\frac{1}{\sqrt{2}^k}$ )
- **frequent** - count using the frequent counter (we must also specify the  $k$  parameter, using the flag **-k**)

### B. Counter

The **counter.py** contains the main logic of the solution, using the **ABC** [3] (Abstract Base Classes) Python module, the **ExactCounter**, **DecreasingProbabilityCounter**, and **FrequentCounter**

classes extend the **Counter** parent class, using an OOP (object-oriented programming) model. Inside the **Counter** class there is the `read_letters()` method, which is responsible for parsing the text file, removing the Project Gutenberg file headers and footers, all stop-words and punctuation marks, as well as converting all letters to uppercase.

```
def read_letters(self):
    with open(self.filename, 'r') as file:
        while True:
            line = file.readline()

            # ignore the Project Gutenberg's
            # file headers
            if line.strip() in [header.
                              value for header in
                              Headers]: break

        while line:
            line = file.readline()

            # ignore the Project Gutenberg's
            # file footers
            if line.strip() in [footer.value
                              for footer in Footers]: break

        for words in line.split():
            # remove all stop-words and
            # punctuation marks
            for word in regex.findall('\p{
            alpha}+', words):
                for letter in word:
                    self.parsed_letters.
                        append(letter.
                                upper())
```

### C. Exact Counter

Inside the **ExactCounter** class we can find the `compute()` method, which is responsible for counting the exact number of occurrences of each letter from the literary work.

### D. Decreasing Probability Counter

The **DecreasingProbabilityCounter** class implements the Decreasing Probability Counter as a probabilistic counter, with the increment being made with probability  $\frac{1}{\sqrt{2}^k}$ , where  $k$  represents the number of occurrences of each letter. If the counter has value  $k$ , the algorithm increases the number of occurrences of the letter with the previously mentioned probability. Due to  $k$  and the probability being inversely proportional, as  $k$  increases, the probability of incrementing the counter will be much smaller. This method allows the counting of a large number of events using a small amount of memory. The estimated value of the counter for each letter can be calculated using the following formula:

$$\frac{\sqrt{2}^k - \sqrt{2} + 1}{\sqrt{2} - 1}$$

### E. Frequent Counter

The **FrequentCount** class implements a Frequent Counter as a Data Stream Algorithm. The goal is to establish an estimate for the frequency of any stream

letter. The frequent count algorithm implemented was the Misra Gries algorithm. This uses a parameter  $k$  that controls the quality of the results given. It uses a `letter: counter` dictionary, with at most  $(k - 1)$  counters, at any time. This algorithm provides, for any letter,  $l$ , a frequency estimate satisfying

$$f_l - \frac{m}{k} \leq f_l^* \leq f_l$$

where  $m$  is the length of the data stream or, in this case, the total number of letters in the text. If some letter has  $f_l > \frac{m}{k}$ , its counter  $A[l]$  will be positive, i.e., no item with frequency  $\frac{m}{k}$  is missed.

### F. Tests

The **tests.py** contains the logic of the tests that are used to compute the results for the Decreasing Probability and Frequent counters, comparing both to the Exact Counter. The obtained results will be discussed in the next section.

## IV. RESULTS AND DISCUSSION

As mentioned in the previous subsection, some tests were developed and used to compare both algorithms. Regarding the Decreasing Probability Counter, the following measures were computed: absolute error (mean, minimum, and maximum), relative error (mean, minimum, and maximum), standard deviation, and variance. All measures were calculated as an average from one hundred trials, and the comparison was made using the estimated count from each letter from the text. Regarding the Frequent Counter, due to being a deterministic algorithm, the results are always the same. So, one run of the tests gives the results needed for comparison. For this counter, the  $k$  most frequent letters were computed. For both tests, the literary work used was *Crime and Punishment*, by Fyodor Dostoevsky in both English and Spanish languages. Inside the `/texts/` folder, *The Metamorphosis*, by Franz Kafka can also be found, in both English and German languages.

### A. Decreasing Probability Counter

The figure 2 shows the results for the *Crime and Punishment*, by Fyodor Dostoevsky text, in English, which has, without the stop-words, 874218 letters.

Letter	Value	Exact Value	Expected Value	Mean	Min	Max	Mean	Min	Max	Standard	Variance
a	10	28111	28106	28091	28076	28121	28076	28061	28111	0.00000	0.00000
b	10	79429	79434	79449	79464	79434	79464	79479	79429	0.00000	0.00000
c	10	28821	28826	28841	28856	28811	28856	28871	28821	0.00000	0.00000
d	10	70808	70813	70828	70843	70798	70843	70858	70808	0.00000	0.00000
e	10	81114	81119	81134	81149	81104	81149	81164	81114	0.00000	0.00000
f	10	61114	61119	61134	61149	61104	61149	61164	61114	0.00000	0.00000
g	10	41114	41119	41134	41149	41104	41149	41164	41114	0.00000	0.00000
h	10	31114	31119	31134	31149	31104	31149	31164	31114	0.00000	0.00000
i	10	21114	21119	21134	21149	21104	21149	21164	21114	0.00000	0.00000
j	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
k	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
l	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
m	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
n	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
o	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
p	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
q	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
r	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
s	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
t	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
u	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
v	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
w	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
x	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
y	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
z	10	11114	11119	11134	11149	11104	11149	11164	11114	0.00000	0.00000
Execution time (for 100 trials):	131.32s										

Fig. 2: Results obtained using the Decreasing Probability Counter, in English

From these results, there are a few conclusions we can take. For each letter, the expected (estimated) value has a high deviation from the real value. This can be confirmed by the mean absolute and relative error columns. Taking, for example, the letter *E*, the mean relative error is  $34.4\%$ , meaning that the accuracy of the results is less than  $70\%$ , which I believe is not a bad result, considering that this algorithm only increments each letter's counter with probability  $\frac{1}{\sqrt{2}^k}$ . There is also a relatively high standard deviation, and thus, variance. On the other hand, memory-wise this method is much more efficient than the exact count. It is worth mentioning that there are some letters whose mean relative error is  $0.0\%$ , translating into a  $100\%$  accuracy.

Letter	Value	Exact Value	Expected Value	Mean Abs. Error	Mean Rel. Error	Standard Dev.	Max. Error	Min. Error	Mean Abs. Error	Mean Rel. Error	Standard Dev.	Max. Error	Min. Error
E	31	187588	114717	27482.9	14.64	28886.0	34.48	1.08	154.18	1.297444	2.486140	34.48	1.08
T	30	98958	97580	2122.7	2.15	12820.8	12.96	12.28	126.25	1.272446	1.486140	21.03	22.28
A	30	72334	73811	6773.0	9.36	24338.0	34.08	9.48	337.35	4.616246	1.433140	34.08	9.48
O	29	59877	66618	2952.1	5.00	9313.0	15.68	5.28	187.25	3.152146	7.332140	15.68	5.28
N	29	55842	55842	12862.8	22.86	12374.0	22.35	8.28	185.25	2.488446	8.138540	22.35	8.28
I	28	52476	49974	1945.0	3.71	3923.0	7.45	8.08	132.25	2.432146	6.298140	7.45	8.08
S	28	44537	44537	12827.5	28.57	3774.0	8.45	14.78	78.05	1.594446	2.343140	8.45	14.78
R	28	48855	48855	12548.0	25.68	7888.0	16.75	1.28	176.25	1.691446	2.486140	16.75	1.28
D	28	38958	38958	13839.5	35.52	411.0	1.05	1.68	386.25	2.852446	3.837140	29.18	1.68
L	28	36388	36388	16867.8	46.38	7791.0	21.43	7.48	286.25	3.902146	3.837140	9.48	7.48
U	27	35238	35238	15440.0	43.81	4317.0	13.48	9.28	168.25	1.533146	2.898446	13.48	9.28
C	27	36774	36774	15430.0	41.98	4323.0	12.15	9.18	137.15	1.558446	2.432140	12.15	9.18
V	27	22881	22881	16423.0	71.82	3194.0	13.22	8.28	181.25	1.533146	8.212140	13.22	8.28
F	26	21588	21588	23823.0	110.33	3447.0	13.78	8.08	168.25	8.212140	8.425140	13.78	8.08
M	26	18889	18889	18877.0	99.98	2347.0	12.45	8.28	161.25	8.416446	8.416446	12.45	8.28
H	26	18845	18845	18845.0	100.00	3984.0	21.18	18.18	151.78	4.816446	2.309446	5.48	18.18
G	26	18812	17378	1445.0	7.68	3124.0	16.45	8.28	139.25	8.748446	2.747140	16.45	8.28
X	26	18811	2222.8	239.0	1.27	2199.0	11.78	2.48	388.25	5.119446	2.486140	11.78	2.48
Q	26	18812	2097.0	1672.0	8.84	1112.0	5.75	8.08	189.25	3.558446	1.228140	5.75	8.08
J	23	7226	7477	2548.0	35.19	1258.0	17.48	3.18	173.78	1.582140	1.188446	17.48	3.18
K	23	7885	7885	1881.0	23.84	897.0	11.78	1.18	17.25	1.621440	8.816446	11.78	1.18
Y	23	4894	4894	2828.8	57.80	37.0	0.76	0.76	186.48	2.788446	7.838446	21.08	8.08
Z	23	4894	1643.0	1212.1	24.77	648.0	11.78	1.18	186.48	1.552140	1.184446	17.15	1.18
W	21	4578	4578	3438.0	75.14	365.0	9.45	8.08	289.48	1.922140	3.888446	18.08	8.08
P	21	3881	3881	1191.0	30.69	648.0	18.75	11.18	149.48	1.534446	2.378446	18.75	11.18
B	21	3492	3492	1445.0	41.38	615.0	12.08	8.18	183.18	1.587140	2.129446	12.08	8.18
0	21	3487	3487	1191.0	34.16	648.0	18.75	11.18	149.48	1.534446	2.378446	18.75	11.18
1	20	2875	2875	2875.0	100.00	972.1	33.84	14.18	143.18	1.238446	1.128446	11.08	14.18
2	19	1545	1545	1545.0	100.00	487.5	31.55	12.78	147.08	6.778446	3.942446	10.48	12.78
3	17	892	892	892.0	100.00	368.7	41.30	2.18	127.25	6.964446	1.178446	24.18	2.18
4	17	892	892	892.0	100.00	368.7	41.30	2.18	127.25	6.964446	1.178446	24.18	2.18
5	17	892	892	892.0	100.00	368.7	41.30	2.18	127.25	6.964446	1.178446	24.18	2.18
6	17	892	892	892.0	100.00	368.7	41.30	2.18	127.25	6.964446	1.178446	24.18	2.18
7	17	892	892	892.0	100.00	368.7	41.30	2.18	127.25	6.964446	1.178446	24.18	2.18
8	17	892	892	892.0	100.00	368.7	41.30	2.18	127.25	6.964446	1.178446	24.18	2.18
9	17	892	892	892.0	100.00	368.7	41.30	2.18	127.25	6.964446	1.178446	24.18	2.18
[Execution time (for 100 trials): 109.13s]													

stream algorithms. Both the decreasing probability counter and the frequent counter can be good solutions to reducing memory usage, even though they have different applications, as the use cases differ between the two. Regarding future work, it would be relevant to compare the algorithms explored with different ones, to have an even better term of comparison and knowledge about this type of counting methods.

#### REFERENCES

- [1] Project Gutenberg. (1971-2021). Welcome to Project Gutenberg. <https://www.gutenberg.org/>
- [2] Python Software Foundation. (Dec 22, 2022). re — Regular expression operations. <https://docs.python.org/3/library/re.html>
- [3] Python Software Foundation. (Dec 22, 2022). abc — Abstract Base Classes. <https://docs.python.org/3/library/abc.html>
- [4] Wikipedia contributors. (Dec 24, 2022). Approximate counting algorithm. Wikipedia. [https://en.wikipedia.org/wiki/Approximate\\_counting\\_algorithm](https://en.wikipedia.org/wiki/Approximate_counting_algorithm)
- [5] Wikipedia contributors. (Dec 24, 2022). Approximate counting algorithm. Wikipedia. [https://en.wikipedia.org/wiki/Approximate\\_counting\\_algorithm](https://en.wikipedia.org/wiki/Approximate_counting_algorithm)
- [6] Wikipedia contributors. (Jan 2, 2023). Streaming algorithm. Wikipedia. [https://en.wikipedia.org/wiki/Streaming\\_algorithm#:~:text=In%20computer%20science%2C%20streaming%20algorithms,passes%20\(typically%20just%20one\)](https://en.wikipedia.org/wiki/Streaming_algorithm#:~:text=In%20computer%20science%2C%20streaming%20algorithms,passes%20(typically%20just%20one))