

Count the Number of Occurrences of Letters in Text Files using Probabilistic Methods

João Carvalho, 89059

Resumo - Neste relatório é abordado e analisado o processo de desenvolvimento e resultados de *counters* de letras num determinado texto usando algoritmos probabilísticos.

O objetivo deste problema é fazer uma estimativa do número de vezes que cada letra aparece num certo texto usando um *counting method* com probabilidades fixa e decrescente.

Ao longo do relatório são analisados e comparados os resultados obtidos entre os dois métodos e tendo por base a contagem exata de cada letra.

I. INTRODUCTION

A implementação de counters probabilísticos é bastante usada em ambientes de grande quantidade de dados, uma vez que usa menos espaço de memória para armazenamento do *counter*.

Neste trabalho foi tratado o desenvolvimento de um *counter* com probabilidade fixa(*FC*) - $1 / 2^k$, especificamente com k de valor 2 ($1/4$), e um *counter* com probabilidade decrescente(*DC*) - $1 / 2^k$ - com k representando o valor do *counter*.

Este métodos foram testados em várias repetições para a obtenção de resultados fidedignos cujos valores são, posteriormente, analisados e comparados em diferentes aspetos recorrendo ao counter exato(*EX*).

II. FIXED COUNTER

Primeiramente será abordado o *counter* com probabilidade fixa. Neste tipo de *counter*, um evento é contabilizado tendo em conta uma certa probabilidade, neste caso $p = 1/4$. Com esta definição é possível perceber que, no resultado final, a contagem do número de vezes que uma letra aparece num texto será aproximadamente $1/4$ do valor exato dessa letra. Em outras palavras, para estimar o número real da contagem do *FC*, apenas se multiplica o valor do counter por 4.

Abaixo, na Fig.1, encontra-se um gráfico com a comparação do *FC* com o *EX* para texto escrito em finlandês(filandes.txt) com um total de 30878 letras(*len*). O *FC*, neste caso, realizou 100 repetições, isto é, efetuou a

contagem das letras 100 vezes, e os valores no gráfico representam o mean value.

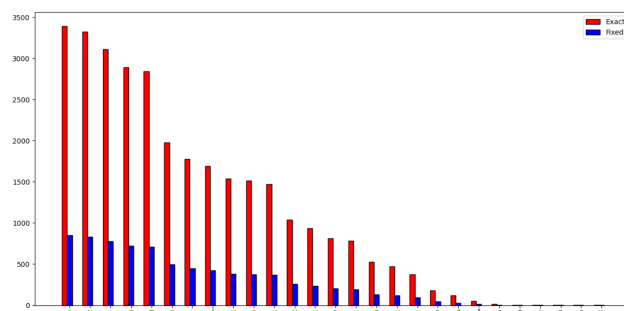


Fig. 1: comparação do FC e EX

O número de repetições é bastante relevante pois, como se trata de um processo probabilístico, um número baixo de repetições pode levar a resultados pouco precisos.

Na Fig.2 está retratado um gráfico que compara a exatidão média em percentagem do *FC* com quatro diferentes números de repetições: 2(vermelho), 10(verde), 100(azul), 1000(preto).

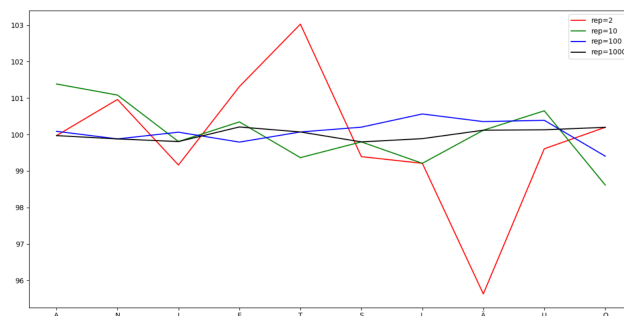


Fig. 2: comparação de exatidão do FC com diferentes números de repetições

A exatidão foi calculada com a divisão do valor médio esperado pelo valor exato. Pode-se observar que quanto maior o número de repetições maior a exatidão(valores mais próximos de 100%).

Neste ponto, será analisado mais profundamente o algoritmo *FC* tendo por base algumas métricas: *Estimated count*(Mean *FC*), *Mean absolute deviation*, *Mean accuracy*, *Max error*, *Min error*, *Mean error*. Estes três últimos com valores absolutos e relativos.

Na Fig.3 está representada uma tabela com estes valores resultantes de um teste com o exemplo de texto testado anteriormente(filandes.txt) e com 100 repetições. As letras

estão organizadas por ordem decrescente do seu número(contagem).

test: 10000											
Char	EC	1 Mean	1 Mean	1 Mean Dev	1 Mean Acc	Abs/Rel Max	Max Error	Abs/Rel Min	Min Error	Abs/Rel Mean	Mean Error
A	3393	3386	82.52	99.8	317	-0.9	1	-0.6	82.8	-0.4	-0.4
C	3388	3386	82.52	99.8	317	-0.9	1	-0.6	82.8	-0.4	-0.4
I	3114	3109	72.86	99.8	246	-0.7	1	-0.1	72.5	-0.3	-0.3
E	2890	2901	69.64	100.0	218	-0.5	2	-0.1	71.2	-0.5	-0.5
S	2822	2835	68.82	100.0	202	-0.5	2	-0.1	68.5	-0.3	-0.3
L	1898	1996	51.96	100.8	220	-11.1	0	-0.6	54.4	-2.7	-2.7
U	1794	1767	59.8	99.6	0	-0.9	2	-0.1	68.6	-3.8	-3.8
O	1694	1683	58.76	99.4	158	-0.9	2	-0.1	55.4	-3.3	-3.3
U	1538	1589	55.76	100.0	182	-11.8	2	-0.1	58.8	-3.8	-3.8
M	1532	1532	56.52	99.7	158	-11.8	2	-0.1	58.8	-3.8	-3.8
K	1075	1088	56.8	100.3	203	-13.8	0	-0.1	56.6	-3.8	-3.8
0	1042	1032	42.92	99.9	150	-14.4	2	-0.2	42.9	-0.1	-0.1
U	917	943	31.3	100.1	107	-11.8	1	-0.1	31.5	-0.5	-0.5
R	811	810	42.0	99.9	147	-11.8	1	-0.1	42.1	-0.5	-0.5
V	784	786	43.76	100.3	104	-13.8	0	-0.0	43.6	-5.2	-5.2
C	725	725	49.36	99.5	97	-11.8	1	-0.1	49.7	-5.1	-5.1
V	473	472	32.6	99.8	87	-11.8	1	-0.2	32.6	-6.9	-6.9
J	376	376	31.4	100.0	80	-22.3	0	-0.0	31.4	-8.4	-8.4
F	363	363	31.4	100.0	80	-22.3	0	-0.0	31.4	-8.4	-8.4
Q	120	123	17.14	102.5	52	-0.3	0	-0.0	17.2	-14.3	-14.3
Q	52	55	19.8	105.8	32	-0.5	0	-0.0	19.8	-20.8	-20.8
Q	12	12	11.2	100.0	10	-0.3	0	-0.0	11.2	-12.8	-12.8
E	6	6	3.76	100.0	10	-16.6	7	-33.3	3.8	-63.3	-63.3
4	4	4	2.6	100.0	0	-20.6	0	-0.9	2.6	-65.0	-65.0
4	4	4	2.6	100.0	0	-20.6	0	-0.9	2.6	-65.0	-65.0
3	3	3	2.48	100.0	5	-16.6	7	-33.3	2.5	-83.3	-83.3
3	1	1	1.58	100.0	5	-30.6	0	-100.0	1.6	-160.0	-160.0

Fig. 3: *FC* com 100 repetições

E, para se poder ter uma comparação de valores com um número diferente de repetições, tem-se na Fig.4 uma tabela de um teste que utilizou o mesmo ficheiro de texto porém com 10000 repetições.

ext. loss: 30078												
Char	EC	Mean F1C	Mean Dev	Mean Acc	Abs/Rel Max	Max Error	Abs/Rel Min	Min Error	Abs/Rel Mean	Mean Error		
A	3393	3391	88.2932	99.9	437	-112.9	1	-0.0	88.3	-2.4		
J	3328	3330	88.2932	99.9	408	-112.9	1	-0.0	79.7	7.7		
I	3114	3115	77.0088	100.0	418	-113.4	2	-0.1	77.0	-2.5		
E	2890	2889	74.0888	100.0	378	-113.1	2	-0.1	76.4	-2.2		
C	2842	2844	74.0888	100.0	394	-113.1	2	-0.1	76.3	-2.3		
N	1908	1908	61.5992	100.0	308	-115.6	0	-0.0	61.6	-1.1		
L	1774	1775	58.1712	100.0	1	-117.2	2	-0.1	58.2	-3.7		
U	1694	1692	55.2728	99.9	278	-115.9	1	-0.1	56.3	-3.3		
U	1538	1537	53.67	99.9	2	-116.8	2	-0.1	53.7	-3.5		
U	1537	1516	53.8888	99.9	3	-120.6	1	-0.1	53.9	-6.9		
O	1475	1475	51.6888	99.9	275	-116.8	1	-0.1	51.7	-3.3		
N	1402	1401	44.888	99.9	222	-121.3	2	-0.2	44.1	-4.2		
H	933	932	42.5152	99.9	223	-121.3	1	-0.1	42.4	-0.8		
H	833	1311	38.8888	100.0	199	-119.7	1	-0.1	39.1	-1.4		
V	784	784	34.3888	100.0	204	-126.0	0	-0.0	38.6	-4.9		
P	528	528	31.728	100.0	100	-126.0	0	-0.0	32.4	-0.9		
J	473	473	25.9776	99.9	139	-126.0	0	-0.2	39.0	-0.7		
V	376	376	26.9276	100.0	124	-133.0	0	-0.0	26.9	-7.2		
D	183	183	20.9276	100.0	9	-133.0	1	-0.1	20.9	-1.4		
O	128	128	14.9688	100.0	76	-163.3	0	-0.0	15.0	-12.8		
A	52	52	9.846	100.0	48	-92.3	0	-0.0	9.8	-3.5		
E	276	276	35.8888	100.0	20	-98.0	0	-0.0	36.0	-18.0		
E	6	6	3.5558	100.0	1	-398.0	2	-39.3	3.4	-60.5		
A	4	4	2.5112	100.0	12	-398.0	0	-9.0	2.5	-62.5		
E	4	4	2.5112	100.0	9	-398.0	1	-11.2	2.5	-62.5		
F	3	3	2.5296	100.0	9	-398.0	1	-13.3	2.5	-83.3		
W	1	1	1.5976	100.0	3	-398.0	1	-100.0	1.5	-158.0		

Fig. 4: Análise do FC com 10000 repetições

A seguir, na Fig. 5, pode-se ver a comparação de valores dos quatro tipos de valor do *counter*: valor exato, estimativa do valor do *FC*, valor do *FC* e a estimativa do valor exato. Teste realizado com 1000 repetições sobre o mesmo ficheiro sendo apresentadas apenas as letras mais comuns(para o *FC*).

```
g_AA$ python3 src/main.py texts/filandes.txt 1000
Text length: 30878
Char -- EC      FC Est. Counter  FC      Mean Est. FC      Mean Acc
A -- 3393      848.25      850.3      3461      100.24
I -- 3328      832.0      831.8      3327      99.97
N -- 3114      778.5      778.4      3114      100.0
E -- 2890      722.5      722.5      2890      100.0
T -- 2842      710.5      711.1      2844      100.07
S -- 1980      495.0      495.3      1981      100.05
L -- 1774      443.5      444.4      1778      100.23
Ä -- 1694      423.5      422.5      1690      99.76
U -- 1538      384.5      384.9      1539      100.07
O -- 1517      379.25      379.4      1517      100.0
K -- 1475      368.75      368.5      1474      99.93
M -- 1042      260.5      260.4      1042      100.0
```

Fig. 5: Comparação de valores dos *counters*

A estimativa do valor exato foi calculada, como já referido, pela multiplicação do valor do *counter* por 4 (fazendo a média posteriormente), já a estimativa do valor do *FC* é o inverso, divide-se o valor exato por 4.

A compreensão destas tabelas será feita no bloco de análise de resultados a seguir no relatório.

III. DECREASING COUNTER

Agora será abordado o *counter* com probabilidade decrescente.

Neste tipo de counter, um evento é contabilizado tendo em conta uma probabilidade que vai decrescendo à medida que aumenta a contagem, neste caso $p = 1 / 2^k$ (k representa o valor do *counter*). Nesta situação também se pode definir como *log counter* pois o valor do *counter* aumenta de forma logarítmica em relação ao número de eventos. Assim, é possível perceber que, no resultado final, a contagem do número de vezes que uma letra aparece num texto será aproximadamente o valor inteiro de $\log_2(n+1)$ com n sendo o valor exato dessa letra. Para estimar o número de surgimentos da letra no texto com base no counter pode-se recorrer à expressão $2^k - 1$.

Abaixo, na Fig.6, encontra-se um gráfico com a comparação do *DC*(100 repetições) com o *EX* para texto escrito em finlandês(*filandes.txt*) com um total de 30878 letras(*len*). Contudo, contrariamente ao exemplo do *FC*(Fig.1), os valores do *DC* são já os valores estimados pois os valores do *counter* são demasiado pequenos para serem visíveis.

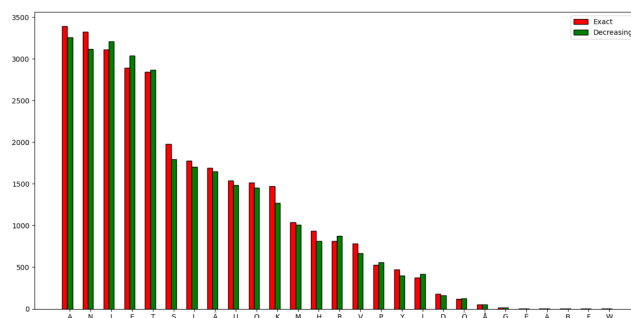


Fig. 6: comparação do DC e EC

Na Fig.7, abaixo, está o gráfico da exatidão para o *DC* com quatro diferentes números de repetições:

2(vermelho), 10(verde), 100(azul), 1000(preto). Estão expostas apenas as letras mais comuns.

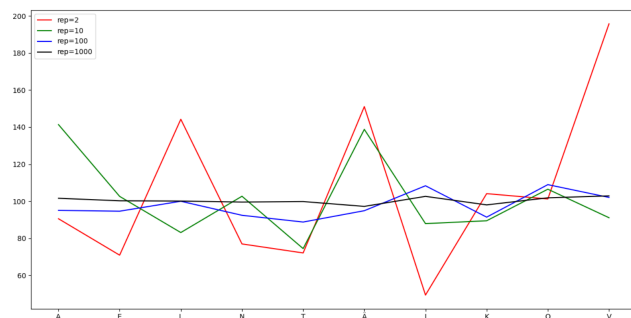


Fig. 7: comparação de exatidão do *FC* com diferentes números de repetições

Fazendo já uma rápida análise de ambos os gráficos de exatidão (Fig.2, Fig.6), pode-se já verificar que o DC

precisa de realizar mais repetições para adquirir uma exatidão razoável, conclusão que, já seria esperada.

Agora, irá analisar-se o *DC* recorrendo à tabela com as métricas já referidas anteriormente. Na Fig.8 pode-se observá-la num teste com 100 repetições do mesmo ficheiro usado anteriormente(filandês.txt).

[illegible]

Fig. 8: *DC* com 100 repetições

E na Fig. 9 encontra-se um teste com 10000 repetições.

Char	EC	Mean FC	Mean Dev	Mean Acc	Abs/Rel Max Error	Abs/Rel Min Error	Abs/Rel Mean Error
A	3393	3411	1768.38	108.53	29374	-865.7	762
C	3378	3276	1818.85	111.35	29385	-962.2	981
E	3114	3149	1623.89	101.12	29653	-952.2	981
I	2890	2913	1569.13	100.88	29977	-1033.8	843
S	2842	2830	1607.57	101.76	33461	-976.3	995
T	1988	2080	898.46	101.01	19403	-717.4	67
L	1774	1783	871.54	100.51	16469	-823.5	797
N	1694	1698	848.57	100.84	16459	-967.1	730
O	1538	1536	789.67	99.87	31229	-838.5	589
U	1517	1522	789.15	100.23	31866	-988.1	499
M	1456	1471	789.15	100.23	31866	-988.1	499
N	1042	1045	464.26	100.39	7149	-686.0	19
H	933	932	403.05	99.89	7258	-777.9	99
P	797	797	398.78	99.87	7407	-777.9	99
V	784	783	399.78	99.87	7407	-844.4	239
P	528	527	237.23	100.61	3567	-676.6	17
Q	463	469	196.16	100.61	3567	-676.6	17
J	376	379	196.16	100.8	3719	-981.9	121
I	184	184	96.06	101.66	1856	-1038.9	84
K	121	120	86.83	100.84	1856	-1038.9	84
S	52	52	25.94	100.49	4659	-882.7	11
K	16	16	7.18	100.6	239	-1498.1	1
A	6	6	2.99	100.77	239	-1498.1	1
G	4	4	1.93	100.8	11	-275.0	1
B	3	3	1.41	100.8	1	-133.3	0
W	3	3	1.41	100.8	1	-133.3	0
X	1	1	0.0	100.0	0	-0.0	0

Fig. 9: *DC* com 10000 repetições

A seguir, na Fig. 10, pode-se ver a comparação de valores dos quatro tipos de valor do *counter*. Teste realizado com 1000 repetições sobre o mesmo ficheiro sendo apresentadas apenas as letras mais comuns(para o DC).

```
g_AA$ python3 src/main.py texts/filandes.txt 1000
Text length: 30878
```

Char	EC	DC	Est. Counter	DC	Mean	Est. DC	Mean Acc
N	3328	11		11.5	3344		100.48
A	3393	11		11.4	3266		96.26
I	3114	11		11.3	3173		101.89
E	2890	11		11.2	2899		100.31
T	2842	11		11.2	2773		97.57
S	1980	10		10.7	2034		102.73
L	1774	10		10.5	1706		96.17
Ä	1694	10		10.4	1655		97.7
O	1517	10		10.3	1592		104.94
U	1538	10		10.3	1553		100.98
K	1475	10		10.2	1427		96.75
M	1042	10		9.7	984		90.43

Fig.10: Comparação de valores dos *counters*

A estimativa do valor exato foi calculada, como já referido, pela expressão $2^k - 1$, já a estimativa do valor do DC é o valor inteiro dado pela expressão $\log_2(n+1)$.

IV. RESULTS ANALYSIS

Numa visão geral das tabelas de métricas pode-se tirar algumas conclusões. Tendo em conta os valores absolutos, *mean deviation*, *abs max error*, *abs mean error*, é possível verificar um **aumento** destes à medida que **aumenta** o *counter* da letra(letras mais comuns), por outro lado, os valores relativos **diminuem** à medida que o *counter* da letra aumenta mas apenas no *FC*, no caso de *DC* os valores relativos parecem estar uniformemente distribuídos pelas letras. Estes valores já eram esperados uma vez que as letras com o maior *counter* passaram, obviamente, mais vezes pelo processo de contagem probabilística, tendo uma maior quantidade de informação. Pode-se ainda observar que, como já referido, o teste com 10000 repetições apresenta uma ***mean accuracy***, ou seja, uma **exatidão mais próxima** dos 100% num geral de todas as letras, aliás, no caso do *FC*, todas as letras obtiveram uma *mean accuracy* de 100% com uma margem de 0.1%.

Falando de **exatidão** pode-se extrair das tabelas que o *FC* tende a ter valores mais exatos quando comparado ao *DC* em testes com o mesmo número de repetições, resultado previsível pois o *FC* usa mais memória(*counters maiores*).

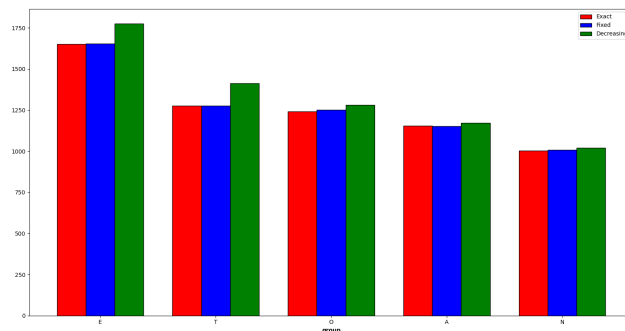
Outra observação importante é relacionada com a diferença de valores de *errors* quando comparados os dois tipos. Consegue-se observar que, o *mean deviation* e o *abs mean error* são **muito maiores no DC**, isto significa que, apesar de ser razoavelmente exata, este tipo de *counter* é muito menos **preciso** que o *FC*.

Uma nota adicional é que, quando se ordena as letras por ordem decrescente do seu valor do *counter*, a lista pode ser diferente quando comparados os *DC* e *FC* com o *EC* sobretudo com números de repetições baixos(e sobretudo o *DC*), contudo com números de repetições mais elevados essa diferença praticamente não existe.

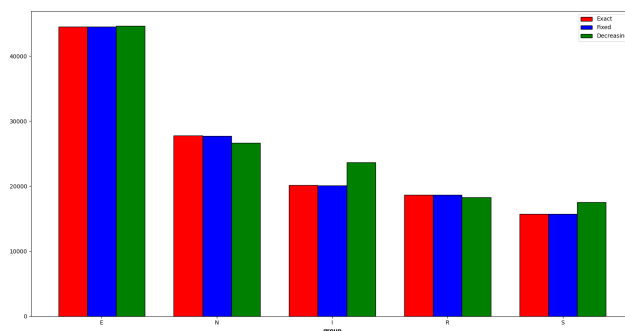
V. EXTRA TESTS

Neste bloco serão apresentados alguns testes destes *counters* com ficheiros de linguagens diferentes.

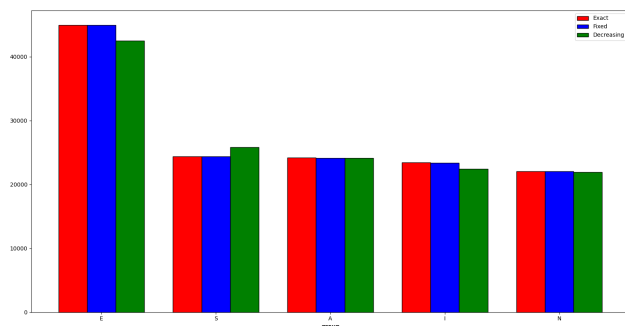
Assim podemos verificar outros valores dos *counters*, além de curiosidades como saber as letras mais comuns de cada linguagem. Todos os testes foram feitos com 100 repetições. É curioso reparar que as letras mais comuns de algumas das línguas mais faladas da europa são muito coincidentes.

Fig.11: Texto em inglês com $len=14296$

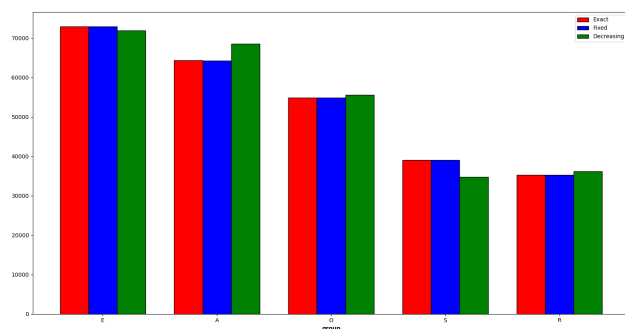
E - T - O - A - N

Fig.12: Texto em alemão com $len=257920$

E - N - I - R - S

Fig.13: Texto em francês com $len=311776$

E - S - A - I - N

Fig.13: Texto em português com $len=533657$

E - A - O - S - R

VI. CONCLUSION

Por fim pode-se concluir que os resultados obtidos foram ao encontro das expectativas prévias. Os algoritmos provaram-se bastante fidedignos com um número razoavelmente alto de repetições.

Assim conclui-se que o processo de contagem probabilística é uma boa técnica no que diz respeito ao menor uso de memória em *Big Data*.

REFERENCES

[1]<https://stackoverflow.com/questions/14576083/can-someone-explain-how-probabilistic-counting-works>

[2]https://elearning.ua.pt/pluginfile.php/2931377/mod_resource/content/0/AA_09_Probabilistic_Counters.pdf