

# TI – zad 3

Kamil Kowalczyk, 136742, L1, czwartek 13:30

Do rozwiązania wszystkich zadań, na wejściu podano korpus Wikipedii.

## Zad 1

Entropia przybliżenia zerowego rzędu:

$$H(X) = -37 * \left(\frac{1}{37} * \log_2\left(\frac{1}{37}\right)\right) \approx 5.21$$

## Zad 2

Wynik działania programu został przedstawiony na poniższym zrzucie ekranu.

filename	S1	SC1	SC2	SC3	W1	WC1	WC2	WC3
norm_wiki_en.txt	4.29   100.00	3.52   82.01	3.02   70.39	2.48   57.87	11.54   100.00	6.39   55.35	2.18   18.85	0.48   4.20
norm_wiki_eo.txt	4.18   100.00	3.34   79.97	2.87   68.76	2.39   57.28	11.56   100.00	6.56   56.72	2.48   21.49	0.63   5.48
norm_wiki_et.txt	4.17   100.00	3.51   84.10	3.13   75.17	2.61   62.61	13.75   100.00	5.42   39.46	0.90   6.58	0.12   0.85
norm_wiki_ht.txt	4.15   100.00	3.11   75.10	2.27   54.83	1.49   35.99	8.17   100.00	3.19   39.10	1.31   16.06	0.81   9.95
norm_wiki_la.txt	4.23   100.00	3.45   81.60	2.82   66.78	2.15   50.90	11.97   100.00	4.40   36.76	1.17   9.75	0.39   3.24
norm_wiki_nv.txt	3.87   100.00	2.95   76.06	2.37   61.10	1.80   46.33	9.15   100.00	3.86   42.21	1.72   18.78	0.90   9.82
norm_wiki_so.txt	4.04   100.00	3.30   81.67	2.84   70.40	2.37   58.77	11.73   100.00	5.40   46.02	1.61   13.71	0.41   3.49
sample1.txt	4.27   100.00	2.92   68.24	2.00   46.81	1.54   36.02	7.75   100.00	7.49   96.61	4.41   56.87	0.60   7.68
sample2.txt	4.13   100.00	3.24   78.49	2.86   69.33	2.33   56.38	11.50   100.00	5.37   46.71	1.57   13.69	0.51   4.41
sample3.txt	3.99   100.00	3.05   76.39	2.47   61.79	1.94   48.58	8.02   100.00	7.35   91.58	3.78   47.13	0.86   10.71
sample4.txt	3.93   100.00	3.18   81.02	2.63   66.86	2.02   51.50	9.06   100.00	5.95   65.67	2.63   29.03	1.26   13.95
sample5.txt	4.25   100.00	4.23   99.42	4.23   99.37	4.18   98.23	17.13   100.00	3.44   20.11	0.23   1.37	0.00   0.02

Zrzut prezentuje tabelę. W pierwszej kolumnie znajduje się nazwa pliku wejściowego. Każdy plik wejściowy jest opisany dwoma wierszami. W pierwszym znajdują się miary entropii, oraz entropii warunkowych wyliczonych dla pliku. Kolumny S1 i W1 reprezentują entropię pierwszego rzędu odpowiednio dla znaków (signs) i dla słów (words). Następnie kolumny opisane jako SCk i WCK oznaczają entropie warunkowe k-tego rzędu odpowiednio dla znaków i słów. Pod każdym wierszem z wynikami dla konkretnego pliku wejściowego znajduje się wiersz pomocniczy, zawierający wartości procentowe. Ukazują one tempo zmieniania się entropii, wraz z rosnącym rzędem. Wartością stuprocentową są zawsze klasyczne wartości entropii. Występujące po nich wartości w kolumnach ukazują jaką część pierwotnej entropii stanowi wyliczona nowa entropia warunkowa k-tego rzędu. Ze względu na to, że dane zostały zgromadzone zarówno dla znaków jak i dla słów, w tabeli kolumny S1, SC1, SC2 i SC3, dotyczą znaków i ich wartości procentowe są wyliczane niezależnie w stosunku do kolumn W1, WC1, WC2, i WC3, które odpowiadają słowom.

Odpowiedzi na pytanie, czy plik zawiera język naturalny:

sample1.txt - FAŁSZ

sample2.txt - PRAWDA

sample3.txt - FAŁSZ

sample4.txt - PRAWDA

sample5.txt - FAŁSZ

Plik sample5.txt jest oczywistym przykładem pliku, który nie zawiera języka naturalnego. Widać to po tempie zmniejszania się wartości entropii warunkowych dla znaków przy rosnących rzędach, które w językach naturalnych jest szybsze. Analogiczna sytuacja występuje w plikach sample1.txt i sample3.txt, z tą różnicą, że tu tempo zmniejszania się entropii warunkowych słów (a nie znaków) jest zbyt niskie. Plik sample2.txt posiada zbliżone wyniki do plików z językami naturalnymi, dlatego bez wątpliwości można go zaklasyfikować do grupy poprawnych plików. Natomiast ciekawym przypadkiem jest sample4.txt. Dla znaków, jego entropie mają podobną tendencję do poprawnych plików. Jednak w przypadku słów, tempo zmniejszania się wartości entropii warunkowych jest zauważalnie wolniejsze. Mimo to, ze względu na poprawne wyniki entropii dotyczących znaków i brak większych anomalii został on uznany za plik zawierający język naturalny.