

Big Data

Lista zadań

Jacek Cichoń, WiT, PWr, 2022/23

1 Wstep

Zadanie 1 — Pobierz plik z kilkoma dramatami Szekspira ze strony wykładu. Wybierz jeden z dramatów.

1. Oczyszczyć wybrany plik. Podzielić go na słowa.
2. Usunąć z niego "Stop Words" i usunąć z niego słowa o długości mniejszej lub równej 2.
3. Zbudować chmurę wyrazów (word cloud) z otrzymanej listy. Można skorzystać np. z serwisu <http://www.wordclouds.com/>

Celem tego zadania jest wygenerowanie mniej więcej takiego obrazka (dla poematu "Pan Tadeusz"):



Zadanie 2 — To jest kontynuacja poprzedniego zadania.

1. Zastosuj część funkcji które napisałeś do realizacji poprzedniego zadania do wyznaczenia indeksów TF.IDF dla wszystkich wyrazów z dokumentów w dramatach Szekspira znajdujących się w pliku ze strony wykładu.
2. Zbuduj chmury wyrazów oparte o TF.IDF dla wszystkich rozważanych dramatów.

Zadanie 3 — Pokaż, że jeśli chcesz jednoznacznie wyreprezentować każdą z liczb ze zbioru $\{0, 1, \dots, n\}$ za pomocą b bitów to $b \geq \lceil \log_2(n+1) \rceil$.

Zadanie 4 — Pokaż, że jeśli $x = \sum_{k=0}^s a_k 2^k$, gdzie $a_i \in \{0, 1\}$ oraz $a_s = 1$ to $s = \lceil \log_2(x+1) \rceil$

Zadanie 5 — Rozważmy następującą modyfikację licznika Morrisa: ustalamy liczbę $\alpha > 0$ oraz rozważamy tak oprogramowany licznik:

```

init :: C=0
onInc :: if  $\left( random() < \left( \frac{1}{1+\alpha} \right)^C \right)$  then C = C+1
onGet :: return (?????)

```

Niech C_n oznacza wartość zmiennej C po n wywołaniach metody `onInc`.

1. Wyznacz $E[(1 + \alpha)^{C_n}]$

2. Uzupełnij funkcję onGet tak aby otrzymać nieobciążony estymator liczby użyć metody onInc.

c.d.n.

Powodzenia,
Jacek Cichoń