



СОФИЙСКИ УНИВЕРСИТЕТ „СВ. КЛИМЕНТ ОХРИДСКИ“
ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА

КУРСОВ ПРОЕКТ ПО СИСТЕМИ, ОСНОВАНИ НА ЗНАНИЯ

Тема:

Разпознаване на авторство (класификатор на автори). Класификатор, който различава текстове на Алеко Константинов от текстове на Захари Стоянов.

Студенти:

Боян Иванов Дончев, група I, фн: 72029

Атанас Живков Атанасов, група I, фн: 72028

София, януари 2023 г.

1. Формулировка на задачата

Задачата реализира класифициране на текстове на Захари Стоянов и Алеко Константинов. Задачата приема текстов файл и извежда вероятността той да е написан от Захари Стоянов или Алеко Константинов. Реализацията на проекта е написана на javascript. За класификацията е необходима база от текстове на двамата автори.

2. Използвани алгоритми

Всяка дума от входния файл се сравнява с думите от базата текстове на двамата автори и се изпълнява следния алгоритъм:

- Ако се среща в текстовете на единия автор повече от колкото на другия, увеличаваме процентната вероятност на съответния автор.
- Ако се среща само в текстовете на един от двамата автори, тогава увеличаваме вероятността двойно повече.
- Ако се среща в текстовете на двамата автори по равно, променяме вероятността пропорционално на дотук достигнатия резултат.
- Ако не се среща в нито един от текстовете, вероятността не се променя.

3. Описание на програмната реализация

Функцията **countWords()** преобразува текст в честотна таблица на думите му. От текста се пропускат всички специални символи и се игнорира капитализацията на символите.

```
1 function countWords(text) {
2   var wordArr =
3     text.replace(/[.,\/#!?%&^&#x2013;{}~()]/g, '').
4     split(' ')
5
6   const res = {}
7
8   wordArr.forEach(word => {
9     word = word.toLowerCase()
10    if (res[word]) res[word] += 1
11    else res[word] = 1
12  });
13
14  return res
15 }
```

Функцията **classify()** създава честотни таблици от текстовете на двамата автори и входния текст и прилага алгоритъма описан в точка 2.

За пресмятане на процентната вероятност се използват два брояча, чиито стойности се увеличават спрямо алгоритъма.

Накрая се превръщат в проценти и се записват в глобална променлива.

```
1 function classify() {
2   const alekoMap = countWords(aleko)
3   const zahariMap = countWords(zahari)
4   const inputMap = countWords(input)
5
6   let res = {aleko: 0, zahari: 0}
7
8   Object.keys(inputMap).forEach(curr => {
9     if (!alekoMap[curr] && !zahariMap[curr]) {
10      return
11    }
12
13    if (!alekoMap[curr]) {
14      res.zahari += 2*inputMap[curr]
15      return
16    }
17
18    if (!zahariMap[curr]) {
19      res.aleko += 2*inputMap[curr]
20      return
21    }
22
23    if (alekoMap[curr] < zahariMap[curr]) res.zahari +=
    = inputMap[curr]
24    else if (alekoMap[curr] > zahariMap[curr]) res.ale
    ko += inputMap[curr]
25    else {
26      res.aleko++
27      res.zahari++
28    }
29  });
30
31  let sum = (res.aleko + res.zahari) || 1
32  output.aleko = Math.round( res.aleko / sum * 100 )
33  output.zahari = Math.round( res.zahari / sum * 100 )
34 }
```

Функцията **print()** изписва форматиран текст с резултата от класифицирането, спрямо глобалната променлива.

```
1 function print() {
2   if (output.aleko == output.zahari) {
3     console.log("\nThe input can't be classified betwee
    n Aleko Konstantinov and Zahari Stoyanov.\n")
4     return
5   }
6   console.log("\nThe input text is ${Math.max(output.al
    eko, output.zahari)}% more likely to be authored by ${o
    utput.aleko > output.zahari ? "Aleko Konstantinov then
    by Zahari Stoyanov" : "Zahari Stoyanov then by Aleko K
    onstantinov"}.\n")
7 }
```

4. Примери, илюстриращи работата на програмната система

Ако Алеко е написал „atanas atanas atanas”, а Захари е написал „boyan boyan boyan” и входният текст е:

- „atanas atanas boyan” → “The input text is 67% more likely to be authored by Aleko Konstantinov then by Zahari Stoyanov.”
- “boyan” → “The input text is 100% more likely to be authored by Zahari Stoyanov then by Aleko Konstantinov.”
- “nasko” → “The input can't be classified between Aleko Konstantinov and Zahari Stoyanov.”
- “boyan atanas” → “The input can't be classified between Aleko Konstantinov and Zahari Stoyanov.”

5. Литература

За база от текстове на двамата автори са използвани текстове от <https://chitanka.info/>.