

# Алгоритм

## Обозначения

Пусть:

Функция  $\text{intersect}(A, B)$  возвращает элементы, которые содержатся одновременно и в массиве  $A$ , и в массиве  $B$ .

Функция  $\text{size}(A)$  возвращает количество элементов, которые содержатся в массиве  $A$ .

Функция  $\text{lcs}(A, B)$  возвращает наибольшую общую подпоследовательность слов в  $A$  и  $B$

$\text{sent}$  - массив предложений.

$\text{sent}[i].\text{nGrams}$  - массив  $n$ -грамм в  $i$ -ом предложении.

$\text{classes}$  - массив классов

$\text{classes}[i].\text{nGrams}$  - массив  $n$ -грамм в  $i$ -ом классе.

$\text{classes}[i].\text{sent}$  - массив предложений в  $i$ -ом классе.

---

### Algorithm 1 Поиск неточных повторов

---

```
1: for  $i = 1$  to  $\text{size}(\text{sent})$  do
2:    $\text{curSent} = \text{sent}[i]$ 
3:    $\text{bestOverlap} = 0$ 
4:    $\text{bestClass} = \text{NULL}$ 
5:   for  $j = 1$  to  $\text{size}(\text{classes})$  do
6:      $\text{curClass} = \text{classes}[j]$ 
7:      $\text{curIntersect} = \text{intersect}(\text{curSent.nGrams}, \text{curClass.nGrams})$ 
8:      $\text{curOverlap} = \text{size}(\text{curIntersect}) / \text{size}(\text{curSent.nGrams})$ 
9:     if  $\text{curOverlap} > \text{bestOverlap}$  then
10:       $\text{bestOverlap} = \text{curOverlap}$ 
11:       $\text{bestClass} = \text{curClass}$ 
12:     end if
13:   end for
14:   if  $\text{bestOverlap} < 0.5$  then
15:     Создать новый класс  $\text{newClass}$ 
16:      $\text{newClass.nGrams} += \text{curSent.nGrams}$ 
17:      $\text{newClass.sent} += \text{curSent}$ 
18:   else
19:      $\text{bestClass.nGrams} += \text{curSent.nGrams}$ 
20:      $\text{bestClass.sent} += \text{curSent}$ 
21:   end if
22: end for
```

---