
Skript pro extrakci klíčových slov

Artem Gorodilov

1. Abstrakt

Skript přijímá text v kódování UTF-8 jako standardní vstup a jako standardní výstup poskytuje textový soubor s klíčovými slovy seřazenými podle důležitosti (četnost výskytu v textu a celkový počet výskytů).

2. Popis skriptu

Skript je napsán v jazyce Perl v5.34.0 [1] (funguje i v novějších verzích).

Skript je rozdělen do několika částí:

1. Import požadovaných modulů, nastavení kódování UTF-8 a definování oddělovače writerů.
2. Načtení souboru ze standardního vstupu, převod textu na seznam slov a změna velikosti všech slov na malá písmena.
3. Import datového souboru se stop slovy, filtrování seznamu slov od stop slov a čísel.
4. Výpočet četnosti výskytu jednotlivých filtrovaných klíčových slov v textu, seřazení klíčových slov podle četnosti výskytu a seřazení podle celkového počtu výskytů v textu.
5. Výstup seznamu klíčových slov.

2.1. Nastavení skriptu

Nastavení skriptu je v kódu definováno následujícími řádky:

```
1 use strict;
2 use warnings;
3
4 use open qw(:std :utf8);
```

Skript používá moduly `strict` a `warnings` k zajištění správné činnosti skriptu a k zobrazení varování v případě chyb. Modul `open` slouží k nastavení kódování vstupních a výstupních souborů na UTF-8.

2.2. Čtení souboru

Skript načte text ze standardního vstupu a převede jej na seznam slov [2]. Vstupní text je umístěn ve složce *input/*.

Text je rozdělen na slova pomocí mezer, interpunkčních znamének a zalomení řádků. Velká písmena všech slov jsou změněna na malá.

```
1 my $text = <STDIN>;
2
3 my @words = $text =~ /(\w+)/g;
4
5 foreach my $word (@words) {
6     $word = lc $word;
7 }
```

2.3. Filtrování stop slov a čísel

Skript načte stop slova z datového souboru [3] a odfiltruje seznam slov od stop slov a čísel.

```
1 my @stopwords;
2 open my $stopwords_file,
3     '<', 'stopwords-cs.txt' or die "Cannot
4     open stopwords-cs.txt
5     : $!";
6
7 @stopwords = split(/\s+/, <$stopwords_file>);
8
9 my %stopwords_hash = map
10     { $_ => 1 }
11     @stopwords;
12
13 my @filt_words = grep {
14     !$stopwords_hash{$_}
15     && !/\d/ } @words;
```

2.4. Výpočet četnosti výskytu a třídění

Skript vypočítá četnost výskytu každého filtrovaného klíčového slova v textu a seřadí klíčová slova podle četnosti výskytu [4] a podle celkového počtu výskytů v textu.

Prahová hodnota pro počet výskytů klíčového slova v definici podle argumentu uvedeného v inicializačním řádku kódu [5] (vysvětleno v kapitole "3. Dokumentace").

```
1 my %word_counter;
2 foreach my $word (
3     @filt_words) {
4     $word_counter{$word}++;
5 }
6 my @sorted_words = sort {
7     $word_counter{$b} <=>
8     $word_counter{$a} } keys
9     %word_counter;
10 my $frequency_threshold =
    $ARGV[0];
11 my @keywords = grep {
    $word_counter{$_} >
    $frequency_threshold }
    @sorted_words;
```

2.5. Výstup seznamu klíčových slov

Skript vypíše seznam klíčových slov na standardní výstup. Klíčová slova jsou uvedena ve formě seznamu, přičemž každé klíčové slovo je na samostatném řádku.

```
1 foreach my $keyword (
2     @keywords) {
3     print "$keyword\n";
4 }
```

Odkazy

- [1] Perl website. www.perl.com/
- [2] PerlScript to take certain words from text file. <https://stackoverflow.com/questions/66495616/perlscript-to-take-certain-words-from-text-file>
- [3] Stopwords ISO CZ. <https://github.com/stopwords-iso/stopwords-cs/blob/master/stopwords-cs.txt>
- [4] How do I sort by frequency of a value? <https://stackoverflow.com/questions/4519979/how-do-i-sort-by-frequency-of-a-value>
- [5] How can I pass command-line arguments to a Perl program? <https://stackoverflow.com/questions/361752/how-can-i-pass-command-line-arguments-to-a-perl-program>

3. Dokumentace

Skript se spouští z příkazového řádku následujícím příkazem:

```
perl klicova_slova.pl 10 <
input/text.txt > output/keywords.txt
```

Skript přijímá jeden argument - prahovou hodnotu pro počet výskytů klíčového slova v textu. Tento počet závisí na velikosti textu. Pokud je text velký, měla by být prahová hodnota nastavena na vyšší hodnotu. Pokud je text malý, měla by být prahová hodnota nastavena na nižší hodnotu.

Skript načte text ze složky *input/* a zapíše seznam klíčových slov do složky *output/*.

4. Diskuse

Skript funguje správně a poskytuje očekávaný výsledek. Nicméně jasnost klíčových slov závisí především na údajích o stopwords. Uživatel může do datového souboru napsat vhodnější stopslova pro konkrétní text nebo najít jiný zdroj stopslov.