

# 1 Popis riešenia

Nasledujúci skript v jazyku Perl slúži na extrakciu kľúčových slov z textového súboru a ich zoradenie podľa dôležitosti.

```
1 use strict;
2 use warnings;
3
4 #!/usr/bin/perl
5
6 use strict;
7 use warnings;
8 use Lingua::Stem::Snowball;
9 use Getopt::Long;
10 use JSON;
11
12 # Parse command line options
13 my $stopwords_dir = "stopwords";
14 my $lang = "en";
15 my $num_keywords = 10;
16 my $output_format = "text";
17 my $help = 0;
18 GetOptions(
19     'stopwords=s' => \$stopwords_dir,
20     'lang=s' => \$lang,
21     'num_keywords=i' => \$num_keywords,
22     'output_format=s' => \$output_format,
23     'help' => \$help,
24 ) or die "Error parsing command line options\n";
25
26 # Check for help option
27 if ($help) {
28     print "Usage: perl script.pl [options]\n";
29     print "Extracts and ranks keywords from text\n";
30     print "provided on standard input.\n";
31     print "Options:\n";
32     print "  --stopwords DIR      Path to stop words\n";
33     print "                        folder (default: stopwords)\n";
34     print "  --lang LANG          Language for\n";
35     print "                        stemming and stop words (default: en)\n";
36     print "  --num_keywords NUM   Number of keywords\n";
37     print "                        to output (default: 10)\n";
38     print "  --output_format FMT  Output format (text\n";
39     print "                        , json, csv) (default: text)\n";
40     print "  --help              Show this help\n";
41     print "message\n";
```

```

36     exit;
37 }
38
39 # Load stop words from file
40 my $stopwords_file = "$stopwords_dir/$lang.txt";
41 my @stopwords;
42 if (-e $stopwords_file) {
43     open(my $fh, "<", $stopwords_file) or die "Cannot
44         open file $stopwords_file: $!";
45     while (my $line = <$fh>) {
46         chomp $line;
47         push @stopwords, $line;
48     }
49     close($fh);
50 } else {
51     warn "Warning: Stop words file '$stopwords_file'
52         not found. Proceeding without stop words.\n";
53 }
54
55 # Initialize stemmer
56 my %lang_map = (
57     'sk' => 'sk',
58     'cz' => 'cs',
59     'en' => 'en',
60 );
61 my $stemmer_lang = $lang_map{$lang} || 'en';
62 my $stemmer = Lingua::Stem::Snowball->new( lang =>
63     $stemmer_lang );
64
65 # Load input text from standard input
66 my $text = do { local $/; <STDIN> };
67 if (!defined $text || length($text) == 0) {
68     die "Error: No input provided on standard input\n"
69         ;
70 }
71
72 # Tokenize text into individual words
73 my @words = split /\s+/, $text;
74
75 # Remove punctuation and convert to lowercase
76 foreach my $word (@words) {
77     $word =~ s/[:punct:]]//g;
78     $word = lc $word;
79 }
80
81 # Remove stop words and apply stemming

```

```

78 @words = grep { my $w = $_; !grep { $_ eq $w }
    @stopwords } @words;
79 foreach my $word (@words) {
80     $word = $stemmer->stem($word);
81 }
82
83 # Calculate word frequency
84 my %word_freq;
85 foreach my $word (@words) {
86     $word_freq{$word}++;
87 }
88
89 # Calculate word importance using TextRank algorithm
90 my %scores;
91 my $damping_factor = 0.85;
92 my $max_iterations = 50;
93 my $convergence_threshold = 0.0001;
94 my $num_words = scalar @words;
95
96 foreach my $word (keys %word_freq) {
97     $scores{$word} = 1;
98 }
99
100 for (my $i = 0; $i < $max_iterations; $i++) {
101     my %new_scores;
102
103     foreach my $word (keys %word_freq) {
104         $new_scores{$word} = (1 - $damping_factor) /
            $num_words;
105
106         foreach my $adj_word (keys %word_freq) {
107             next if $word eq $adj_word;
108
109             my $adj_frequency = $word_freq{$adj_word};
110             my $adj_links = $word_freq{$word};
111
112             if ($adj_frequency > 0) {
113                 $new_scores{$word} += ($damping_factor
                    * $scores{$adj_word} * $adj_links)
                    / $adj_frequency;
114             }
115         }
116     }
117
118     my $max_diff = 0;
119     foreach my $word (keys %scores) {

```

```

120         my $diff = abs($scores{$word} - $new_scores{
121             $word});
122         if ($diff > $max_diff) {
123             $max_diff = $diff;
124         }
125     }
126     last if $max_diff < $convergence_threshold;
127
128     %scores = %new_scores;
129 }
130
131 # Normalize scores
132 my $max_score = (sort { $b <=> $a } values %scores)
133     [0];
134 foreach my $word (keys %scores) {
135     $scores{$word} /= $max_score;
136 }
137
138 # Sort keywords by importance
139 my @sorted_words = sort { $scores{$b} <=> $scores{$a}
140     } keys %scores;
141
142 # Output keywords to standard output
143 if ($output_format eq "json") {
144     my @output_data;
145     foreach my $word (@sorted_words[0..$num_keywords
146         -1]) {
147         if (defined $word) {
148             push @output_data, { word => $word, score
149                 => $scores{$word} };
150         }
151     }
152     print to_json(\@output_data, { pretty => 1 });
153 } elsif ($output_format eq "csv") {
154     print "word,score\n";
155     foreach my $word (@sorted_words[0..$num_keywords
156         -1]) {
157         if (defined $word) {
158             print "$word,$scores{$word}\n";
159         }
160     }
161 } else {
162     foreach my $word (@sorted_words[0..$num_keywords
163         -1]) {
164         if (defined $word) {

```

```

159         print "$word\n";
160     }
161 }
162 }

```

1. Skript načíta textový súbor ako vstup zo štandardného vstupu.
2. Načítaný text sa tokenizuje na jednotlivé slová.
3. Pre každé slovo sa vykonajú úpravy, ako odstránenie interpunkcie a prevod na malé písmená.
4. Následne sa odstránia stop slová pre zvolený jazyk.
5. Pre každé slovo sa vypočíta frekvencia výskytu.
6. Dôležitosť slov sa vypočíta pomocou algoritmu TextRank.
7. Dôležitosť slov sa normalizuje a zoradia sa podľa hodnoty.
8. Nakoniec sa kľúčové slová vypíšu na štandardný výstup, každé slovo na samostatnom riadku.

Skript využíva moduly `Lingua::Stem::Snowball` pre stemming slov a `Getopt::Long` pre spracovanie príkazového riadku.

## 2 Inštalácia modulov

Pred použitím skriptu sa uistite, že máte nainštalované moduly `Lingua::Stem::Snowball` a `Getopt::Long`. Ak nie sú nainštalované, môžete ich nainštalovať pomocou správcu balíkov CPAN.

```

1 cpan Lingua::Stem::Snowball
2 cpan Getopt::Long

```

## 3 Použitie skriptu

Na spustenie skriptu použite nasledujúci príkaz:

```

1 perl klicova_slova.pl < input.txt

```

Kde `klicova_slova.pl` je názov skriptu a `input.txt` je vstupný textový súbor.