Laboratorium Analizy Procesów Uczenia.

Data wykonania ćwiczenia:	24.05.2024
Rok studiów:	1
Semestr:	1
Grupa studencka:	1b
Grupa laboratoryjna:	-
Ćwiczenie nr	6

Temat: Problemy NLP w uczeniu maszynowym.

Osoby wykonujące ćwiczenia:

1. Gracjan Wackermann

Katedra Informatyki i Automatyki

1. Cel ćwiczenia:

Celem są 1. Wstępna analiza tekstów z pomocą list częstotliwości, chmur słów, n-gramów. 2. Konstruowanie n-gramów i grafów.

2. Zadanie do wykonania:

Zadanie dotyczy analizy tekstu, w tym list częstotliwości słów, budowanie chmury słów, kojarzeń i innych.

```
3. https://en.wikipedia.org/wiki/Poetry
- Wariant nr. 3 -
```

Uzyskany kod:

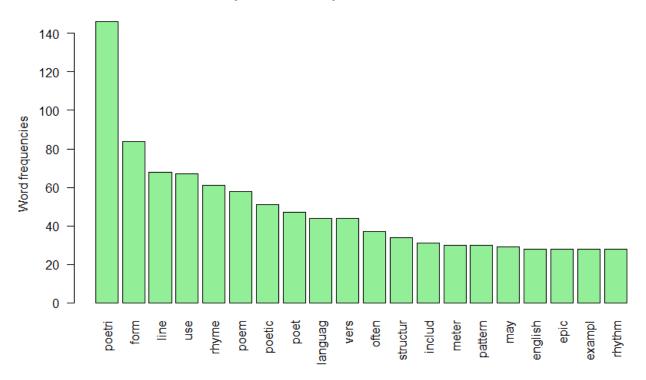
```
# Ładowanie niezbędnych paczek
library("tm")
library("SnowballC")
library("wordcloud")
library("RColorBrewer")
library("syuzhet")
library("ggplot2")
library("tidytext")
library("igraph")
library("ggraph")
library("rvest")
library("dplyr")
library("tidyr")
# Pobranie tekstu z URL
url <- 'https://en.wikipedia.org/wiki/Poetry'</pre>
webpage <- read html(url)</pre>
text <- webpage %>%
 html_nodes("p") %>%
 html text() %>%
 paste(collapse = " ")
# Przekształcenie tekstu do obiektu Corpus
TextDoc <- VCorpus(VectorSource(text))</pre>
# "Wyczyszczanie" tekstu
toSpace <- content transformer(function (x , pattern ) gsub(pattern, " ", x))
TextDoc <- tm map(TextDoc, toSpace, "/")</pre>
TextDoc <- tm map(TextDoc, toSpace, "@")
TextDoc <- tm_map(TextDoc, toSpace, "\\|")</pre>
TextDoc <- tm_map(TextDoc, content_transformer(tolower))</pre>
TextDoc <- tm map(TextDoc, removeNumbers)</pre>
TextDoc <- tm map(TextDoc, removeWords, stopwords("english"))</pre>
TextDoc <- tm map(TextDoc, removePunctuation)</pre>
TextDoc <- tm_map(TextDoc, stripWhitespace)</pre>
TextDoc <- tm map(TextDoc, stemDocument)</pre>
# Budowanie macierzy dokumentu
TextDoc dtm <- TermDocumentMatrix(TextDoc)</pre>
dtm m <- as.matrix(TextDoc dtm)</pre>
dtm v <- sort(rowSums(dtm m), decreasing = TRUE)</pre>
dtm d <- data.frame(word = names(dtm v), freq = dtm v)</pre>
# Wyświetlanie 5 najczęstszych słów
```

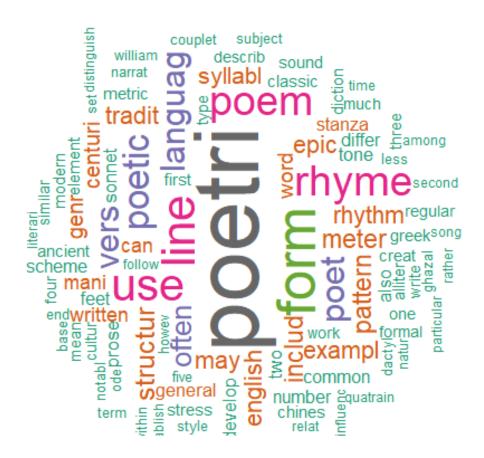
```
head(dtm d, 5)
# Wykres słupkowy najczęstszych słów
barplot(dtm d[1:20,] $freq, las = 2, names.arg = dtm d[1:20,] $word,
        col = "lightgreen",
        main ="Top 20 most frequent words in the article",
        ylab = "Word frequencies")
# Generowanie chmury słów
set.seed(1234)
wordcloud(words = dtm d$word, freq = dtm d$freq, scale=c(5, 0.5),
          min.freq = 1, max.words=100, random.order=FALSE,
          rot.per=0.40, colors=brewer.pal(8, "Dark2"))
# Kojarzenia słów
findAssocs(TextDoc dtm, terms = c("poetry", "form", "literary", "work"), corlimit =
# Analiza sentymentu
syuzhet vector <- get sentiment(text, method="syuzhet")</pre>
summary(syuzhet_vector)
bing_vector <- get_sentiment(text, method="bing")</pre>
summary(bing_vector)
afinn vector <- get sentiment(text, method="afinn")</pre>
summary (afinn vector)
# Analiza emocji
d <- get nrc sentiment(as.vector(dtm d$word))</pre>
td <- data.frame(t(d))
td_new <- data.frame(rowSums(td[1:56]))</pre>
names(td new)[1] <- "count"
td new <- cbind("sentiment" = rownames(td new), td new)
rownames(td new) <- NULL
td new2 <- td new[1:8,]
# Wykres liczby słów związanych z każdym uczuciem
ggplot(td new2, aes(x = sentiment, y = count, fill = sentiment)) +
  geom bar(stat = "identity") +
  ggtitle("Survey sentiments") +
  ylab("count")
# Bigramy
text df <- tibble(line = 1, text = text)</pre>
tidy_text <- text_df %>%
  unnest tokens (word, text)
data(stop words)
tidy text <- tidy text %>%
  anti_join(stop_words, by = "word")
# Lista częstotliwości słów
tidy_text %>%
 count (word, sort = TRUE)
# Bigramy
text bigrams <- text df %>%
  unnest tokens (bigram, text, token = "ngrams", n = 2)
text bigrams %>%
 count(bigram, sort = TRUE)
bigrams separated <- text bigrams %>%
 separate(bigram, c("word1", "word2"), sep = " ")
bigrams filtered <- bigrams separated %>%
  filter(!word1 %in% stop_words$word) %>%
```

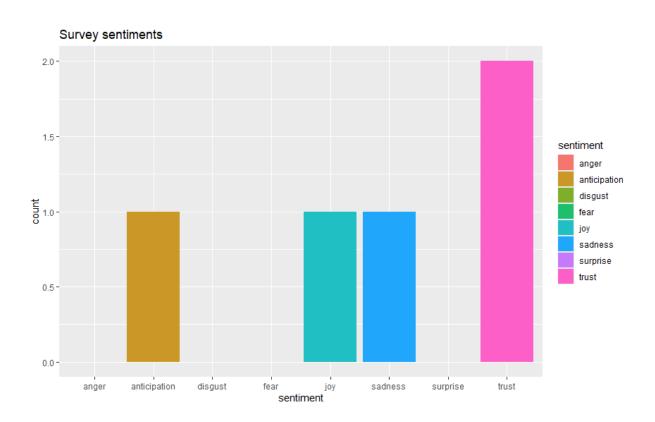
```
filter(!word2 %in% stop words$word)
bigram counts <- bigrams filtered %>%
  count(word1, word2, sort = TRUE)
bigrams united <- bigrams filtered %>%
  unite(bigram, word1, word2, sep = " ")
# Konstruowanie grafów
bigram_graph <- bigram_counts %>%
  filter(word1 == "poetry" | word2 == "poetry") %>%
  graph_from_data_frame()
# Wyświetlanie grafów
ggraph(bigram_graph, layout = "fr") +
  geom_edge_link(aes(edge_alpha = n), show.legend = TRUE, arrow = arrow(type =
"closed", length = unit(.\overline{15}, "inches")), end cap = circle(.07, 'inches')) +
 geom_node_point(color = "lightblue", size = 5) +
  geom_node_text(aes(label = name), position = "identity") +
  theme void()
```

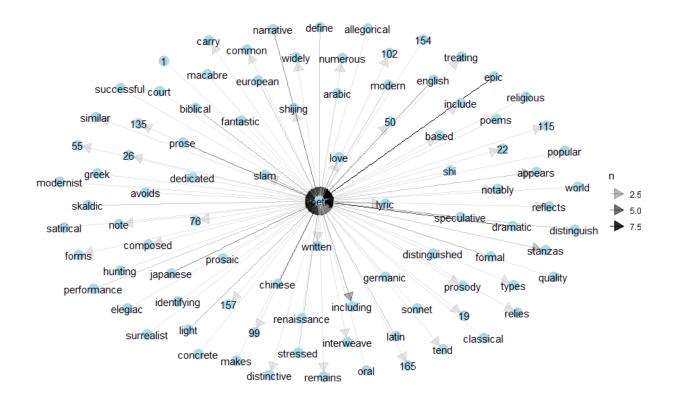
Wyniki z konsol:











3. Wnioski:

- Na podstawie wyników uzyskanych z analizy tekstu artykułu na Wikipedii dotyczącego poezji, można wyciągnąć następujące wnioski:

Najczęściej występujące słowa to: poetry, form, line, use, rhyme

Wskazuje to na to, że artykuł skupia się na definicji i strukturze poezji, omawiając różne formy, linie, użycia oraz rymy w poezji. Chmura słów wizualizuje te wyniki, pokazując, że słowo "poetry" jest zdecydowanie najczęściej używanym słowem w tekście, co nie jest zaskakujące, biorąc pod uwagę temat artykułu.

Próba znalezienia kojarzeń słów dla poetry, form, literary, work przy korelacji 0.5 nie zwróciła żadnych wyników (numeric(0)). Może to oznaczać, że żadne z tych słów nie są wystarczająco silnie powiązane z innymi słowami w tekście przy ustawionym progu korelacji.

Ogólnie rzecz biorąc, analiza tekstu artykułu na Wikipedii dotyczącego poezji pokazuje, że jest on skoncentrowany na omawianiu definicji, form i struktur poezji. Artykuł jest napisany w neutralnym tonie, co jest zgodne z oczekiwaniami wobec treści encyklopedycznych. Najczęściej używane słowa i bigramy są zgodne z tematem, a analiza emocji i sentymentu potwierdza obiektywność i informacyjny charakter tekstu.

Link do repozytorium: https://github.com/fireinx/apu