



**Kolegium Nauk Przyrodniczych  
Uniwersytet Rzeszowski**

**Przedmiot:**

**Wybrane Zagadnienia Współczesnej  
Informatyki**

**Tytuł projektu:**

***Aplikacja webowa do akwizycji i  
analizy danych z kanału RSS portalu  
gofin.pl (prawo pracy).***

**Wykonał:**

**Kamil Filar**, Informatyka, rok III, lab1

**Prowadzący: dr hab. inż. Krzysztof Pancerz, prof. UR**

**Rzeszów 2020**

## Spis treści:

1. Opis serwisu/aplikacji internetowej .....	3
1.1 Główne funkcjonalności .....	3
2. Szczegółowy opis funkcjonalności/wyglądu programu .....	4
2.1 Graficzny interfejs użytkownika .....	5
2.2 Apache Solr .....	7
2.3 Występowanie danego słowa w wybranej tabeli .....	7
2.4 Lematyzacja .....	8
2.5 Przygotowanie danych .....	9
2.6 Chmura słów .....	9
2.7 Częstotliwość występowania danego słowa .....	11
2.8 Klasteryzacja danych .....	12

## 1. Opis serwisu/aplikacji internetowej

Nazwa aplikacji: Aplikacja webowa do akwizycji i analizy danych z kanału RSS portalu gofin.pl (prawo pracy).

Aplikacja została stworzona z wykorzystaniem środowiska R oraz Apache Solr. Aplikacja webowa oparta jest o zaawansowane GUI (możliwość ustawiania parametrów) zrealizowane z wykorzystaniem pakietu Shiny. Ponadto posiada funkcjonalności takie jak:

- Możliwość pobierania wiadomości i magazynowania ich w Apache Solr.
- Możliwość analizy częstości występowania słów i wizualizacji chmury słów w wiadomościach zmagazynowanych w Apache Solr.
- Możliwość klasteryzacji wiadomości zmagazynowanych w Apache Solr oraz wizualizacji klastrów.

### 1.1 Główne funkcjonalności

1. Pobieranie danych z kanału RSS.
2. Magazynowanie ich w Apache Solr.
3. Pobieranie danych z Apache Solr.
4. Lematyzacja tekstu.
5. Analiza występowania danego słowa w zadanym zakresie danych.
6. Tworzenie chmury słów.
7. Analiza częstości występowania słów.
8. Klasteryzacja danych.

## 2. Szczegółowy opis funkcjonalności/wyglądu programu

### 2.1 Graficzny interfejs użytkownika

Graficzny interfejs użytkownika został zrealizowany z wykorzystaniem pakietu Shiny. Jego szczegółowy wygląd został zaprezentowany poniżej:

```
150 #UI aplikacji
151 ui <- fluidPage(
152   includeCSS("style.css"),
153   # Application title
154   tags$div(class="Title",
155     tags$p("Aplikacja webowa do akwizycji i analizy danych"),
156     tags$p("z kanału RSS portalu"),
157     tags$p("Gofin.pl - prawo pracy")
158   ),
159   tags$div(class="Line1"),
160   tags$div(class="Line2"),
161   sidebarLayout(
162     sidebarPanel(
163       tags$p("Procentowa zawartosc konkretnego slowa", class="SubTitle"),
164       selectInput(inputId="selectTab2",label=h4("wybierz dane"),choices = c("Tytuł"="Tytuł","Zawartosc"="Zawartosc"),
165         selected = "Tytuł",multiple = F),
166       textInput("słowo", h4("Słowo do wyszukania: ")),
167       selectInput(inputId="color1",label=h4("wybierz kolor"),
168         choices = c("Czerwony"="Czerwony","Niebieski"="Niebieski","Zielony"="Zielony"),
169         selected = "Zielony",multiple = F),
170       actionButton("buttonA1","Start"),
171     ),
172     mainPanel(
173       plotOutput("wykresA1")
174     )
175   ),
176   ),
177   ),
178   ),
179   ),
```

Rys. 1

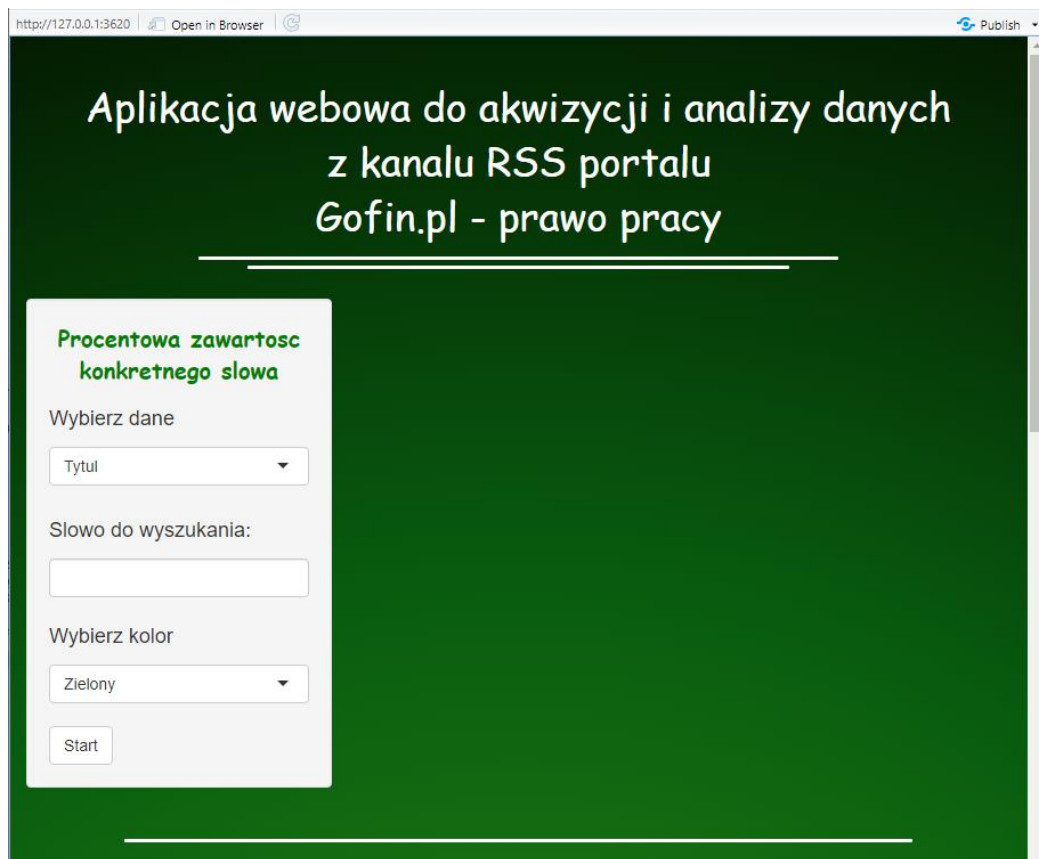
```
194 numericInput("grid", h4("Odstep miedzy literami"), value = 5, min = 0, max = 25, step = 5),
```

Rys. 2

Rys. 1 przedstawia początek kodu UI. W linii 153 znajduje się załączony plik zawierający CSS, który w stanowi o wyglądzie aplikacji. W liniach od 156 do 160 znajduje się tytuł projektu, a poniżej tego są dwie linie służące do separacji między blokami w projekcie. Poniżej znajduje się pierwszy panel wchodzący w skład aplikacji. W aplikacji znajdują się cztery takie panele umożliwiające użytkownikowi analizę danych. Składają się z takich elementów jak:

- selectInput – służy do wybierania danych (168 linia Rys. 1)
- textInput – służy do wczytywania danych tekstowych (170 linia Rys. 1)
- actionButton – przycisk, który służy do uruchomienia funkcji (174 linia Rys. 1)
- numericInput – służy do wczytywania liczb z zakresu (194 linia Rys. 2)

Wygląd UI został przedstawiony na Rys. 3, Rys. 4, Rys. 5 poniżej:



Rys. 3

http://127.0.0.1:3620 Open in Browser Publish

### Chmura slow

Wybierz dane

Tytul

Wybierz ksztalt

Gwiazda

Rozmiar

0,5

Odstep miedzy literami

5

Start

---

### Analiza czestosci wystepowania slowa

Wybierz dane

Tytul

Slowo do wyszukania:

Start

Wczytane slowo:

Ile razy wystepuje:

Rys. 4

### Klasteryzacja

Wybierz dane

Tytul

Metryka

euclidean

Ilosc klastrow

2

Start

---

Kamil Filar, Informatyka, rok III, lab1

Rys. 5

## 2.2 Apache Solr

```
19 #POBRANIE DANYCH ZE STRONY
20
21 aktualnosci <- tidyfeed(feed="http://www.rss.gofin.pl/prawopracy.xml")
22 #view(aktualnosci)
23 #WCZYTANIE WYBRANYCH DANYCH DO OBIEKTOW (SENSOWNYCH)
24 tytul_data <- aktualnosci$entry_title
25 zawartosc_data <- aktualnosci$entry_content
26 URL_data <- aktualnosci$entry_url
27
28 #view(aktualnosci) PODGLAD DANYCH (POMOCNICZO)
29
30 #DZIALANIE W SOLR
31 polaczenie <- SolrClient$new(host="127.0.0.1", port = 8983, path = "/solr/ProjektWZWI/select");
32
33 dane = data.frame(matrix(ncol=4, nrow = 10));
34 colnames(dane)[1] <- "id";
35 colnames(dane)[2] <- "Tytul";
36 colnames(dane)[3] <- "Zawartosc";
37 colnames(dane)[4] <- "URL"
38
39 dane$id<-c(1,2,3,4,5,6,7,8,9,10);
40 dane$Tytul<-c(tytul_data);
41 dane$Zawartosc<-c(zawartosc_data);
42 dane$URL<-c(URL_data);
43
44 solrium::add(x=dane, conn=polaczenie, name="ProjektWZWI", commit=TRUE);
45
```

Rys. 6

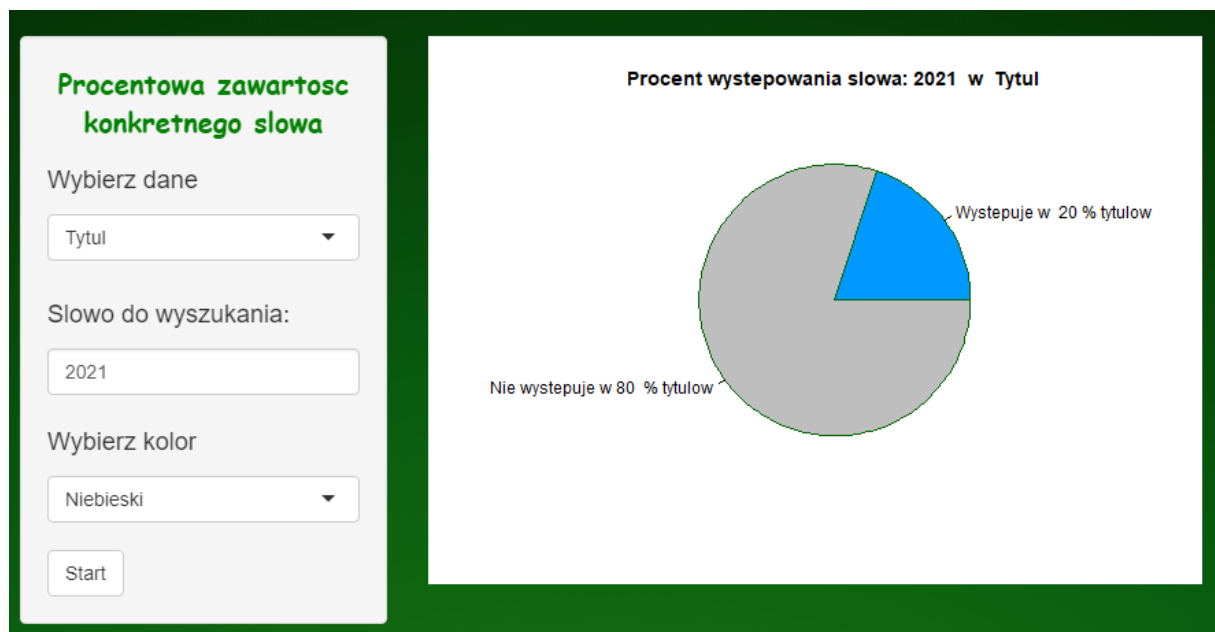
Na Rys. 6 została przedstawiona obsługa Apache Solr. W linii 21 wczytuje dane z kanału RSS za pomocą funkcji `tidyfeed`, następnie w liniach od 24-26 przypisuje obiektom (w celu ułatwienia sobie pracy) dane wczytane z kanału RSS, następnie w linii 34 tworzy połączenie z węzłem ProjektWZWI, tworzy obiekt do przechowywania danych i w linii 44 dodaje do Apache Solr moje dane. Pobieranie danych z Solr-a zostało opisane przy funkcjach, które to wykorzystują.

## 2.3 Występowanie danego słowa w wybranej tabeli

```
49 #FUNKCJA ZWRACA ILE % WIERZSY TABELI tabela2 ZAWIERA WCZYTANE slowo
50 analiza1<-function(slowo, tabela2){
51   dane1<-solr_search(conn = polaczenie, params = list(q=paste(tabela2,":",slowo),fl=paste(tabela2), rows=-1));
52   sprawdz_duplikaty <- duplicated(dane1)
53   view(sprawdz_duplikaty)
54   dane_koncowe <- unique(dane1)
55   ilosc_kolumn <- nrow(dane_koncowe)
56
57   return(ilosc_kolumn);
58 }
```

Rys. 7

Funkcja `analiza1` służy do zliczenia ilości kolumn, w jakich zawiera się dane słowo. Funkcja przyjmuje dwa parametry: *slowo* oraz *tabela2*. Parametr *slowo* jest wczytywany w UI przez użytkownika i za pomocą logiki serwera zostaje przekazane do funkcji jako parametr przy wywołaniu, podobnie jest z drugim parametrem. Funkcja sprawdza czy są duplikaty (linia 52-54) w wczytanych danych z Apache Solr (samo wczytywanie jest zrealizowane za pomocą funkcji `solr_search` w linii 51) i wybiera tylko unikatowe wartości. Zwrócone wartości to ilość kolumn. Na Rys. 8 (poniżej) został zaprezentowany wygląd wyniku działania funkcji:



Rys. 8

## 2.4 Lematyzacja tekstu

```
59 ~ 1 lematyzacja<-function(tekst){
60   #FUNKCJA ODPOWIEDZIALNA ZA LEMATYZACJE TEKSTU (KOPIA Z LAB4) sprawdzanie duplikatow!
61   parametry<-list(lpmn="any2txt|wcrft2", text=tekst, user="Filarkamil04@gmail.com");
62   odpowiedz<-POST("http://ws.clarin-pl.eu/nlprest2/base/process", body=parametry, encode="json", verbose());
63   zawartosc<-content(odpowiedz, "text", encoding="UTF-8");
64   xml<-xmlParse(zawartosc, encoding="UTF-8");
65   slowa<-xpathSApply(xml, '//chunkList/chunk/sentence/tok/lex/base', xmlValue, encoding="UTF-8");
66   return(paste(slowa, collapse=" "));
67 }
68 ~ }
```

Rys. 9

Funkcja służąca do lematyzacji tekstu została przedstawiona na Rys. 9. Jest to funkcja, która została przedstawiona w laboratorium 4 podczas zajęć.



## 2.5 Przygotowanie danych do analizy

```
70 #FUNKCJA ODPOWIEDZIALNA ZA PRZYGOTOWANIE DANYCH DO ANALIZY
71 przygotowanie_danych<-function(dane){
72   stop<-as.vector(unlist(read.csv(file="stop_words_pl.txt", header=FALSE, sep=",", fileEncoding="UTF-8")));
73
74   dokumenty<-Corpus(VectorSource(str_enc_toutf8(dane)));
75
76   dokumenty<-tm_map(dokumenty, removePunctuation, preserve_intra_word_dashes=TRUE);
77   dokumenty<-tm_map(dokumenty, removeNumbers);
78   dokumenty<-tm_map(dokumenty, removeWords, stop);
79
80   usun.znaki<-function(x) gsub("[â€Žâ€Žâ€Ž]", "", x);
81   dokumenty<-tm_map(dokumenty, usun.znaki);
82
83   for(d in 1:length(dokumenty))
84   {
85     dokumenty[[d]]$content<-lematyzacja(dokumenty[[d]]$content);
86     dokumenty[[d]]$content<-stri_enc_toutf8(dokumenty[[d]]$content);
87   }
88
89   tdm1<-TermDocumentMatrix(dokumenty);
90   m1<-as.matrix(tdm1);
91   v<-sort(rowSums(m1),decreasing=TRUE);
92   d<-data.frame(words=names(v), freq=v);
93
94   return(d);
95 }
```

Rys. 10

Na Rys. 10 została przedstawiona funkcja (wykorzystana z laboratorium 4), która przyjmuje parametr wejściowy *dane*, które są poddawane działaniom wewnątrz funkcji. Funkcja zwraca tabelę z uporządkowanymi danymi w kolumnie słowa oraz częstotliwość ich występowania.

## 2.6 Chmura słów

```
97 #FUNKCJA ODPOWIEDZIALNA ZA TWORZENIE CHMURY SŁÓW
98 chmura<-function(tabela){
99   dane<-solr_search(conn = polaczenie, params = list(q=paste(tabela,":*"),fl=paste(tabela), rows=-1));
100   sprawdz_duplikaty_chmura <- duplicated(dane)
101   view(sprawdz_duplikaty)
102   dane_koncowe <- przygotowanie_danych(dane)
103   return(dane_koncowe);
104 }

98 #FUNKCJA ODPOWIEDZIALNA ZA TWORZENIE CHMURY SŁÓW
99 chmura<-function(tabela){
100   dane<-solr_search(conn = polaczenie, params = list(q=paste(tabela,":*"),fl=paste(tabela), rows=-1));
101   sprawdz_duplikaty_chmura <- duplicated(dane)
102   view(sprawdz_duplikaty)
103   dane_z_solra <- przygotowanie_danych(dane)
104   dane_koncowe <- unique(dane_z_solra)
105   return(dane_koncowe);
106 }
```

Rys. 11

Funkcja *chmura* przyjmuje parametr *tabela* wczytany z UI przez użytkownika. Parametr wskazuje na tabelę z której solr pobiera dane do dalszego przetworzenia (linia 100 Rys.11). Następnie dane są sprawdzane pod kątem występowania duplikatów, a kolejnym krokiem jest zwrócenie unikatowych danych końcowych, które trafiają do przetworzenia w logice serwera (Rys. 12) gdzie w linii 300 została zastosowana funkcja tworząca chmurę słów (*worldcloud2*).

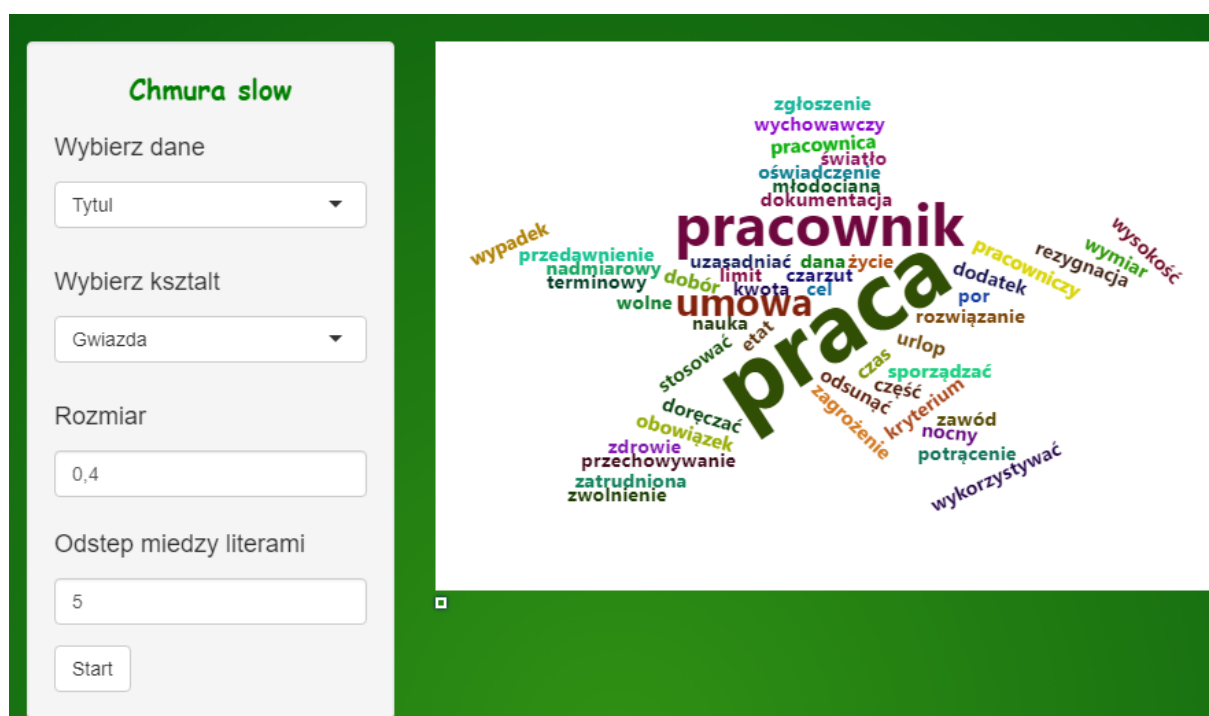
```

286 ▾ observeEvent(input$buttonA2,{
287
288 ▾   output$wykresA2 <- renderwordcloud2({
289
290     dane <- chmura(input$selectTab1);
291     |
292 ▾     if(input$selectShape=="Gwiazda"){
293       varshape = "star"
294 ▾     }else if(input$selectShape=="Diament"){
295       varshape = "diamond"
296 ▾     }else if(input$selectShape=="Trojkat"){
297       varshape = "triangle-forward"
298 ▾     }
299
300     wordcloud2(dane, size=input$size, gridsize = input$grid,
301               color = "random-dark", backgroundColor = "white",
302               shape = varshape)
303 ▾   })
304
305 ▾ });

```

Rys. 12

Wynik działania tej funkcji można zobaczyć na Rys. 13.



Rys. 14

## 2.7 Częstotliwość występowania danego słowa

```

124 #FUNKCJA PRZYGOTOWUJE DANE DO WYPISANIA ILOSCI DANEGO SLOWA I STWORZENIA WYKRESU 10 NAJPOPULARNIEJSZYCH
125 ~ ilosc_slow<-function(tabela2){
126
127     stop<-as.vector(unlist(read.csv(file="stop_words_pl.txt", header=FALSE, sep="," , fileEncoding="UTF-8")));
128
129     dane<-solr_search(conn = polaczenie, params = list(q=paste(tabela2,":*"),fl=paste(tabela2), rows=-1));
130     sprawdz_duplikaty_ilosc_slow <- duplicated(dane)
131     view(sprawdz_duplikaty)
132
133     dokumenty<-Corpus(vectorSource(str_enc_toutf8(dane)));
134     dokumenty<-tm_map(dokumenty, removePunctuation, preserve_intra_word_dashes=TRUE);
135     dokumenty<-tm_map(dokumenty, removeNumbers);
136     dokumenty<-tm_map(dokumenty, removeWords, stop);
137     usun_znaki<-function(x) gsub("[\u00c4\u00e4\u00e4\u00e4]", "", x);
138     dokumenty<-tm_map(dokumenty, usun_znaki);
139
140     for(d in 1:length(dokumenty))
141     {
142         dokumenty[[d]]$content<-lematyzacja(dokumenty[[d]]$content);
143         dokumenty[[d]]$content<-str_enc_toutf8(dokumenty[[d]]$content);
144     }
145
146     tdm1<-TermDocumentMatrix(dokumenty);
147     m1<-as.matrix(tdm1);
148     v<-sort(rowSums(m1),decreasing=TRUE);
149     return(v);
150 ~ }

```

Rys. 15

Funkcja która przygotowuje dane do wyświetlenia danego słowa została przedstawiona na Rys. 15. Wykorzystuje funkcje lematyzującą tekst i zwraca posortowana i zliczona ilość słów. Następnie trafia do logiki serwera (Rys. 16) gdzie dane są przetwarzane. W liniach od 317 do 322 jest zawarty wykres, który wyświetla 10 najbardziej popularnych słów (wg. Występowania). W linii 328 Rys. 16 znajduje się tabela, która jest wyświetlana w UI.

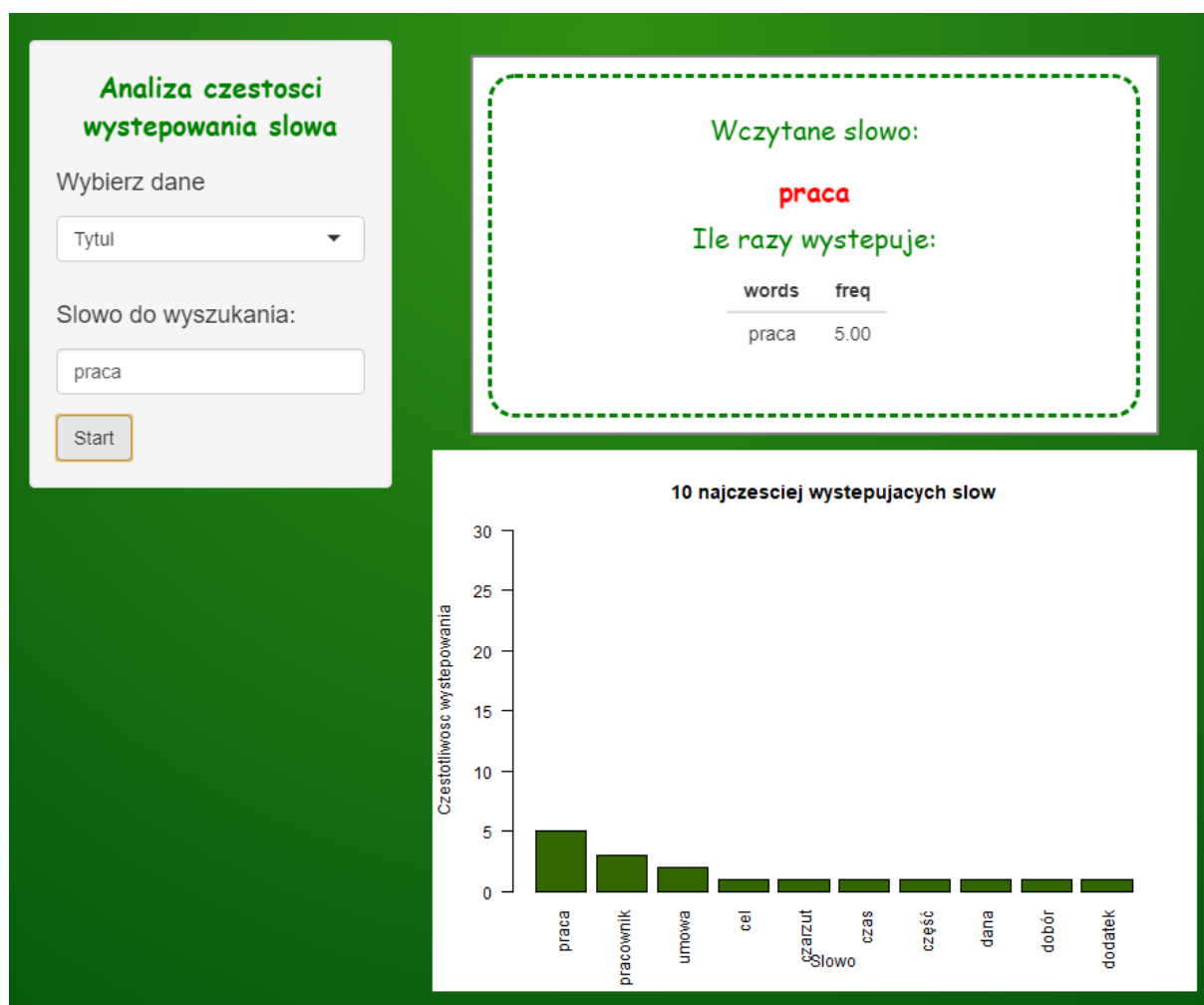
```

307 observeEvent(input$buttonA3,{
308   #WCZYTANY NUMER
309   wypisane_slowo2 <- eventReactive(input$buttonA3, {
310     input$s1owo2
311   })
312   output$s1owooutput <- renderText({
313     wypisane_slowo2()
314   })
315   #WCZYTANA ILOSC SLOW
316   dane <- data.frame(words=names(ilosc_slow(input$selectTab3)), freq=ilosc_slow(input$selectTab3))
317   output$wykresA3 <- renderPlot({
318     barplot(dane[1:10,]$freq, las = 2, names.arg = dane[1:10,]$word,
319             col = "#336600", main = "10 najczesciej wystepujacych slow",
320             ylab = "Czestotliwosc wystepowania", ylim = c(0,30),
321             xlab = "Slovo")
322   })
323
324   inputWORD <- input$s1owo2
325   outputWORD <- filter(dane,
326                        words == toString(inputWORD)
327                        )
328   output$table <- renderTable({outputWORD})
329 });

```

Rys. 16

Wygląd panelu z analizą częstotliwości występowania słowa został przedstawiony na Rys. 17.



Rys. 17

## 2.8 Klasteryzacja danych

```

109 klasteryzacja <- function(tabela3, metryka, numberOfK){
110
111   dane<-solr_search(conn = polaczenie, params = list(q=paste(tabela3,":*"),fl=paste(tabela3), rows=-1));
112   sprawdz_duplikaty_klasteryzacja <- duplicated(dane)
113   view(sprawdz_duplikaty)
114
115   przetworzone_dane <- przygotowanie_danych(dane)
116   dane_z_solra <- przygotowanie_danych(dane)
117   dane_koncowe <- unique(dane_z_solra)
118
119   klasteryzacja.h<-agnes(dane_koncowe, metric=paste(metryka), method="average");
120   klaster<-cutree(klasteryzacja.h, k=paste(numberOfK));
121   tabela.wynikowa.h<-cbind(dane_koncowe, klaster);
122
123   plot(klasteryzacja.h, which.plots=2);
124
125   return(tabela.wynikowa.h);
126 }

```

Rys. 18

W projekcie wykorzystałem klasteryzację hierarchiczną. Funkcja klasteryzacja przyjmuje trzy parametry *tabela3*, *metryka*, *numberOfK*. Te parametry są wczytywane w UI przez użytkownika. Funkcja zwraca tabele z wynikami, która jest przetwarzana (Rys. 19) i wyświetlana (Rys. 20).

```

331 observeEvent(input$buttonA4,{
332   output$table2<- renderTable(klasteryzacja(input$selectTab4, input$selectTab5, input$numberOfK))
333 })

```

Rys. 19

### Klasteryzacja

Wybierz dane

Tytuł ▼

Metryka

euclidean ▼

Ilość klastrow

3

Start

words	freq	klaster
praca	5.00	1
pracownik	3.00	1
umowa	2.00	2
cel	1.00	3
czarzut	1.00	3
czas	1.00	3
część	1.00	3
dana	1.00	3
dobór	1.00	3
dodatek	1.00	3
dokumentacja	1.00	3

Rys. 20