

Zadanie zaliczeniowe 1

Hanna Kranas

6 maja 2017

Contents

Przygotowanie danych	1
Usuwanie trudnych znaków	1
Funkcja przygotowująca dane	2
Kod i symulacje	2
Symbole	2
HMM i BaumWelch - model i trenowanie	3
Pełna symulacja literowa	3
Klastrowanie sufiksów	3
Pełna symulacja sufiksowa	4
Wyniki	4
Litery z polskimi znakami	4
Trzy stany	5
Litery bez polskich znaków	6
Słowa - sufiksy, klastrowane	8
Pięć stanów - więcej sufiksów:	8
Pięć stanów - mniej sufiksów:	9
Symulowanie na ulubionej poezji wytrenowanymi modelami	11
Dżabbersmok	12
Dziewczyna	16
Do prostego człowieka	20

Przygotowanie danych

Usuwanie trudnych znaków

Ponieważ ani R ani Python nie radziły sobie z usuwaniem niektórych znaków, usunięte zostały w notepadzie.

```

Console D:/studiallstopien/SADII/
> daj_symbole(obserwacja <- unlist(strsplit(wczytaj_i_przygotuj(), '')))
[1] "a" "e" "d" "m" " " "i" "c" "k" "w" "z" "p" "n" "t" "u" "s" "y" "l" "o"
[19] "j" "ę" "g" "r" "ó" "ł" "ż" "ś" "ą" "ź" "b" "ć" "h" "f" "ń" "x" "v" "à"
[37] "q" "1" "8" "2" "æ"
> iconv(daj_symbole(obserwacja <- unlist(strsplit(wczytaj_i_przygotuj(), ''))), to='ASCII//TRANSLIT')
[1] "a" "e" "d" "m" " " "i" "c" "k" "w" "z" "p" "n"
[13] "t" "u" "s" "y" "l" "o" "j" "AT" "g" "r" "Al" "L,"
[25] "LL" "L>" "A." "Ls" "b" "A?" "h" "f" "L\" "x" "v" "A"
[37] "q" "1" "8" "2" "A|"
> "æ"=="a"
[1] TRUE

```

```

>>> 'æ' == 'a'
True
>>> ord('æ')
97
>>> ord('a')
97
>>>

```

Owoż te wszystkie rzeczy mając na uwadze,

Ja, reprezentujący województwa wł
 Moją konfederacką ogłaszam wam la
 że Jacek wierną służbą i cesarską
 Zniósł infamiji plamę, powraca do
 I znowu się w rząd prawych patryj
 Więc kto będzie śmiał Jacka zmarł
 Wspomnieć kiedy o dawnej zagładzo
 Ten podpadnie za karę takiego wyr
 Gravis notæ macula, wedle słów St
 Karzących tak militem jak i skarta
 Co by siał infamiją na obywatela;

Szukaj	Zamień	Szukaj w plikach	Oznacz
Szukany tekst: æ			
Zamień na: a			
<input type="checkbox"/> Znajdź tylko całe wyrazy <input type="checkbox"/> Uwzględniaj wielkość liter <input checked="" type="checkbox"/> Wróć na początek pliku			

Funkcja przygotowująca dane

Przygotowanie danych wykonane jednak zostało głównie w R:

```

wczytaj_i_przygotuj <- function(nazwa_pliku='pan-tadeusz.txt', polskie = TRUE){
  docs <- SimpleCorpus(DirSource("D:/studiallstopien/SADII/SAD-zal-1",encoding = 'UTF-8'),control = list())
  docs <- tm_map(docs,removePunctuation)
  docs <- tm_map(docs,removeNumbers)
  docs <- tm_map(docs,stripWhitespace)
  plik <- tolower(as.character(docs[[nazwa_pliku]]))
  plik <- gsub ("[-....-\\t-]", "", plik)
  plik <- plik[plik!=""] #usuwaam puste
  if(polskie){
    cat('POLSKIE LITERY! æę!')
  } else {
    plik <- iconv(plik,from="UTF-8",to="ASCII//TRANSLIT") #bez polskich
  }
  return(plik)
}

```

Kod i symulacje

Symbole

Po przygotowaniu pliku, potrzebna nam będzie możliwość wyciągnięcia symboli z obserwacji:

```

daj_symbole <- function(obserwacja){
  cat('dlugosc obserwacji: ')

```

```

cat(length(obserwacja))
#ta funkcja ekstrahuje wszystkie unikalne elementy - czyli dla nas symbole
symbole <- sort(unique(obserwacja))
return(symbole)
}

```

HMM i BaumWelch - model i trenowanie

Tworzenie HMM i trenowanie algorytmem Bauma-Welcha odbywa się w ten sposób:

```

prob <- function (x) {x / sum (x)} # funkcja pomocnicza

symuluj <- function(obserwacja,stany){
  symbole <- daj_symbole(obserwacja)
  macierz_emisji <- c() #szukujemy macierz emisji
  for(i in 1:length(stany)){
    x <- rep(1/length(symbole),length(symbole))
    x <- x + rnorm(length(symbole),0,1/1000000) #dajemy małe zaburzenie
    x <- abs(x)/sum(abs(x)) #i poprawiamy zeby sie sumowały do jedynki dalej
    macierz_emisji <- c(macierz_emisji,x)
  }
  #tworzymy hmm'a z nieco-losowymi prawdopodobieństwami początkowymi i przejść
  hmm <- initHMM (stany, symbole, startProbs=(prob(runif(length(stany)))),
    transProbs=apply(matrix(runif(length(stany)*length(stany)),
      length(stany)),1,prob),
    emissionProbs=matrix(macierz_emisji,nrow =length(stany),
      ncol=length(symbole),byrow=TRUE))

  #uczymy algorytmem Bauma-Welcha
  wynik <- baumWelch(hmm,obserwacja)
  return(wynik)
}

```

Pełna symulacja literowa

Pełna symulacja dla liter z przygotowanego pliku wygląda tak:

```

lityry <- function(przygotowany_plik=wczytaj_i_przygotuj()){
  obserwacja <- unlist(strsplit(przygotowany_plik, ''))
  obserwacja <- obserwacja[!is.na(obserwacja)]
  stany <- c('1','2') # dwa stany
  wynik <- symuluj(obserwacja,stany)
  return(wynik)
}

```

Klastrowanie sufiksów

W przypadku słów konieczne było klastrowanie, ponieważ inaczej algorytm Bauma-Welcha nie dawał sobie rady. Sufiksy max-3literowe poklastrowane zostały po ich długości oraz umiejscowieniu samogłosek/spółgłosek, łącznie 14 klastrów. Wykonano też drugie klastrowanie, w którym sufiksy jednoliterowe połączono w jedną grupę ('1'), i dwuliterowe też ('2').

```

klastruj <- function(obserwacja){
  obserwacja <- gsub ("^[aeyiou][^aeyiou][aeyiou]$", "ara", obserwacja)
  obserwacja <- gsub ("^[aeyiou][aeyiou][aeyiou]$", "aaa", obserwacja)
  obserwacja <- gsub ("^[aeyiou][^aeyiou][^aeyiou]$", "arr", obserwacja)
  obserwacja <- gsub ("^[aeyiou][aeyiou][^aeyiou]$", "aar", obserwacja)
  obserwacja <- gsub ("^[^aeyiou][^aeyiou][aeyiou]$", "rra", obserwacja)
  obserwacja <- gsub ("^[^aeyiou][aeyiou][aeyiou]$", "raa", obserwacja)
  obserwacja <- gsub ("^[^aeyiou][aeyiou][^aeyiou]$", "rar", obserwacja)
  obserwacja <- gsub ("^[^aeyiou][^aeyiou][^aeyiou]$", "rrr", obserwacja)
  obserwacja <- gsub ("^[^aeyiou]$", "r", obserwacja) #1
  obserwacja <- gsub ("^[aeyiou]$", "a", obserwacja) #1
  obserwacja <- gsub ("^[^aeyiou][^aeyiou]$", "rr", obserwacja) #2
  obserwacja <- gsub ("^[aeyiou][aeyiou]$", "aa", obserwacja) #2
  obserwacja <- gsub ("^[aeyiou][^aeyiou]$", "ar", obserwacja) #2
  obserwacja <- gsub ("^[^aeyiou][aeyiou]$", "ra", obserwacja) #2
  return(obserwacja)
}

```

Pełna symulacja sufiksowa

Sposób symulacji dla słów jest następujący:

```

slova <- function(przygotowany_plik=wczytaj_i_przygotuj()){
  #słowa wymagają trochę więcej przygotowań
  plik <- unlist(strsplit(przygotowany_plik," "))
  obserwacja <- unlist(str_sub(plik[plik!=""],start=-3)) #wyciągamy max-3literowe sufiksy
  obserwacja <- obserwacja[!is.na(obserwacja)]
  obserwacja <- klastruj(obserwacja)
  print(table(obserwacja)) #wypisuję dla informacji o tym ile jakich sufiksów
  stany <- c('1','2','3','4','5') # 5 stanów
  wynik <- symuluj(obserwacja,stany)
  return(wynik)
}

```

Wyniki

W poniższych wynikach można zobaczyć analizy dla liter (znaki polskie lub bez) z dwoma stanami, liter ze znakami polskimi z trzema stanami (tylko 20000 znaków) oraz sufiksów z pięcioma.

Litery z polskimi znakami

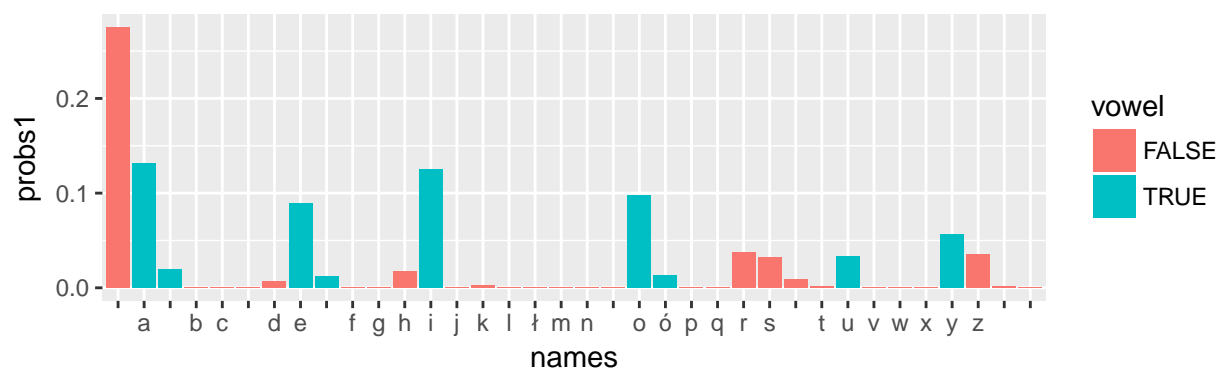
Widać wyraźnie że pierwsza grupa jest samogłoskowa. Na wykresach przedstawione są prawdopodobieństwa emisji poszczególnych symboli z wytrenowanego modelu HMM. Przy polskich znakach ładnie dzielą się stany, ale widać że są problemy przy literach które rzadko występują w tekście.

```

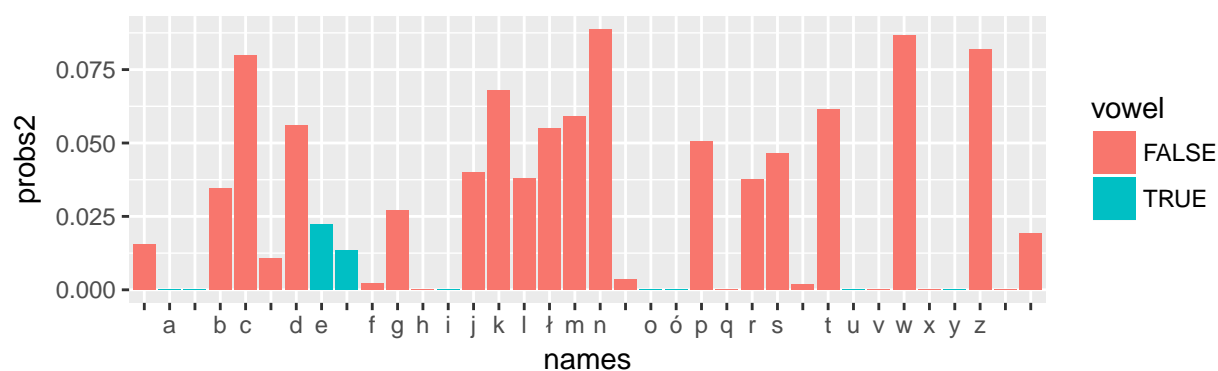
## [1] "To jest grupa pierwsza:"
##  a ą e h i o ó ś u y ż
## [1] "To jest grupa druga:"
## b c ć d ę f g j k l ł m n ń p q r s t v w x z ź

```

Z polskimi znakami: Pierwszy stan – samogłoski

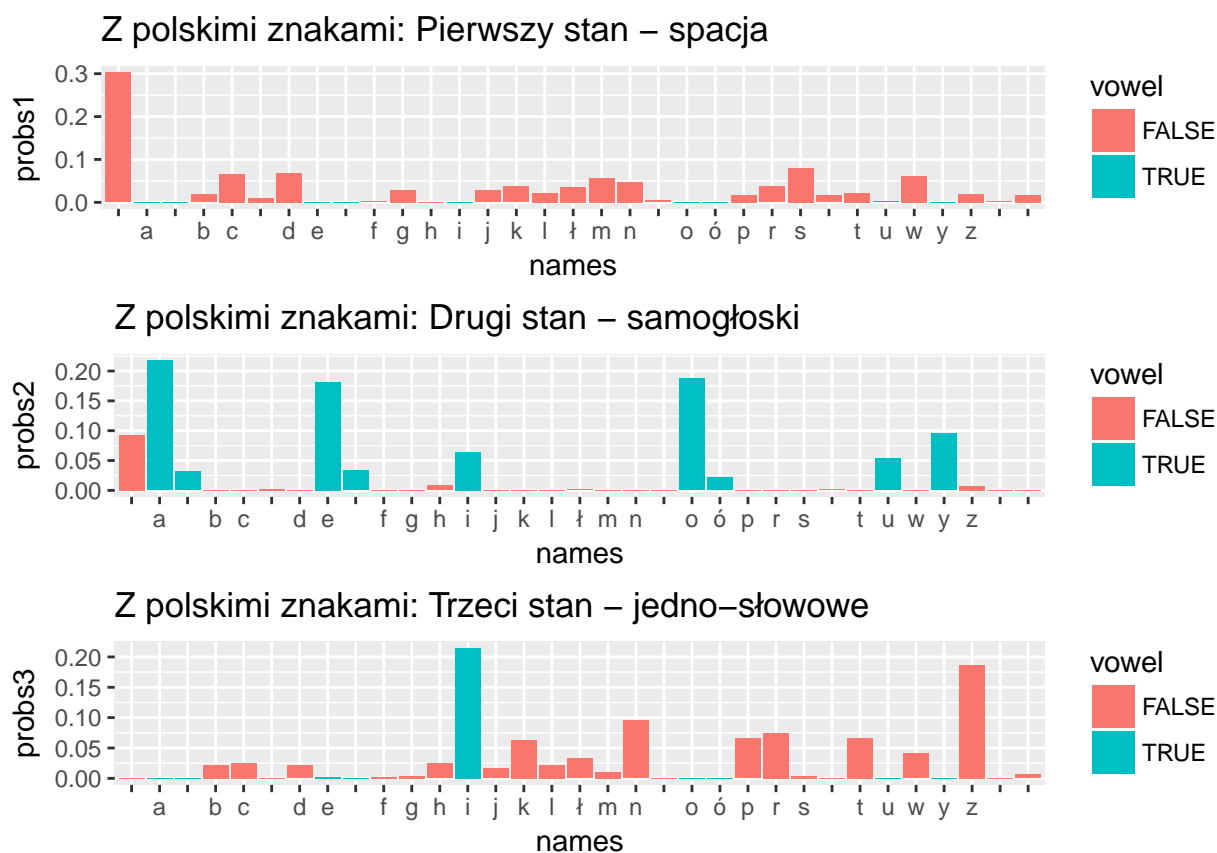


Z polskimi znakami: Drugi stan – spółgłoski



Trzy stany

Z ciekawości wykonano analizę dla większej (3) liczby stanów. Dla liter z polskimi znakami i 3 stanami widać stan “spacjowy”, stan samogłoskowy, oraz stan trzeci, który na pierwszy rzut oka nie jest jasny, ale może jest to związane z często występującymi literami “jedno-słowowymi” (i,z)? Ponieważ nie była jasna interpretacja (oraz ograniczone możliwości obliczeniowe) nie stosowano analiz na większych liczbach stanów.

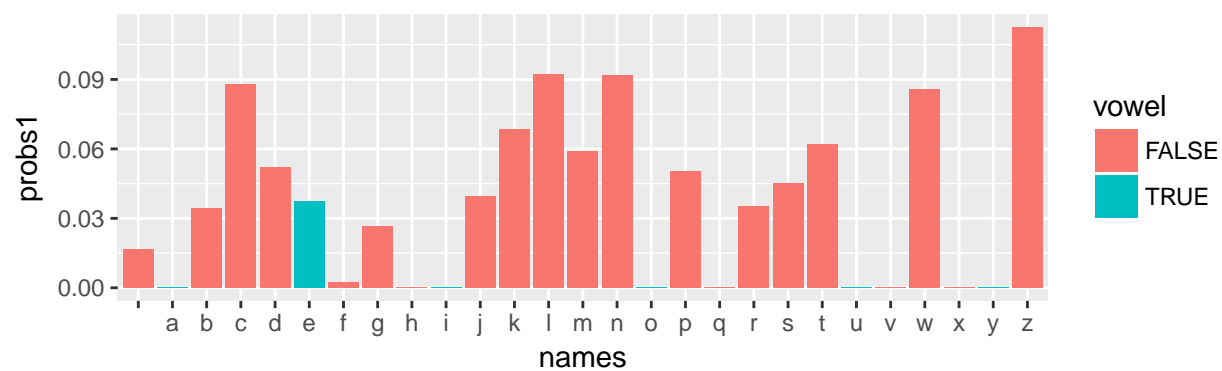


Litery bez polskich znaków

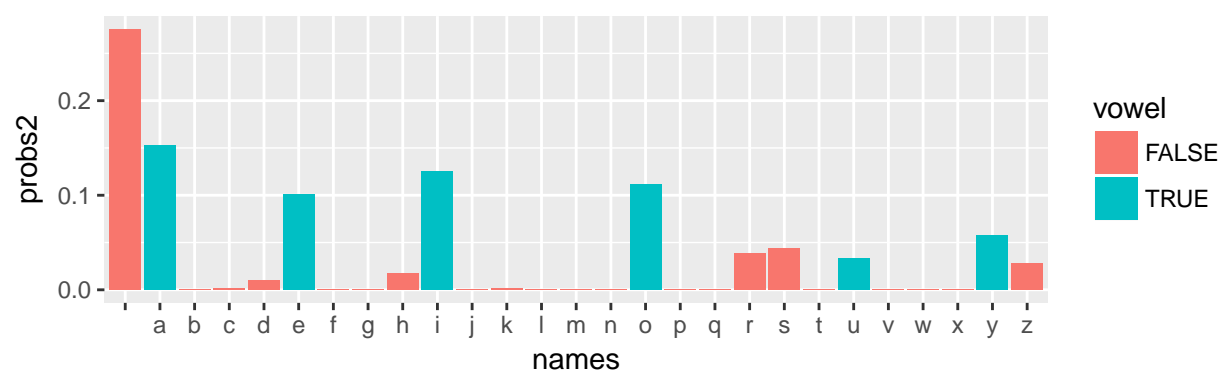
Widać wyraźnie że druga grupa jest samogłoskowa. Na wykresach przedstawione są prawdopodobieństwa emisji poszczególnych symboli z wytrenowanego modelu HMM.

```
## [1] "To jest grupa pierwsza:"
## b c d f g j k l m n p q s t v w x z
## [1] "To jest grupa druga:"
## a e h i o r u y
```

Bez polskich znaków: Pierwszy stan – samogłoski



Bez polskich znaków: Drugi stan – spółgłoski

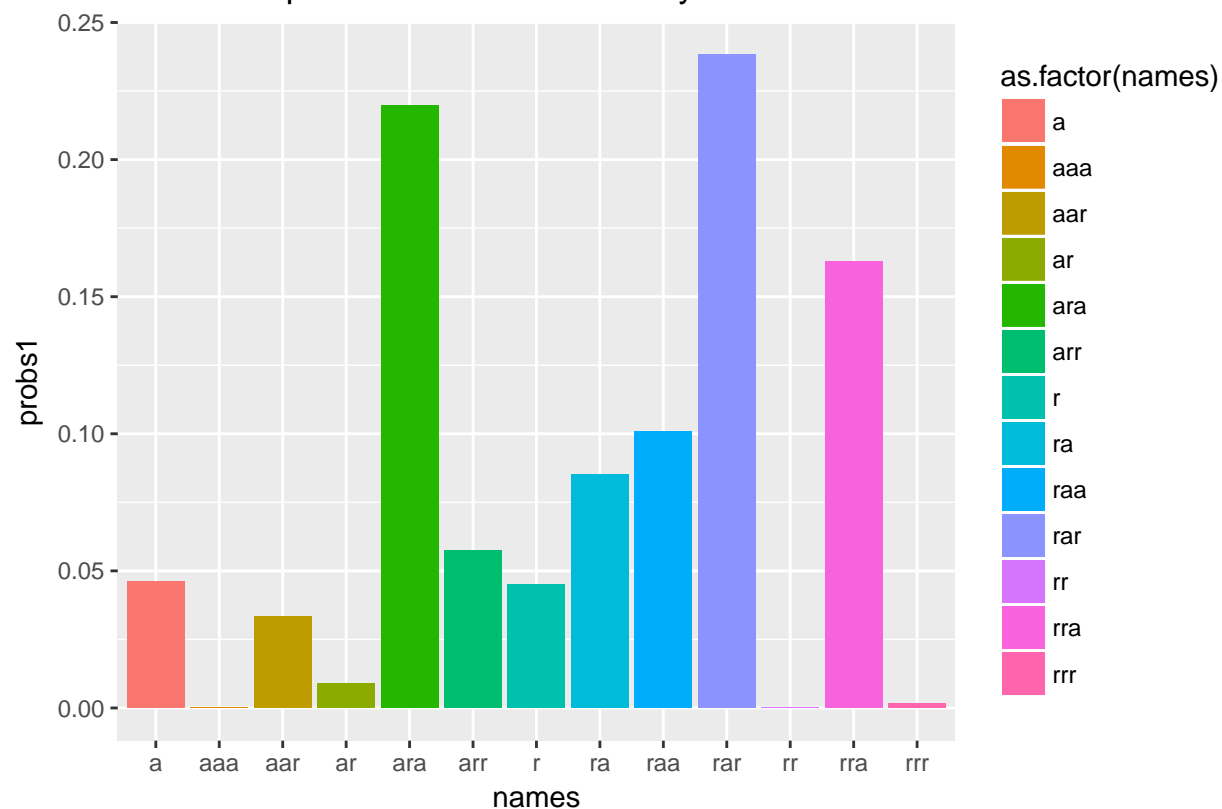


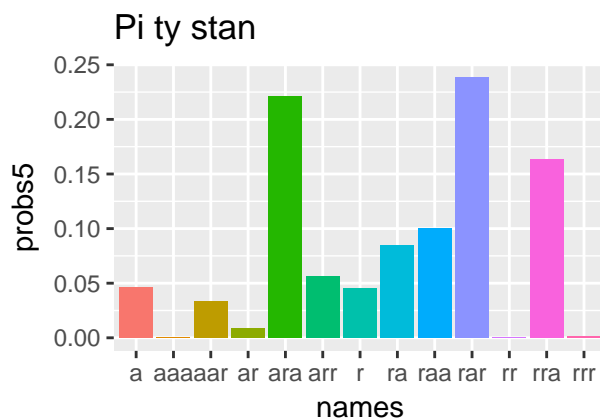
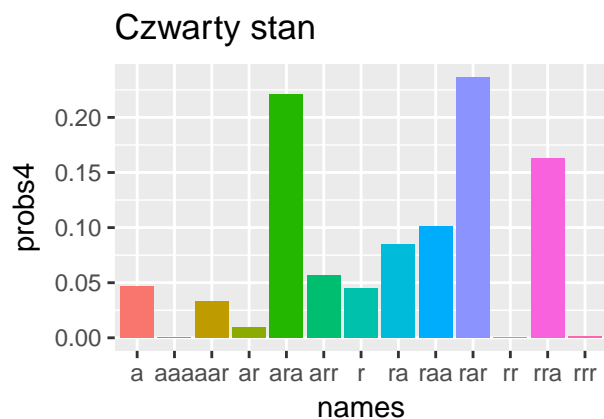
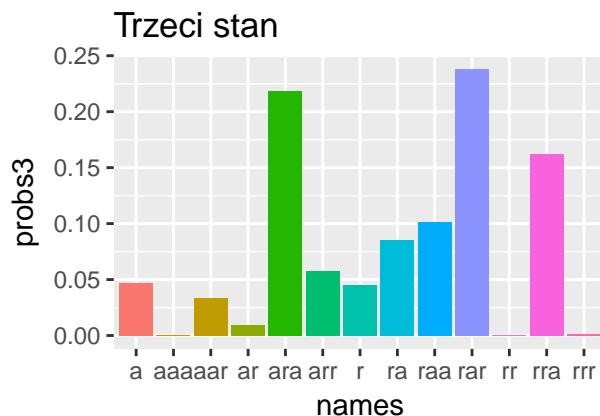
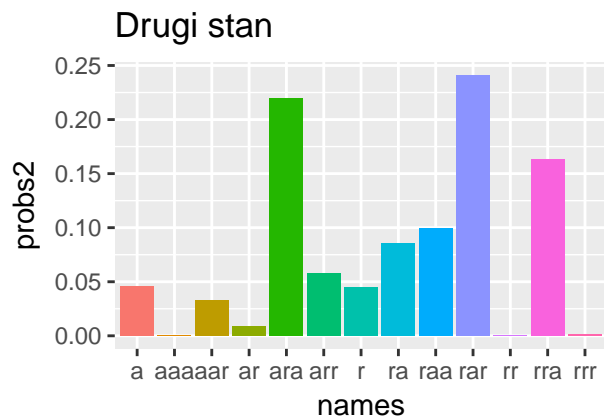
Słowa - sufiksy, klastrowane

Pięć stanów - więcej sufiksów:

W każdym ze stanów prawdopodobieństwa emisji danych symboli są prawie identyczne.

Słowa bez polskich znaków: Pierwszy stan

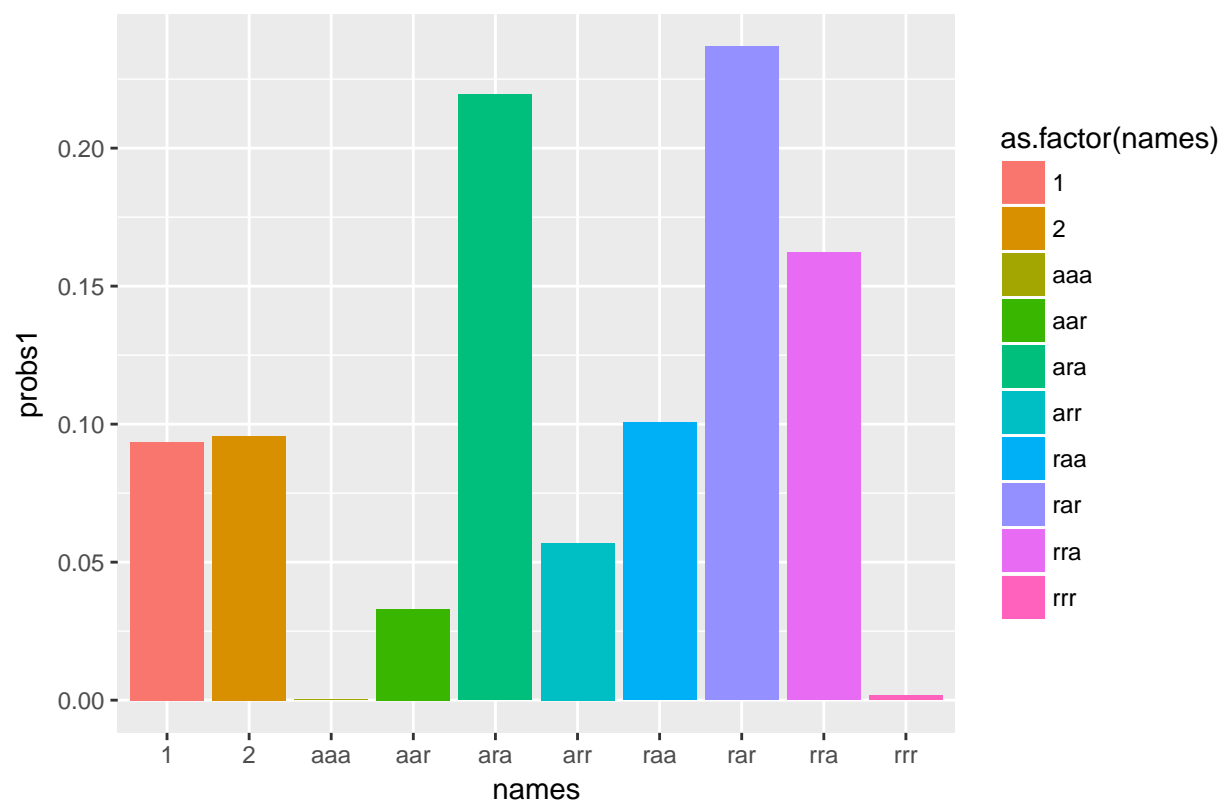


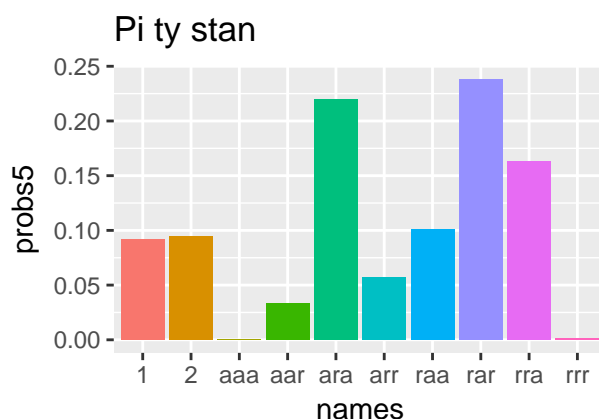
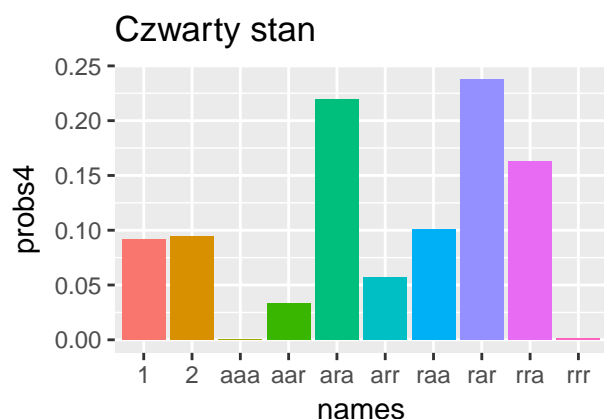
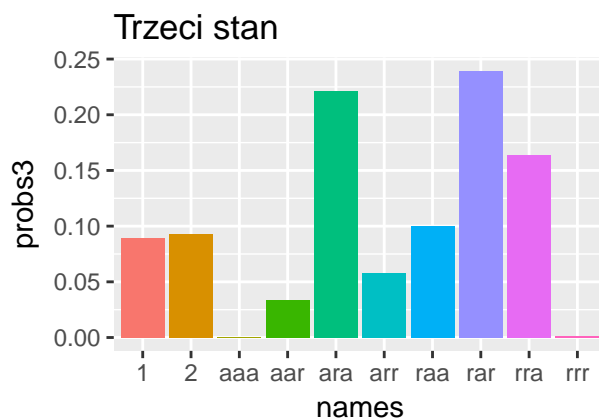
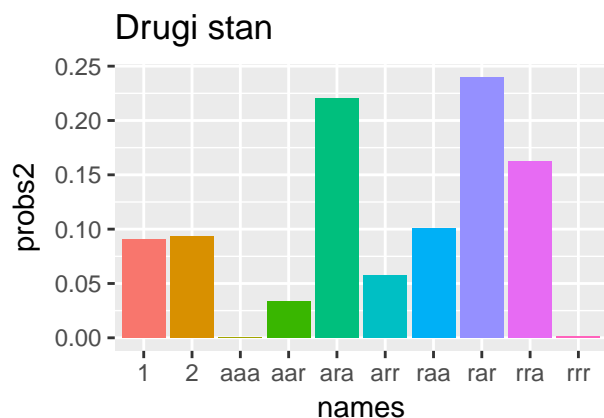


Pięć stanów - mniej sufiksów:

Wykonano dodatkowe klastrowanie, z mniejszą liczbą sufiksów, żeby sprawdzić czy coś to zmieni. Niestety znów w każdym ze stanów prawdopodobieństwa emisji danych symboli są prawie identyczne, można więc uznać że badanie sufiksów nie dało efektu.

Słowa bez polskich znaków: Pierwszy stan





Symulowanie na ulubionej poezji wytrenowanymi modelami

Wybrano 3 różne wiersze:

1. Jabberwocky, Lewis Carroll - przekład Macieja Słomczyńskiego (nie do końca polski język, ale uznano to za tym bardziej ciekawy test)
2. Dziewczyna, Bolesław Leśmian
3. Do prostego człowieka, Julian Tuwim

Dla każdego wiersza i rodzaju analizy literowej (polskie/bez polskich) wykonano dwa wykresy - średniego prawdopodobieństwa wystąpienia danej grupy do samogłoskowego stanu oraz prawdopodobieństw spółgłosek i samogłosek w stanie samogłoskowym. Jak można zauważyć prawdopodobieństwo należenia samogłosek do stanu samogłoskowego we wszystkich trzech wierszach jest bardzo wysoka, co sugeruje że zastosowany HMM został dobrze wytrenowany.

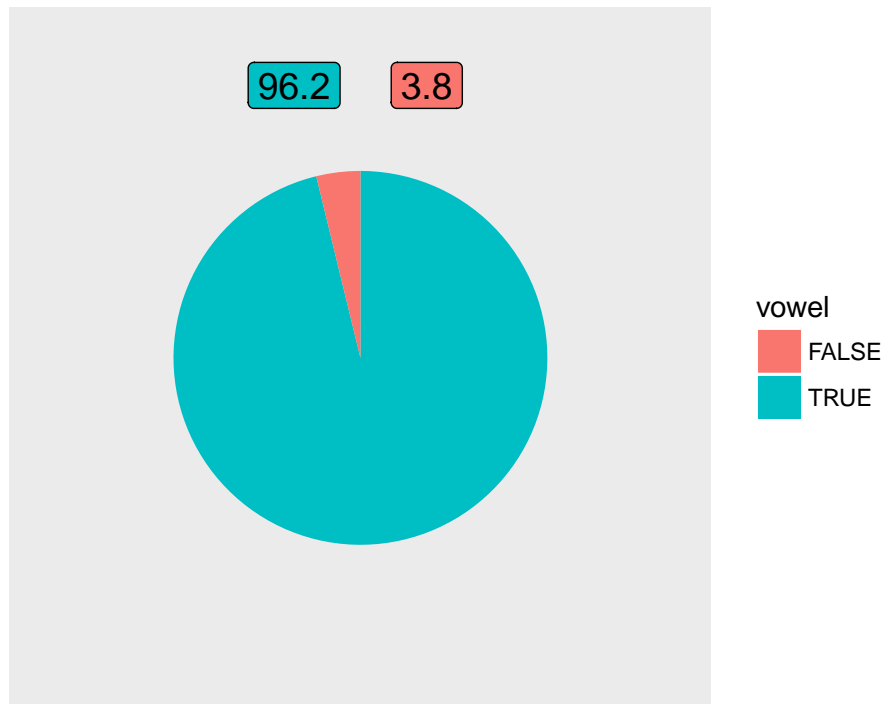
Z powodu słabych wyników dla dzielenia sufiksów do stanów nie wykonano tych analiz na wierszach.

Wiersze załączono do paczki z raportem.

Dżabbersmok

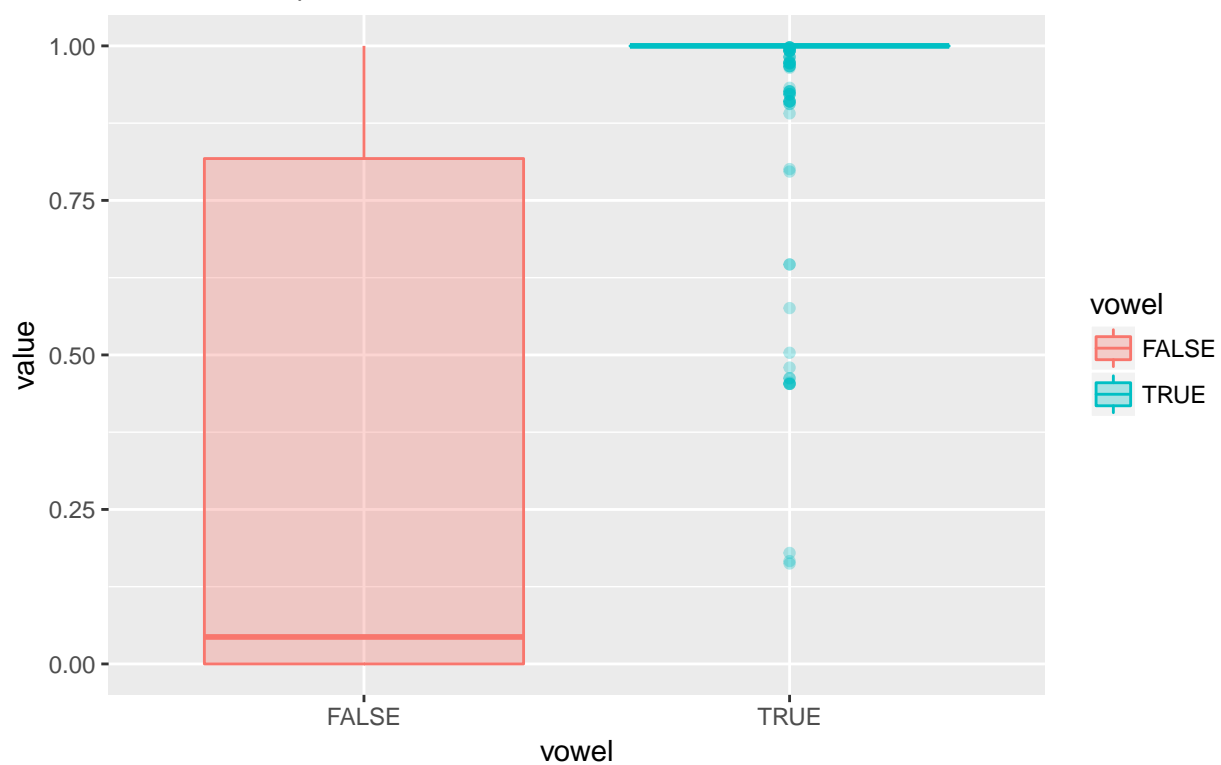
rednie prawdopodobieństwo przynależenia samogłoski
do stanu samogłoskowego

Dżabbersmok – polskie



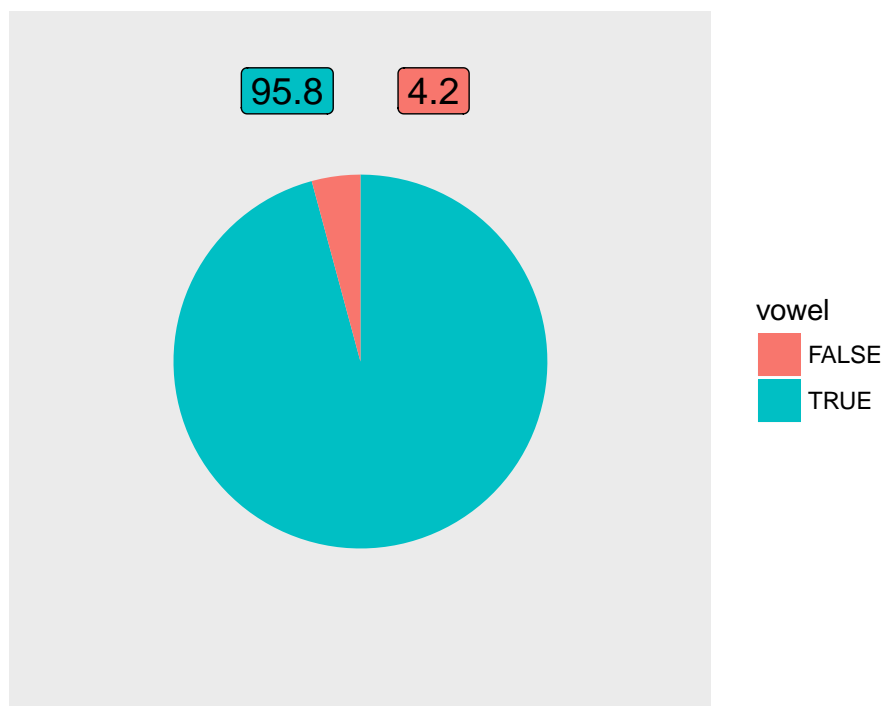
Prawdopodobie stwa spółgłosek i samogłosek w stanie samogłoskowym

D abbersmok – polskie



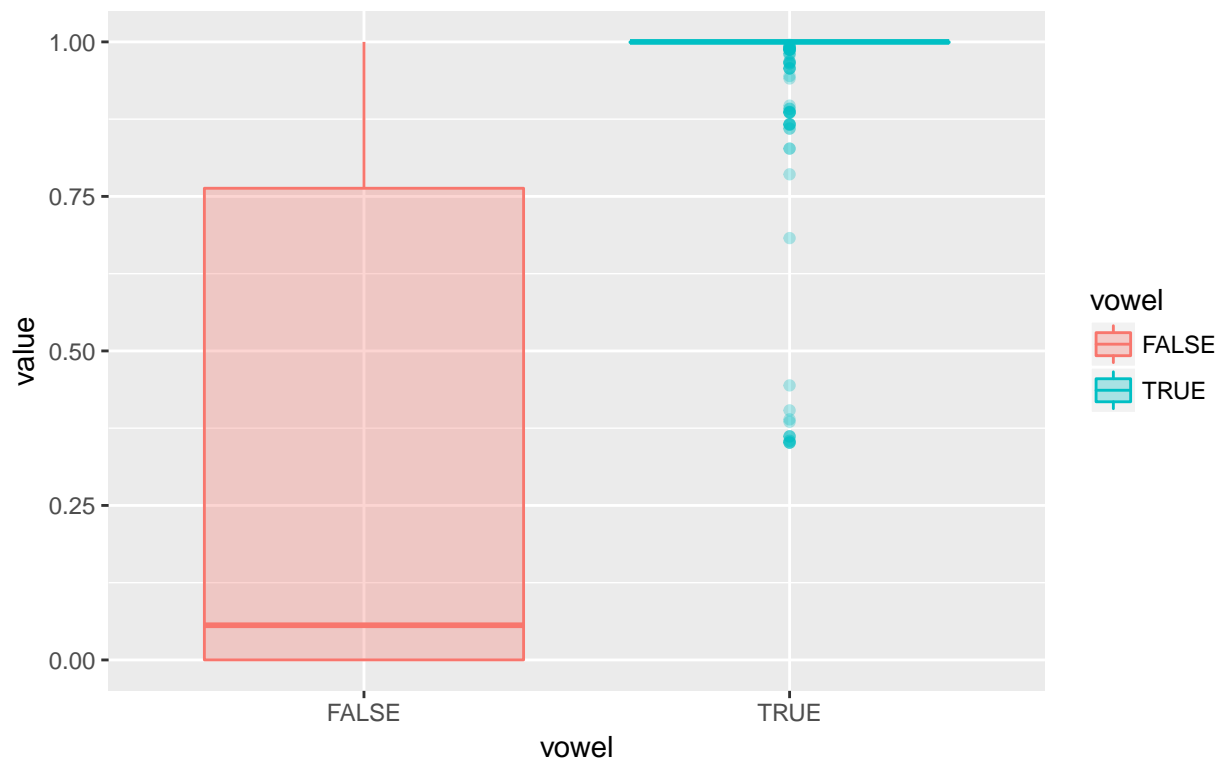
rednie prawdopodobieństwo przynależenia samogłoski
do stanu samogłoskowego

D abbersmok – bez polskich



Prawdopodobie stwa spółgłosek i samogłosek w stanie samogłoskowym

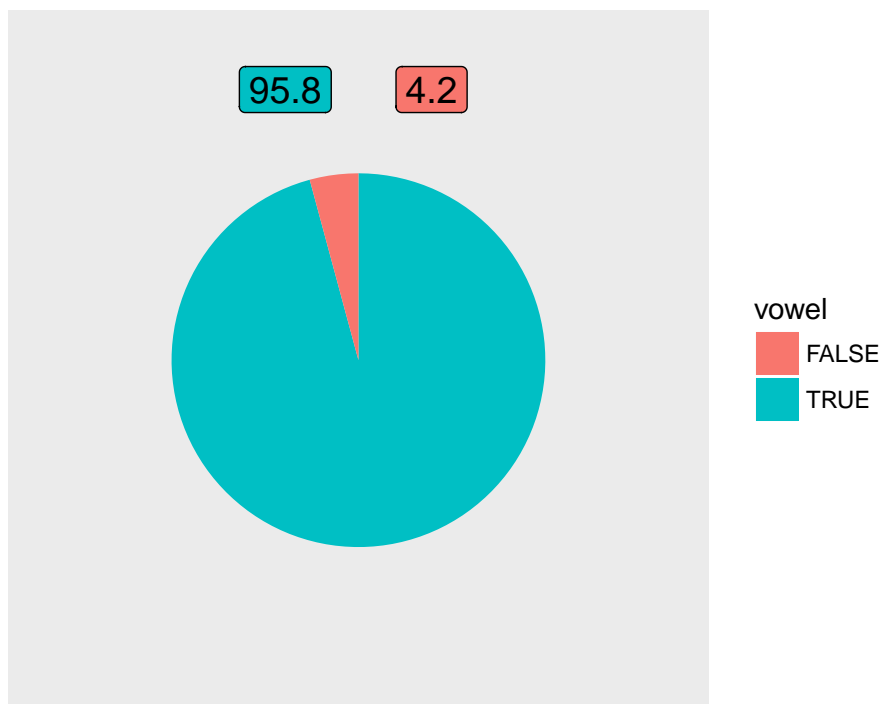
D abbersmok – bez polskich



Dziewczyna

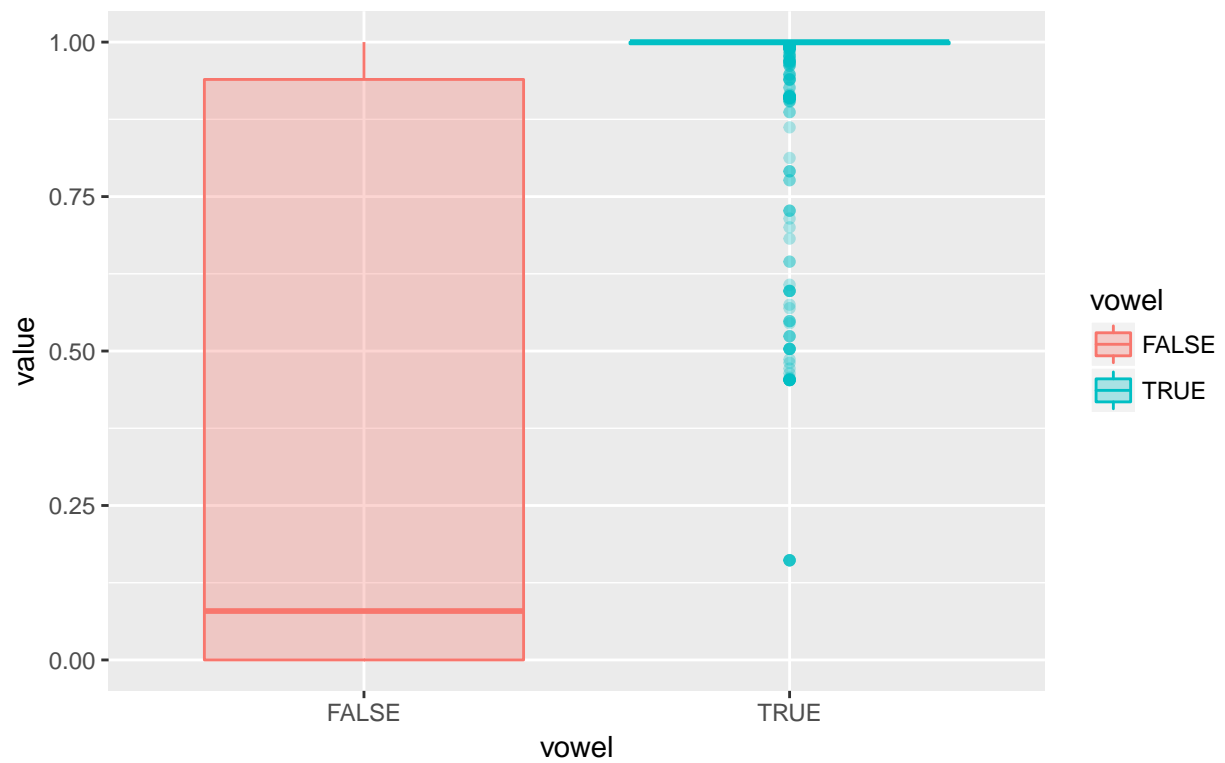
rednie prawdopodobieństwo przynależenia samogłoski
do stanu samogłoskowego

Dziewczyna – polskie



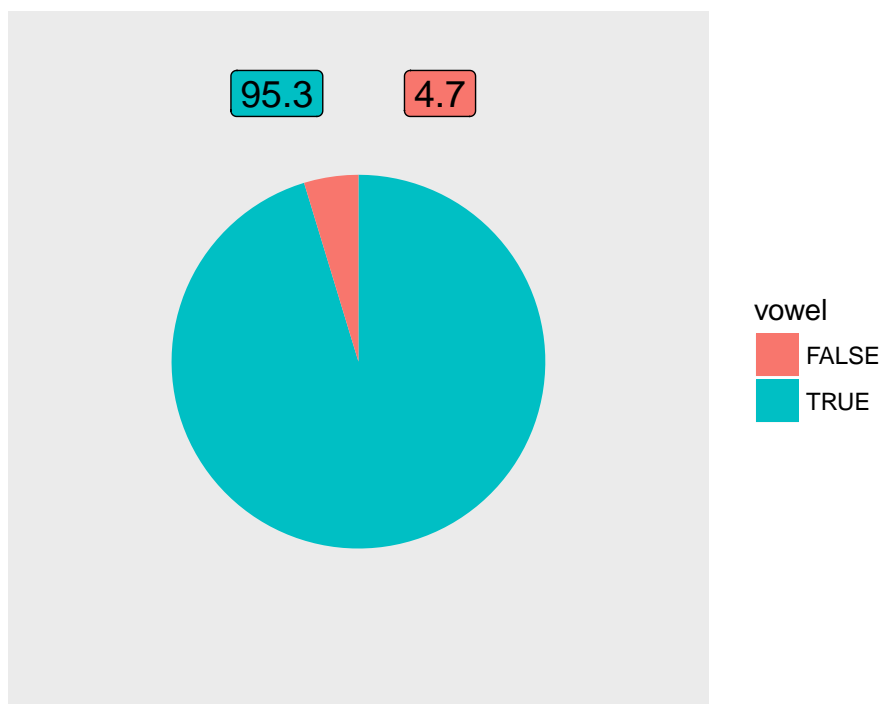
Prawdopodobieństwo spółgłosek i samogłosek w stanie samogłoskowym

Dziewczyna – polskie



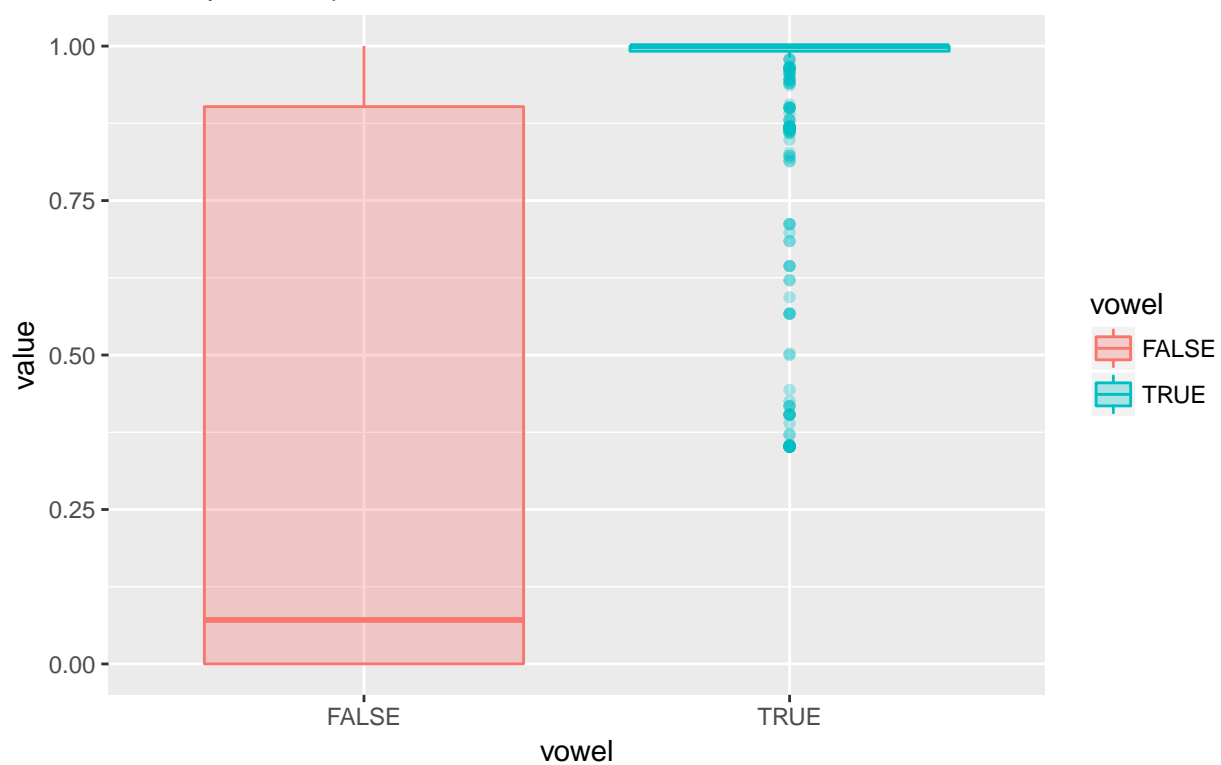
rednie prawdopodobieństwo przynależenia samogłoski
do stanu samogłoskowego

Dziewczyna – bez polskich



Prawdopodobie stwa spółgłosek i samogłosek w stanie samogłoskowym

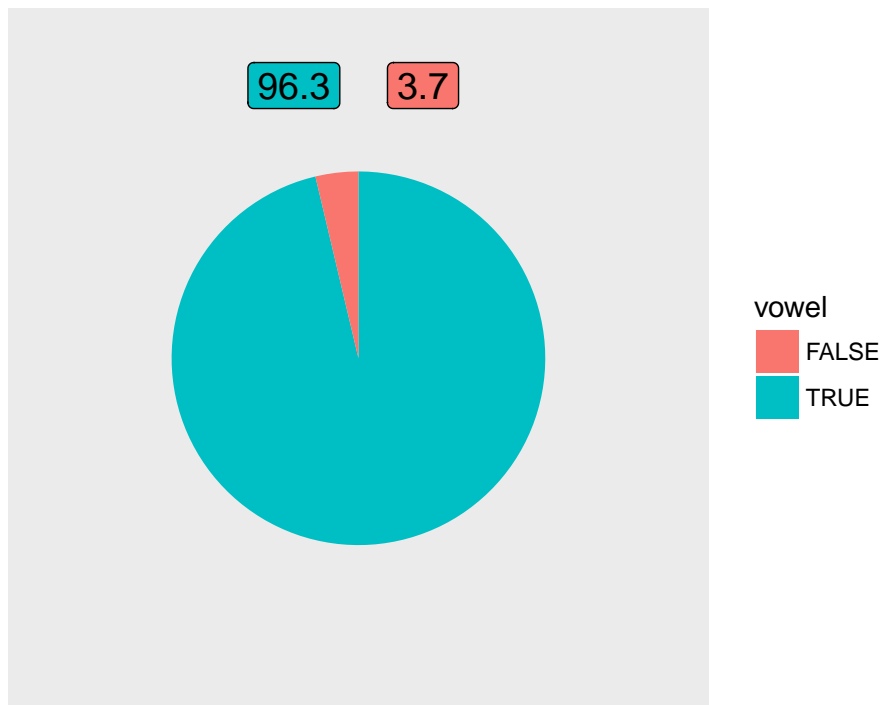
Dziewczyna – bez polskich



Do prostego człowieka

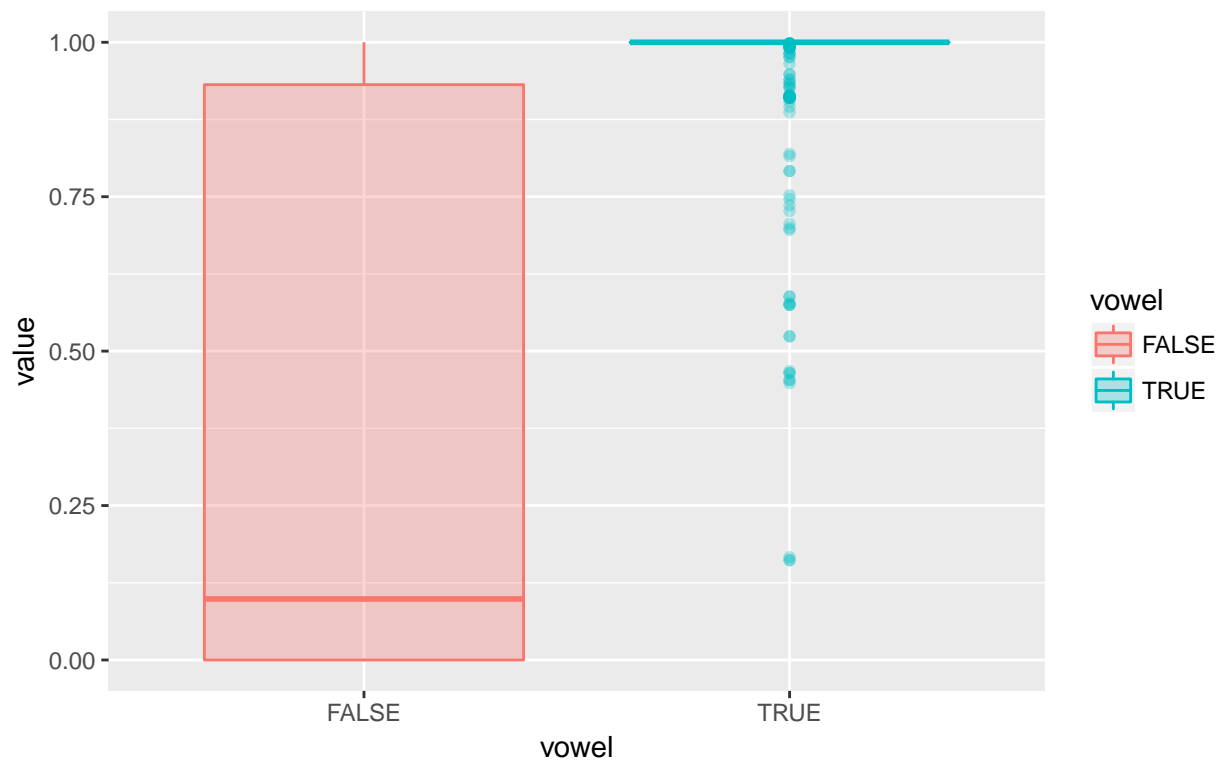
rednie prawdopodobieństwo przynależenia samogłoski
do stanu samogłoskowego

Do prostego człowieka – polskie



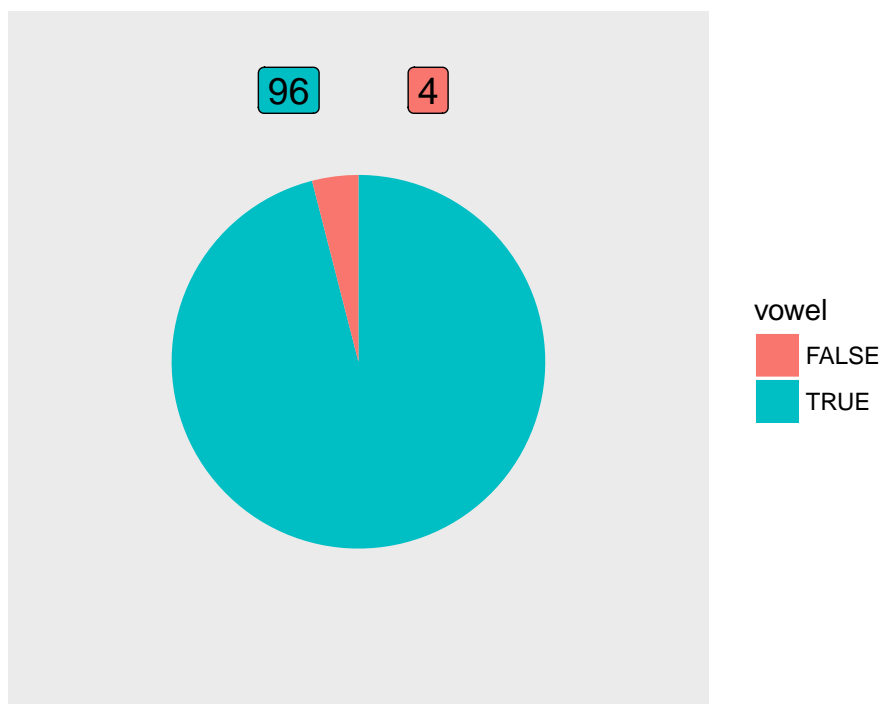
Prawdopodobie stwa spółgłosek i samogłosek w stanie samogłoskowym

Do prostego człowieka – polskie



rednie prawdopodobieństwo przynależenia samogłoski
do stanu samogłoskowego

Do prostego człowieka – bez polskich



Prawdopodobie stwa spółgłosek i samogłosek w stanie samogłoskowym

Do prostego człowieka – bez polskich

