

Метода анализа текста в R

Алексей Горгадзе
Анастасия Кузнецова

NET-RESEARCH.NET

Чистка данных

```
library(tm); library(tidytext)
```

```
text <- gsub("[^[:alnum:]]", " ", text)
```

```
text <- gsub("[a-zA-Z0-9]+", "", text)
```

```
text <- tolower(text)
```

```
text <- removeNumbers(text)
```

```
text <- removePunctuation(text)
```

```
text <- removeWords(text, stopwords("russian"))
```

```
text <- removeWords(text,  
stoplist)
```

В стоп слова входят
(stoplist):

- слишком частотные
- слишком редкие
- слишком короткие
- не существительные
- имена собственные

Лемматизация

Приведение словоформы к лемме (к инфинитиву)

MyStem (Яндекс) - производит морфологический анализ текста на русском языке

text.tmp <- system2("mystem", c("-c", "-l", "-d"), input=docs\$text, stdout=TRUE) (должен быть установлен MyStem)

кошками -> кошка

Стемминг (урезание слова до основы):

кошками -> кошк

Форматы текстовых данных

```
corpus1 <- Corpus(VectorSource(text), readerControl=list(language='ru'))
```

```
tdm.matrix <- TermDocumentMatrix(corpus1) / dtm.matrix <- DocumentTermMatrix(corpus1)
```

```
words_matrix <- as.matrix(tdm.matrix)
```

term-document matrix



	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
ад	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0
ангел	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
бояться	1	0	0	0	0	0	0	0	0	0	0	0	2	0	0
внутри	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
возникать	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
вопрос	3	0	0	0	0	0	0	0	1	0	0	0	3	0	0
впечатление	1	0	0	0	0	0	0	0	0	1	2	0	0	0	0
второй	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
вырастать	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
герой	1	0	0	0	0	3	0	2	1	2	1	0	0	0	0
голова	1	0	1	0	0	0	0	0	0	1	0	0	0	0	0
грязь	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Name	Type	Value
wCorpus	list [2392] (S3: SimpleCorpus,	List of length 2392
1	list [2] (S3: PlainTextDocumer	List of length 2
content	character [1]	'какой-то пустота образовываться внутри я после завершение...
meta	list [7] (S3: TextDocumentMet	List of length 7
author	character [0]	
timestamp	list [1] (S3: POSIXlt, POSIXt)	List of length 1
[[1]]	double [1]	31.43707
description	character [0]	
heading	character [0]	
id	character [1]	'1'
language	character [1]	'ru'
origin	character [0]	
2	list [2] (S3: PlainTextDocumer	List of length 2
3	list [2] (S3: PlainTextDocumer	List of length 2

Частотность слов

```
words_freq <- sort(rowSums(words_matrix), decreasing=TRUE)
```

```
words_freq <- data.frame(freq = words_freq, word = names(words_freq))
```

Облака слов

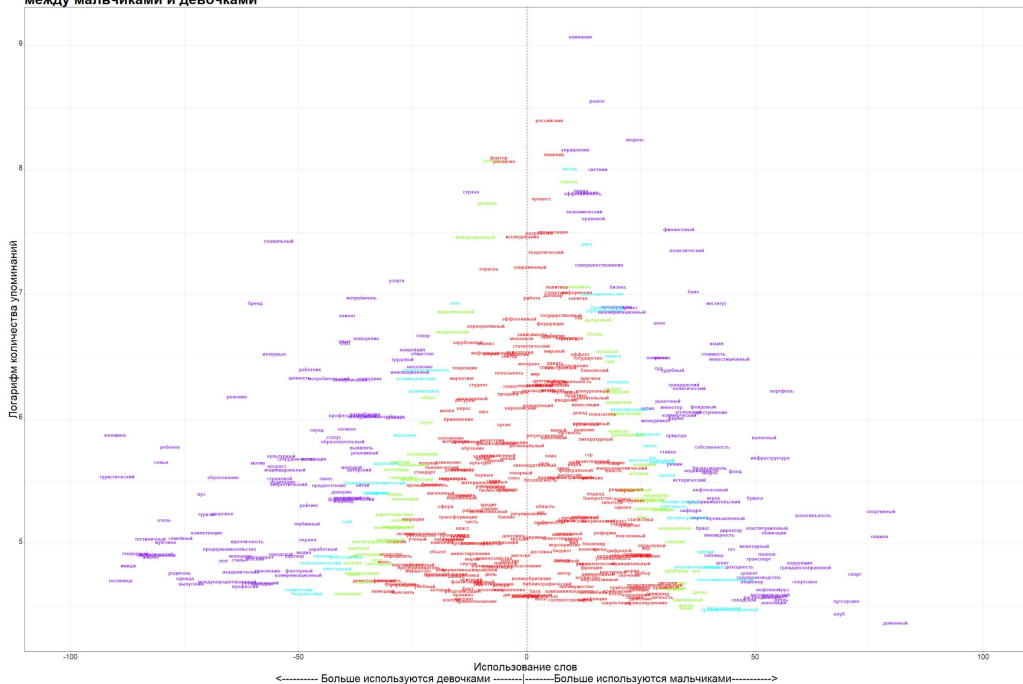


```
wordcloud(words = words_freq$word, freq =  
words_freq$freq, scale=c(2,.2), min.freq = 5,  
max.words=Inf, random.order=FALSE, rot.per=0.1,  
ordered.colors=FALSE,  
random.color=TRUE,colors=pal2)
```


Сравнение частотностей

```
wordsCompare$prop <- wordsCompare$Male/sum(wordsCompare$Male)
wordsCompare$prop2 <- wordsCompare$Female/sum(wordsCompare$Female)
```

Сравнение используемых слов на факультетах экономики, социологии, юриспруденции и менеджмента между мальчиками и девочками



```
# Broke down the z score formula a little to understand how it worked
a <- wordsCompare$prop
b <- wordsCompare$prop2
c <- wordsCompare$Male
d <- wordsCompare$Female
e <- sum(c)
f <- sum(d)
```

```
# z score formula - adds column for z scores
wordsCompare$z <- (a - b) / ((sqrt(((sum(c) * a) + (sum(d) * b)) / (sum(c) +
sum(d)) * (1 - ((sum(c) * a) + (sum(d) * b)) / (sum(c) +
sum(d)))) * (sqrt(sum(c) + sum(d)) / (sum(c) *
sum(d))))))
```

```
# calculate percentage reduction:
wordsCompare$dif1 <- -100 * (1 - wordsCompare$prop/wordsCompare$prop2)
# calculate percentage increase
wordsCompare$dif2 <- 100 * (1 - wordsCompare$prop2/wordsCompare$prop)
```

```
require(ggplot2)
```

```
png("eco_dif_m_f_words_full_size_byDif.png", width=3000,height=1500)
ggplot(wordsCompare3, aes(dif, log(abs(Male + Female))), size =
1,label=Row.names, colour = z2)+
  scale_colour_gradientn(name="Z Score", colours=c("#80FF00FF",
"#00FFFFFF", "#8000FFFF")) +
  geom_text(fontface = 2, alpha = .8) +
  #scale_size(range = c(3, 12)) +
  ylab("Логарифм от количества упоминаний") +
  xlab("Использование слов \n <-----Больше используются
девушками -----|-----Больше используются мальчиками----->")+
  geom_vline(xintercept=0, colour = "red", linetype=2)+
  theme_bw() + #theme(legend.position = "none") +
  ggtitle("Сравнение используемых слов на факультете экономики
\n между мальчиками и девушками")
dev.off()
```


Коллокации

library(quanteda)

collocations <-

textstat_collocations(text, size = 2:3)



rank	collocation	count	length	lambda	z
1	главный герой	51	2	5.594665	26.268972
2	друг друг	36	2	5.457287	23.975265
3	самый дело	30	2	4.817666	20.631267
9	очень понравиться	21	2	3.518055	14.267273
4	главный героиня	19	2	6.262405	18.490280
8	год назад	15	2	6.308578	14.505164
5	сей пора	13	2	7.027718	16.648241

LSA - семантическая близость слов

```
library(lsa)
tdm<-as.TermDocumentMatrix(dtmw)
Isa_space<-lsa(tdm, dims=dimcalc_share())
Isa_word_space<-lsa(dtmw, dims=dimcalc_share())
```

```
tdm_lsa<-as.textmatrix(Isa_space)
tdm_word_lsa<-as.textmatrix(lsa_word_space)
```

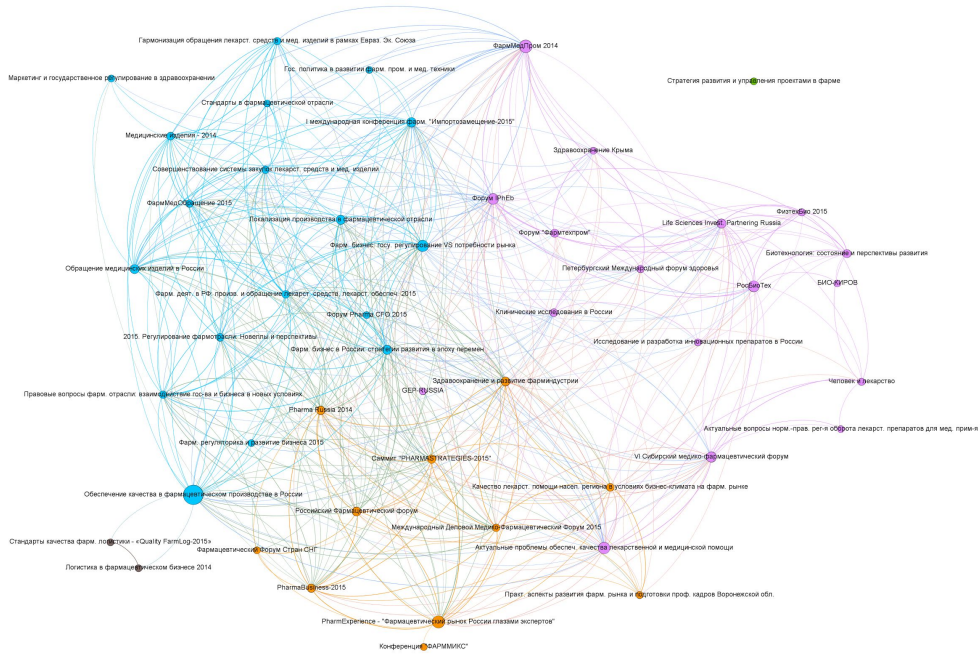
```
tdm_lsa[1:5,1:5] # что присвоено "tdm_lsa" (какое-то значение)
as.matrix(tdm)[1:5,1:5]
t.locs<-Isa_space$tk %>% diag(lsa_space$sk)
plot(t.locs,type="n")
text(t.locs, labels=rownames(Isa_space$tk))
```

```
Isa_space2<-lsa(tdm, dims=2)
t2.locs<-Isa_space2$tk %>% diag(lsa_space2$sk)
plot(t2.locs,type="n")
text(t2.locs, labels=rownames(Isa_space2$tk))
```

```
Isa.distances<-cosine(tdm_lsa) # косинусное расстояние между текстами в
LSA-пространстве
rownames(Isa.distances) <- farm$name
colnames(Isa.distances) <- farm$name
Isa.distances[upper.tri(Isa.distances)] <- NA
diag(Isa.distances)=NA
Isa.matrix <- melt(Isa.distances)
colnames(Isa.matrix) <- c("Source","Target", "Weight")
Isa.matrix<-Isa.matrix[Isa.matrix$Weight > 0, ]
Isa.matrix<-Isa.matrix[!(is.na(Isa.matrix$Weight)), ]
```

```
Isa.matrix$Weight2 <- Isa.matrix$Weight
Isa.matrix$Type <- "Undirected"
write.csv(Isa.matrix, "graph_farmo_lsa.csv")
View(order(Isa.matrix$Weight))
```

```
write.csv(labels_farma, "labels_farma.csv")
```

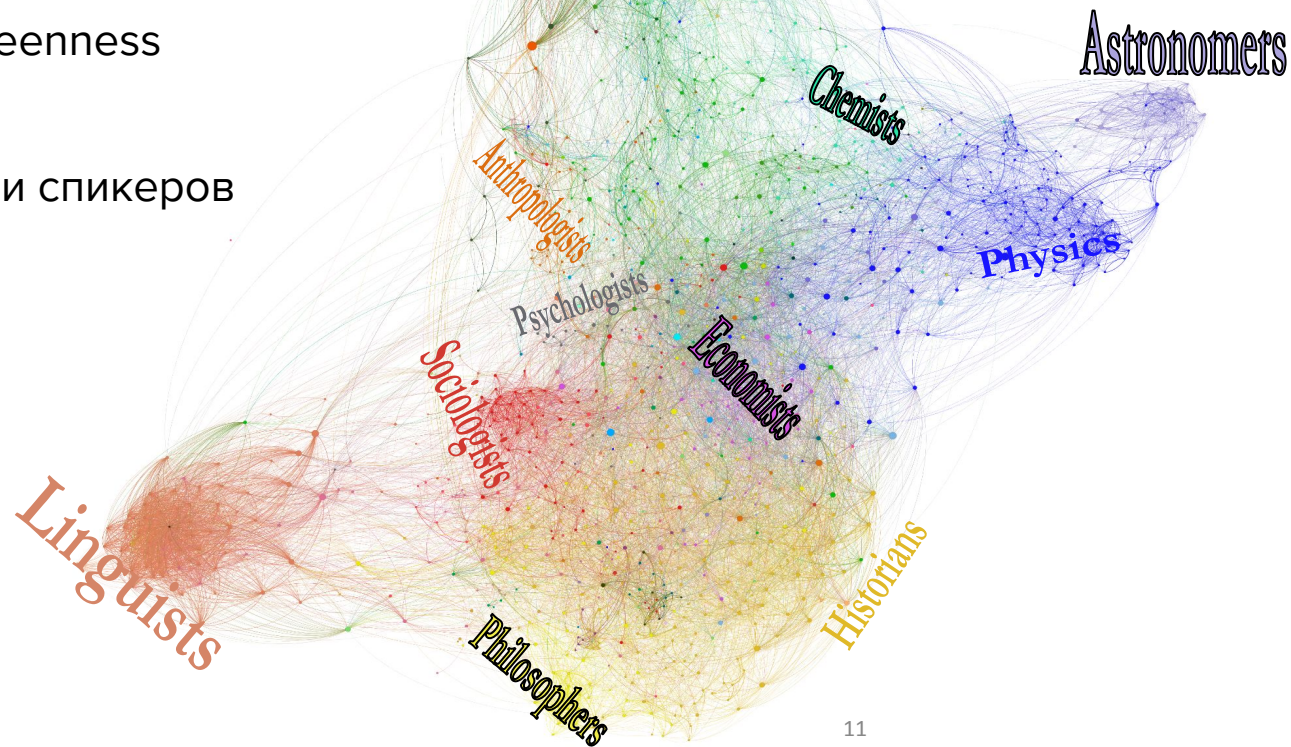


LSA - PostNauka materials

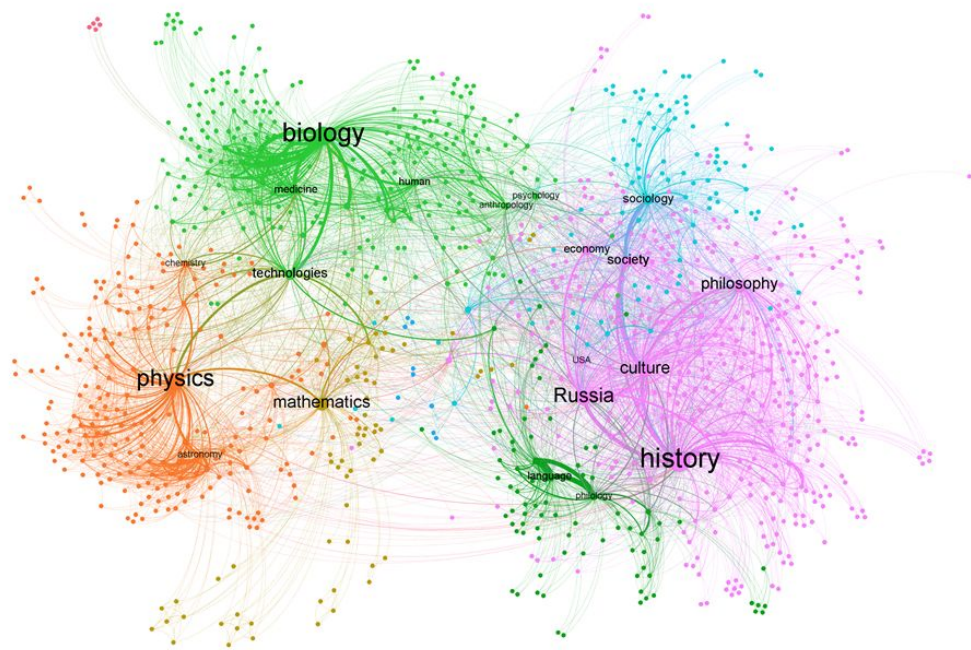
Размер узлов - Betweenness
Centrality

Цвет - специализации спикеров

Связь - LSA метрики

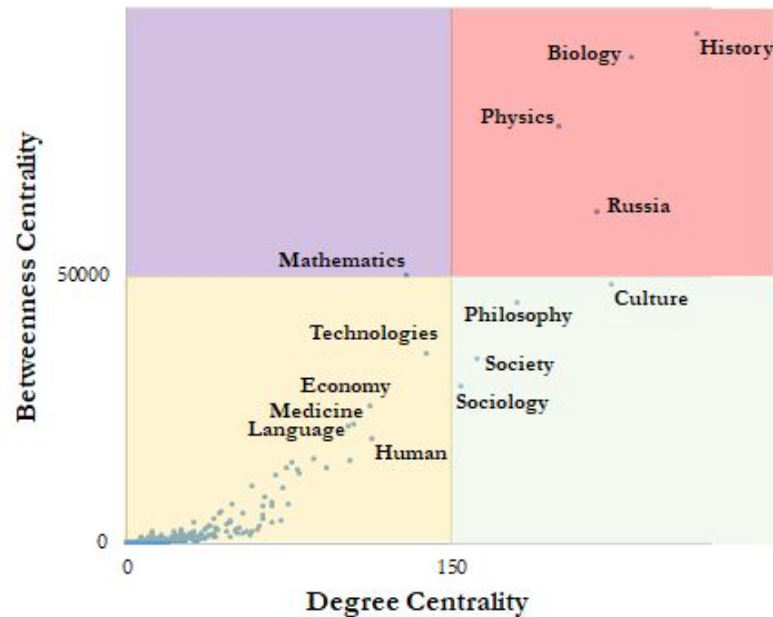


LSA - PostNauka materials



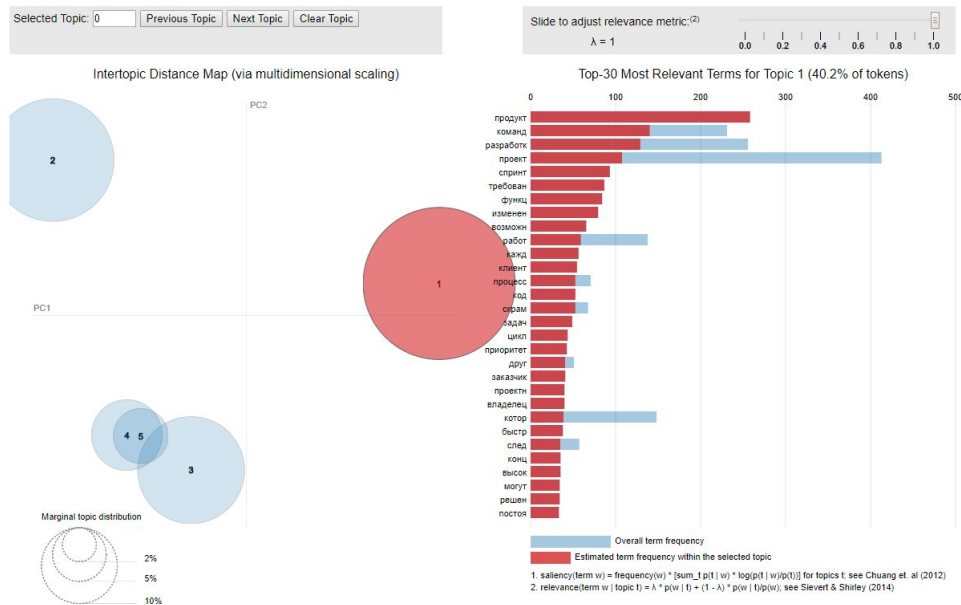
Key-words
centralization structure

- - Globally Central
- - Locally Central
- - Gatekeeper
- - Marginal



LDA — к каким темам относится документ

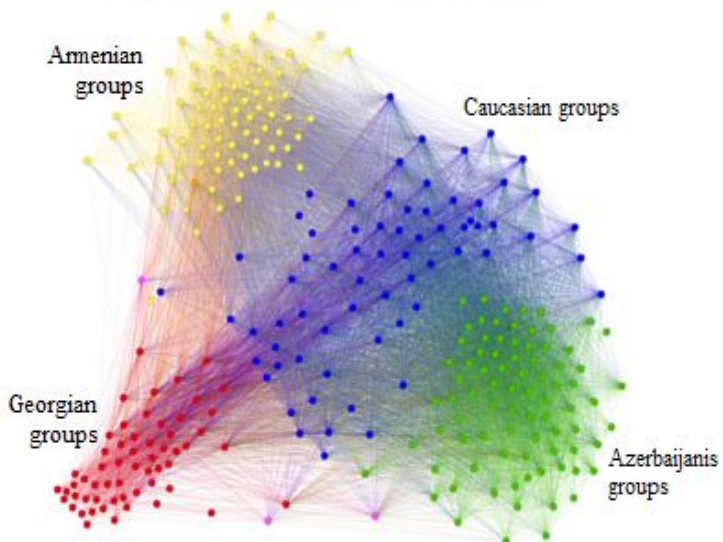
```
library(mallet)
mallet.instances <- mallet.import(id.array = tb$Название.раздела, text.array = corp, stoplist.file = "./data/stopwords.txt")
## настраиваем параметры модели и загружаем данные
topic.model <- MalletLDA(num.topics=5) # количество тем
topic.model$loadDocuments(mallet.instances)
topic.model$setAlphaOptimization(20, 50) # оптимизация гиперпараметров
## собираем статистику: словарь и частотность
vocabulary <- topic.model$getVocabulary() # словарь корпуса
word.freqs <- mallet.word.freqs(topic.model) # таблица частотности слов
## вершина частотного списка (по документной частоте)
head(word.freqs[order(word.freqs$doc.freq, decreasing=T),],30)
## параметр — количество итераций
topic.model$train(1000)
## выбор наилучшей темы для каждого токена
topic.model$maximize(10)
### LDA: выгрузка результатов
## таблица распределения тем по документам
doc.topics <- mallet.doc.topics(topic.model, smoothed=TRUE, normalized=TRUE)
## таблица распределения слов по темам
topic.words <- mallet.topic.words(topic.model, smoothed=TRUE, normalized=TRUE)
## метки для тем (по трем главным словам)
topic.labels <- mallet.topic.labels(topic.model, topic.words, 5)
```



[Ссылка на интерактивный граф](#)

LDA: этнические группы в ВК

Network of groups based on overlapping membership (Jaccard index)



Family values and gender stereotypes

