

Digital humanities

Cvičení: analýza textů

Jindřich Marek

Práce s textem z digitální knihovny

- vyhledání knihy v digitální knihovně KNAV
- stažení metadat
- z metadat přes API extrakce textu
- lemmatizace textu
- vizualizace textu
 - wordcloud

Knihy



<https://kramerius.lib.cas.cz/view/uuid:f24658de-c158-43b0-a41b-5d1761c85cb2?page=uuid:b41cee53-72f0-415b-983f-04f630e9c046>

Digitální knihovna Akademie věd ČR

Hledat v celé digitální knihovně


Sbírký Procházet Informace FAQ Přihlásit

Hledat v dokumentu

Strana 1 z 164

LENKA VESELÁ-PRUDKOVÁ

Židé a česká společnost v zrcadle literatury



ISBN 80-7106-430-0

Autor
[Veselá-Prudková, Lenka](#)

Nakladatelské údaje
Praha: Lidové noviny, 2003

Typ dokumentu
[Kniha](#)

Klíčová slova
[Sociální procesy](#)
[Antijudaismus](#)
[Antisemitismus](#)
[Literární historie](#)
[Občanská společnost](#)
[Písemné památky](#)
[Židé](#)

<https://kramerius.lib.cas.cz/search/api/v5.0/item/uuid:f24658de-c158-43b0-a41b-5d1761c85cb2/children>

```
JSON  Surová data  Hlavičky
Uložit  Kopírovat  Sbalit vše  Rozbalit vše  Filtr JSON
0:
  datanode: true
  pid: "uuid:b41cee53-72f0-415b-983f-04f630e9c046"
  model: "page"
  details:
    type: "FrontCover"
    pagenumber: "[1a] \n"
    title: "[1a]"
    root_title: "Židé a česká společnost v zrcadle literatury"
    root_pid: "uuid:f24658de-c158-43b0-a41b-5d1761c85cb2"
    policy: "public"
1:
  datanode: true
  pid: "uuid:87fc1ed8-6d2f-4608-b992-68a2d9f180ca"
  model: "page"
  details:
    type: "FrontEndSheet"
    pagenumber: "[1b] \n"
    title: "[1b]"
    root_title: "Židé a česká společnost v zrcadle literatury"
    root_pid: "uuid:f24658de-c158-43b0-a41b-5d1761c85cb2"
    policy: "public"
2:
  datanode: true
  pid: "uuid:a76eb020-fbbe-4f77-83f5-0d6cbe1b3181"
  model: "page"
  details:
    pagenumber: "[1] \n"
    title: "[1]"
    root_title: "Židé a česká společnost v zrcadle literatury"
    root_pid: "uuid:f24658de-c158-43b0-a41b-5d1761c85cb2"
    policy: "public"
3:
  datanode: true
  pid: "uuid:729cbe79-a2f4-42c8-a71c-913229f58e9d"
  model: "page"
  details:
    type: "FrontJacket"
    pagenumber: "[2] \n"
    title: "[2]"
    root_title: "Židé a česká společnost v zrcadle literatury"
    root_pid: "uuid:f24658de-c158-43b0-a41b-5d1761c85cb2"
    policy: "public"
4:
```

<https://kramerus.lib.cas.cz/search/api/v5.0/item/uuid:b41cee53-72f0-415b-983f-04f630e9c046/streams>

JSON Surová data Hlavičky

Uložit Kopírovat Sbalit vše Rozbalit vše Filtr JSON

```
▼ BIBLIO_MODS:  
  label: "Metadata Object Description"  
  mimeType: "text/xml"  
▼ IMG_PREVIEW:  
  label: "Preview of this object"  
  mimeType: "image/jpeg"  
▼ IMG_FULL:  
  label: "Presentable version of RAW"  
  mimeType: "image/jpeg"  
▼ ALTO:  
  label: "ALTO for this object"  
  mimeType: "text/xml"  
▼ TEXT_OCR:  
  label: "OCR for this object"  
  mimeType: "text/plain"  
▼ IMG_THUMB:  
  label: "Thumbnail of this object"  
  mimeType: "image/jpeg"  
▼ DC:  
  label: "Dublin Core Record for this object"  
  mimeType: "text/xml"
```

https://kramerus.lib.cas.cz/search/api/v5.0/item/uuid:b41cee53-72f0-415b-983f-04f630e9c046/streams/IMG_FULL

LENKA VESELÁ-PRUDKOVÁ

Židé a česká společnost
v zrcadle literatury



KNIZNICE
DĚJIN
A SOUČASNOSTI

https://kramerus.lib.cas.cz/search/api/v5.0/item/uuid:b41cee53-72f0-415b-983f-04f630e9c046/streams/TEXT_OCR

LENKA VESELÁ-PRU OKOVÁ
Židé a česká společnost
v zrcadle literatury
A SOUČASNOSTI



stazeni_textu_Vesela_Zide.R (1. část)

```
# Load the jsonlite and httr libraries for handling JSON data and making HTTP requests.
library(jsonlite)
library(httr)

# Define the URL for the Kramerius API.
kramerius5_api <- "https://kramerius.lib.cas.cz/search/api/v5.0/item/"

# Define the UUID of the title for which you want to retrieve OCR text.
UUID_titul <- "uuid:f24658de-c158-43b0-a41b-5d1761c85cb2"

# Construct the URL for retrieving children items of the specified title.
request <- paste0(kramerius5_api, UUID_titul, "/children")

# Send a GET request to the Kramerius API to retrieve the JSON response.
response <- GET(request)

# Parse the JSON response into R data structures.
parsed_json <- fromJSON(content(response, "text"))
```

stazeni_textu_Vesela_Zide.R (2. část)

```
# Extract URLs for retrieving OCR text streams.
urls <- paste0(kramerius5_api, parsed_json$pid, "/streams/TEXT_OCR")

# Retrieve OCR text from each URL and store it in a list.
text_list <- lapply(urls, function(url) {
  response <- GET(url)
  content(response, "text", encoding = "UTF-8")
})

# Convert the list of OCR text into a character vector.
text_list <- as.character(text_list)

# Write the OCR text to a text file named "stazeny_text_Vesela_Zide.txt" in the current working
directory.
writeLines(text_list, "stazeny_text_Vesela_Zide.txt")
```


wordcloud_Vesela_Zide.R – 1. část (lemmatizace)

```
# Load necessary libraries for data
manipulation and visualization.
library(readxl)
library(readr)
library(dplyr)
library(udpipe)
library(RColorBrewer)
library(wordcloud2)
library(htmlwidgets)

# Get the directory of the currently executing
script.
this_dir <- dirname(parent.frame(2)$ofile)
setwd(this_dir)
```

```
# Download and load the UDPipe model for Czech
language.
udmodel <- udpipes_download_model(language =
"czech")
udmodel <- udpipes_load_model(file =
udmodel$file_model)

# Read the OCR text from the downloaded file.
file <- "stazeny_text_Vesela_Zide.txt"
y <- read_lines(file, skip = 0, n_max = -1L)

# Perform tokenization, tagging, and
lemmatization using UDPipe.
x <- udpipes_annotate(udmodel, y)
lv <- as.data.frame(x, detailed = TRUE)

# Convert lemma column to lowercase.
lv$lemma <- tolower(lv$lemma)
```

wordcloud_Vesela_Zide.R – 2. část (vyčištění lemmat, wordcloud)

```
# Read Czech stopwords from file.
stopwords <- readLines("stopwords_cz.txt")

# Define a function to remove stopwords from text.
remove_stopwords <- function(text, stopwords) {
  words <- unlist(strsplit(text, "\\s+"))
  filtered_words <- words[!tolower(words) %in% stopwords]
  cleaned_text <- paste(filtered_words, collapse = " ")
  return(cleaned_text)
}

# Apply the remove_stopwords function to remove stopwords.
lv$lemma <- sapply(lv$lemma, remove_stopwords, stopwords)

# Remove leading and trailing whitespaces.
lv$lemma <- trimws(lv$lemma)
```

```
# Convert the lemma column to a factor.
lv$lemma <- as.factor(lv$lemma)

# Count the occurrences of each lemma and arrange in descending
order.
lv_framed <- lv %>% count(lemma, sort = TRUE)

# Convert lemma column to character and filter out lemmas with
less than 3 characters.
lv_framed$lemma <- as.character(lv_framed$lemma)
df <- lv_framed %>% filter(nchar(lemma) > 2)

# Create a word cloud using wordcloud2 package.
hw <- wordcloud2(data = df, size = 1.6, color = 'random-dark')

# Save the word cloud as an HTML file.
saveWidget(hw, "wordcloud_vysledek_Vesela_Zide.html",
selfcontained = FALSE)
```


tematicke_modelovani.R – 1. část (načtení knihoven, stažení modelu)

```
# Load required libraries
libraries <- c(
  "quanteda", "quanteda.textstats", "quanteda.textplots", "topicmodels", "readtext",
  "tidytext", "dplyr", "jsonlite", "udpipe", "stopwords", "readxl", "stringr",
  "LDAvis", "wordcloud", "ggplot2", "gridExtra", "RColorBrewer", "forcats", "readr"
)
invisible(lapply(libraries, library, character.only = TRUE))

# Set working directory
this_dir <- dirname(parent.frame(2)$ofile)
setwd(this_dir)

# Download and load UD model for Czech
udmodel <- udpipes_download_model(language = "czech")
udmodel <- udpipes_load_model(file = udmodel$file_model)
```


tematicke_modelovani.R – 2. část (funkce: pouze podstatná jména)

```
# Function to process text and save nouns to a file
process_text <- function(file, output_file) {
  text <- read_lines(file, skip = 0, n_max = -1L)
  annotations <- udpipe_annotate(udmodel, text)
  nouns <- as.data.frame(annotations, detailed = TRUE) %>%
    filter(upos == "NOUN") %>%
    mutate(lemma = tolower(lemma)) %>%
    pull(lemma)
  writeLines(nouns, output_file)
}

# Process text files and extract nouns
process_text("stazeny_text_Vesela_Zide.txt", "stazeny_text_Vesela_Zide_nouns.txt")
process_text("stazeny_text_magie.txt", "stazeny_text_magie_nouns.txt")
```

tematicke_modelovani.R – 3. část (očištění a tokenizace)

```
# Load stopwords
stopwords_cs <- stopwords::stopwords("cs", source = "stopwords-iso")
stopwords_cz <- c(stopwords_cs, readLines("stopwords_cz.txt"))

# Read text files into corpus
y <- readtext(c("stazeny_text_Vesela_Zide_nouns.txt", "stazeny_text_magie_nouns.txt"))
data_txt <- corpus(y)

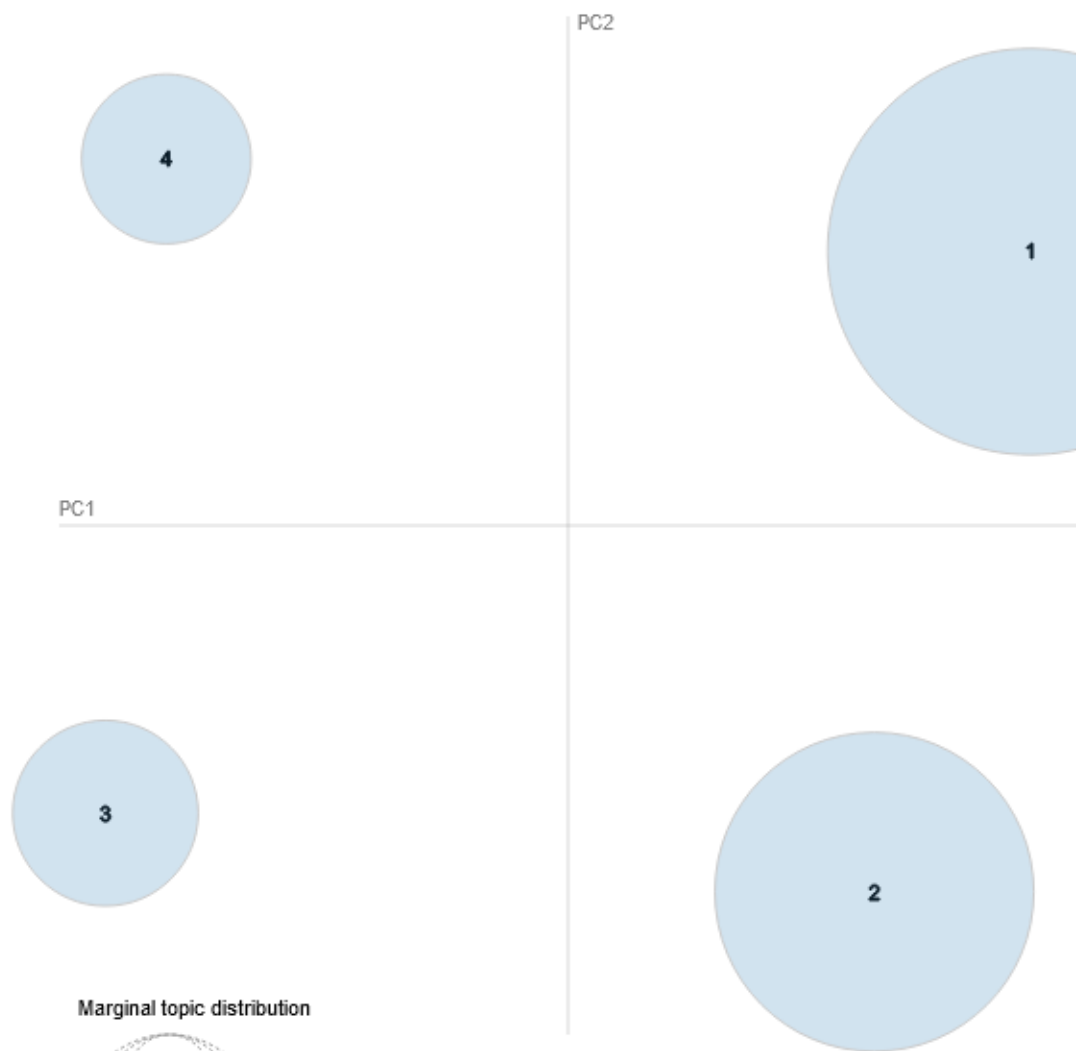
# Tokenization and preprocessing
tokens <- tokens(data_txt) %>%
  tokens_remove(min_nchar = 3) %>%
  tokens_remove(pattern = c("*-*", "und", "der", "býti")) %>%
  tokens_remove(stopwords_cz)
```

tematicke_modelovani.R – 4. část (tematické modelování)

```
# Create document-feature matrix  
dfm <- dfm(tokens)  
dfm_df <- as.data.frame(t(dfm))
```

```
# Topic modeling  
dtm <- convert(dfm, to = "topicmodels")  
set.seed(1234)  
topic_model <- LDA(dtm, method = "Gibbs", k = 4, control = list(alpha = 0.1))
```

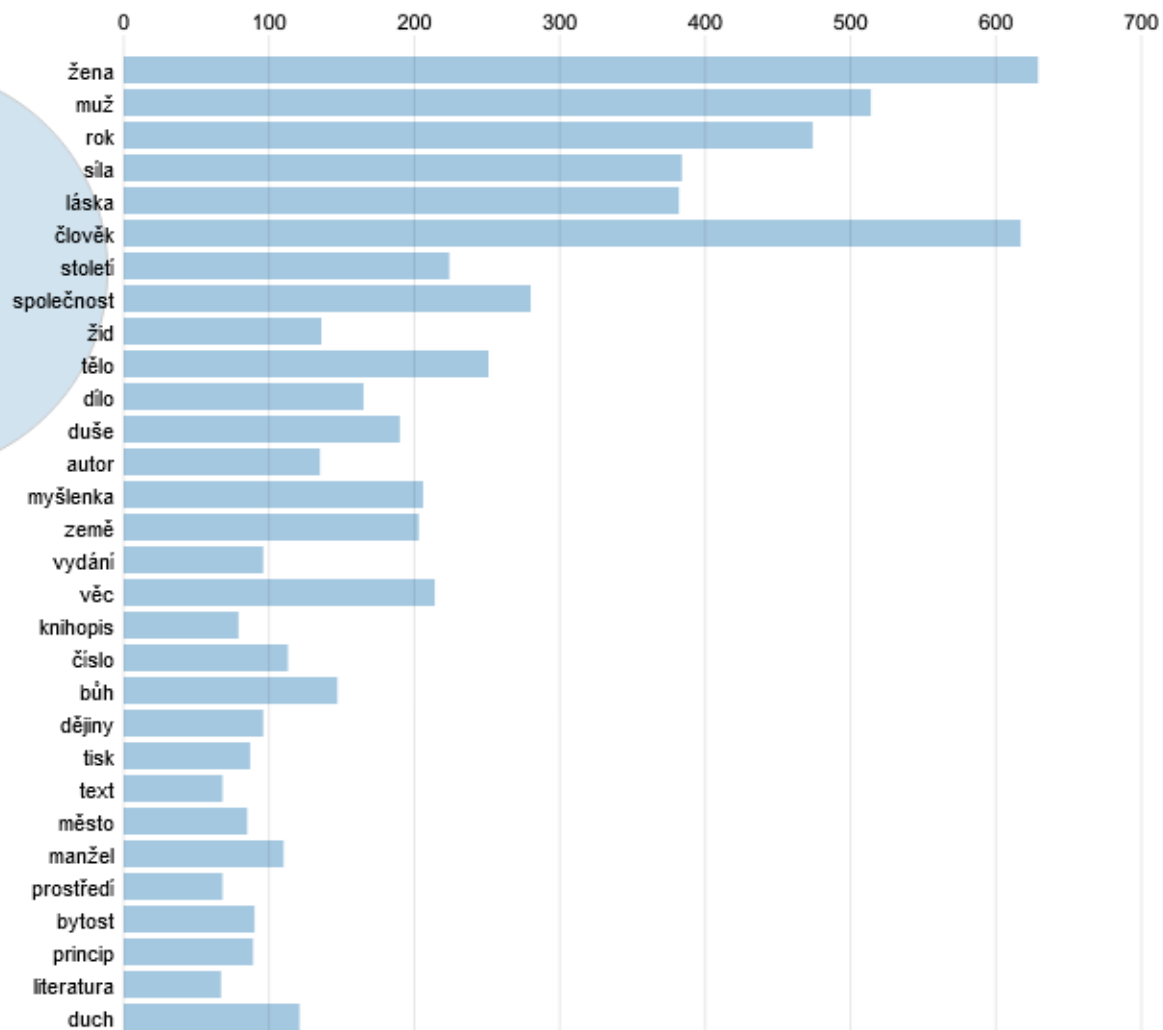

Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



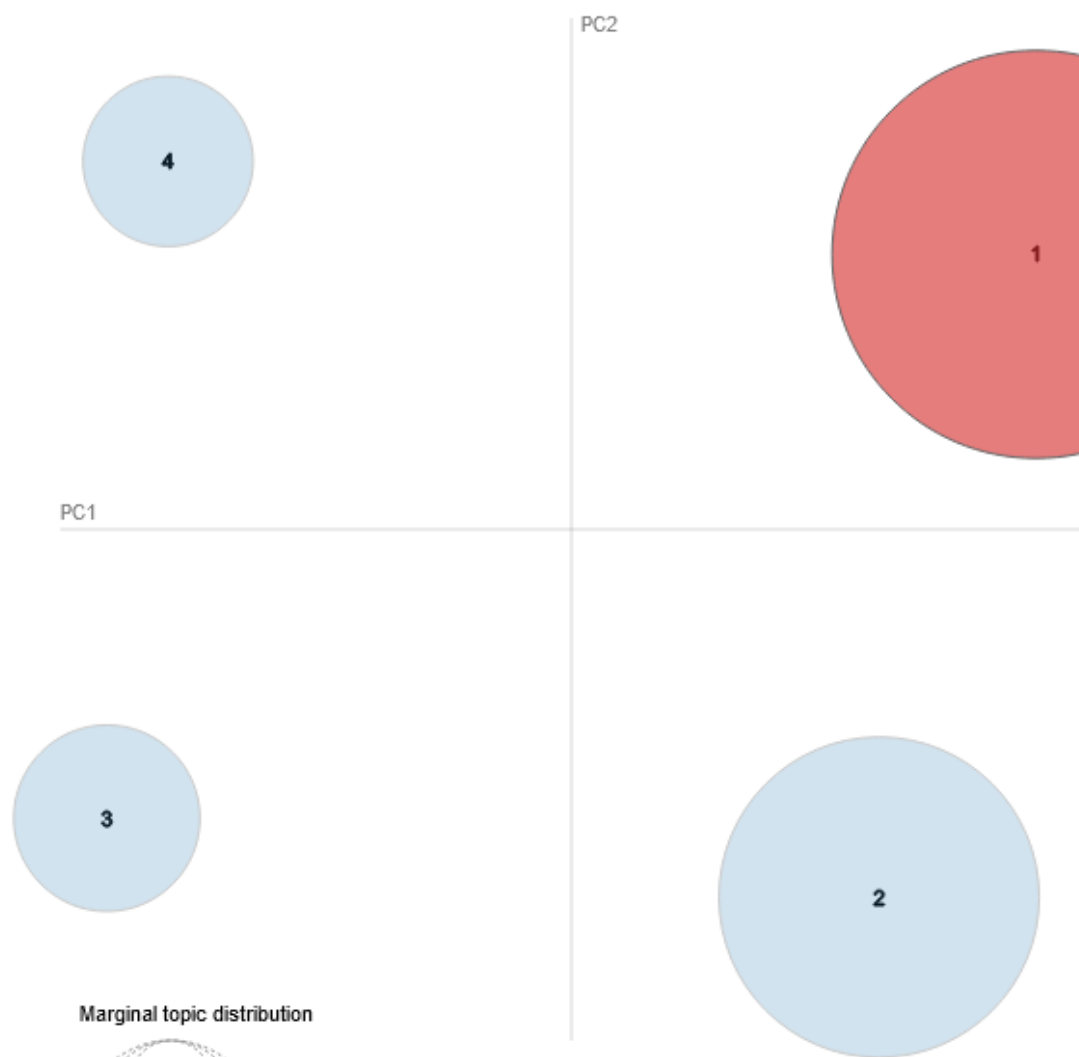
Top-30 Most Salient Terms¹



Overall term frequency
 Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
 2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)

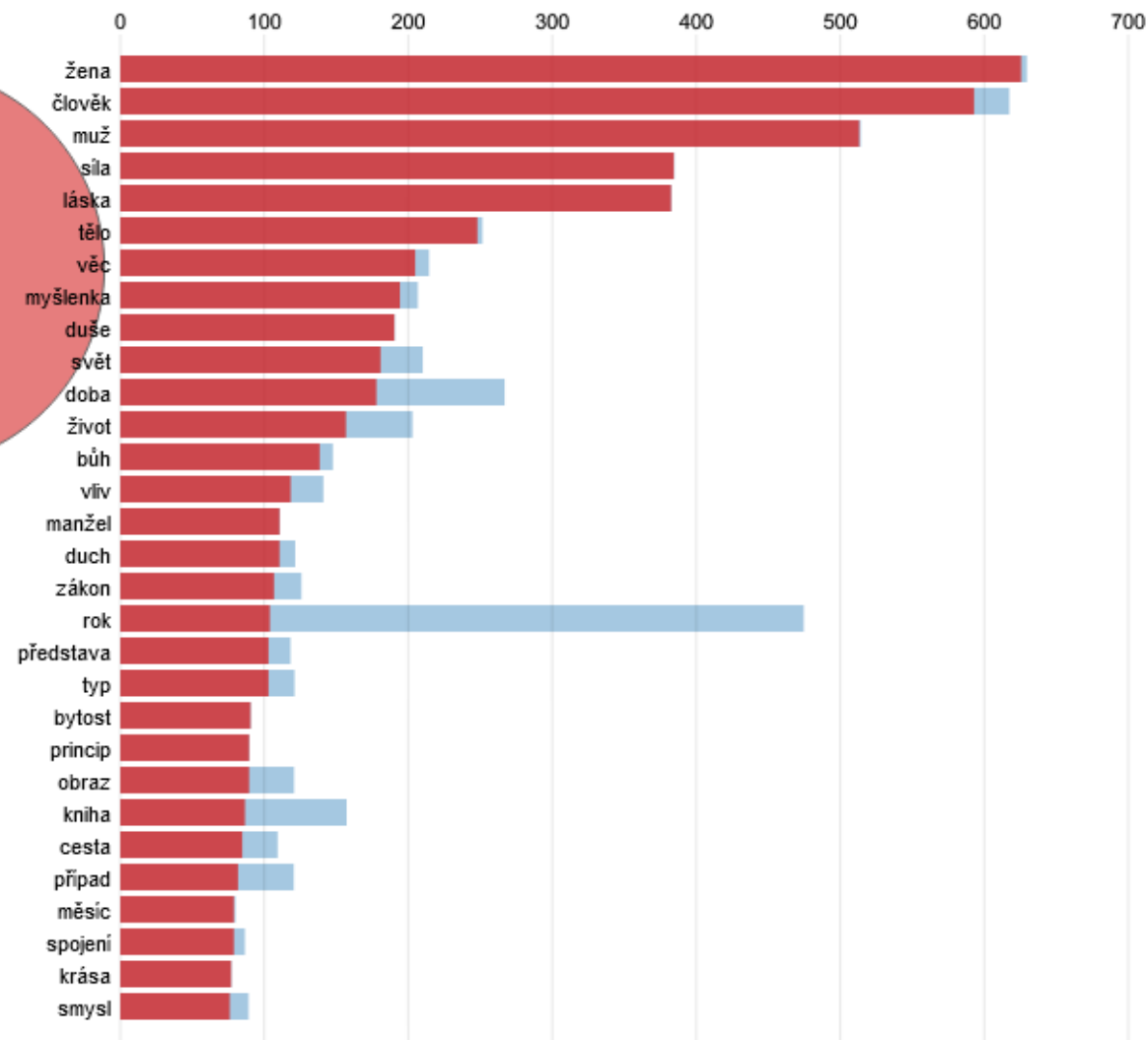
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



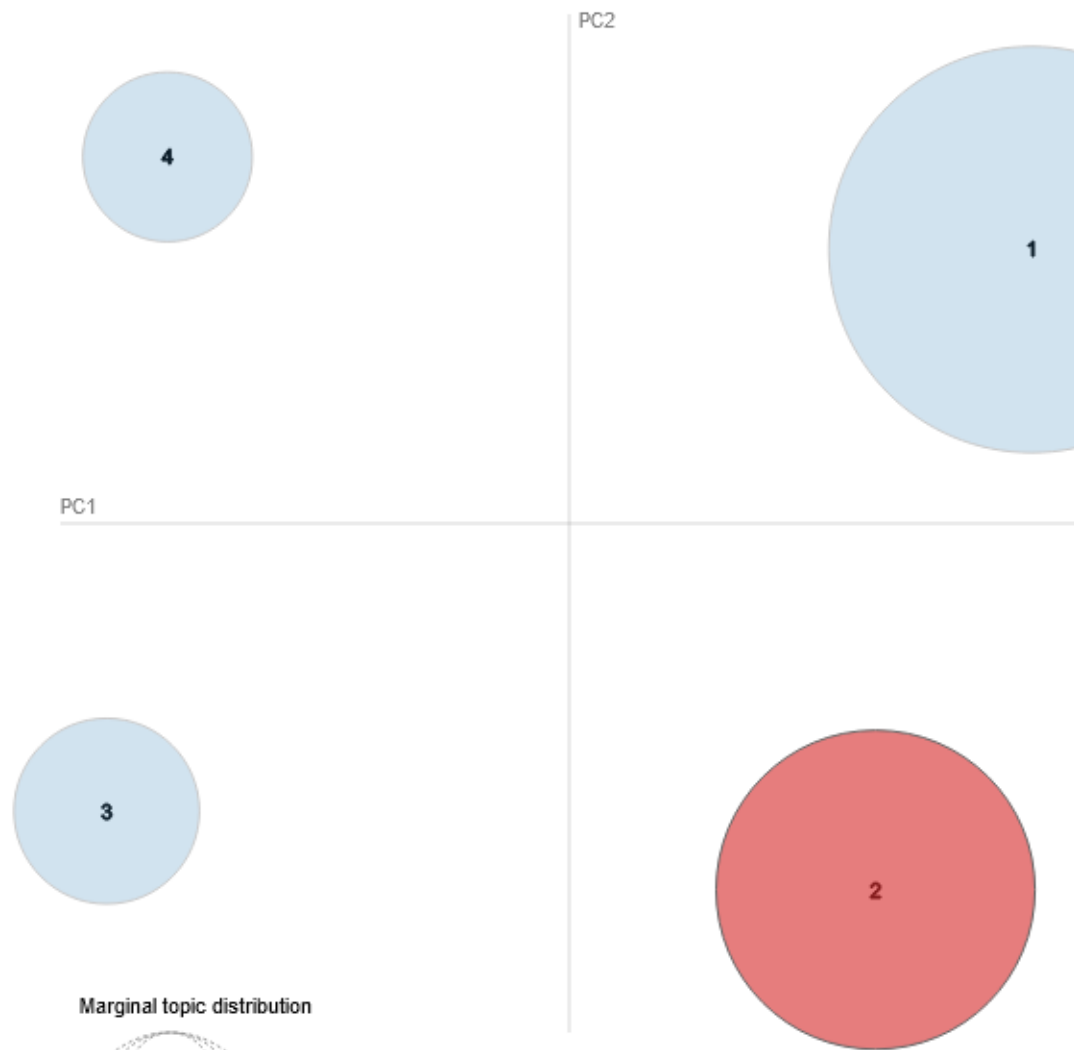
Top-30 Most Relevant Terms for Topic 1 (50% of tokens)



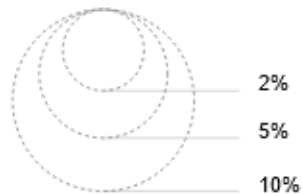
Overall term frequency (blue bar)
 Estimated term frequency within the selected topic (red bar)

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
 2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)

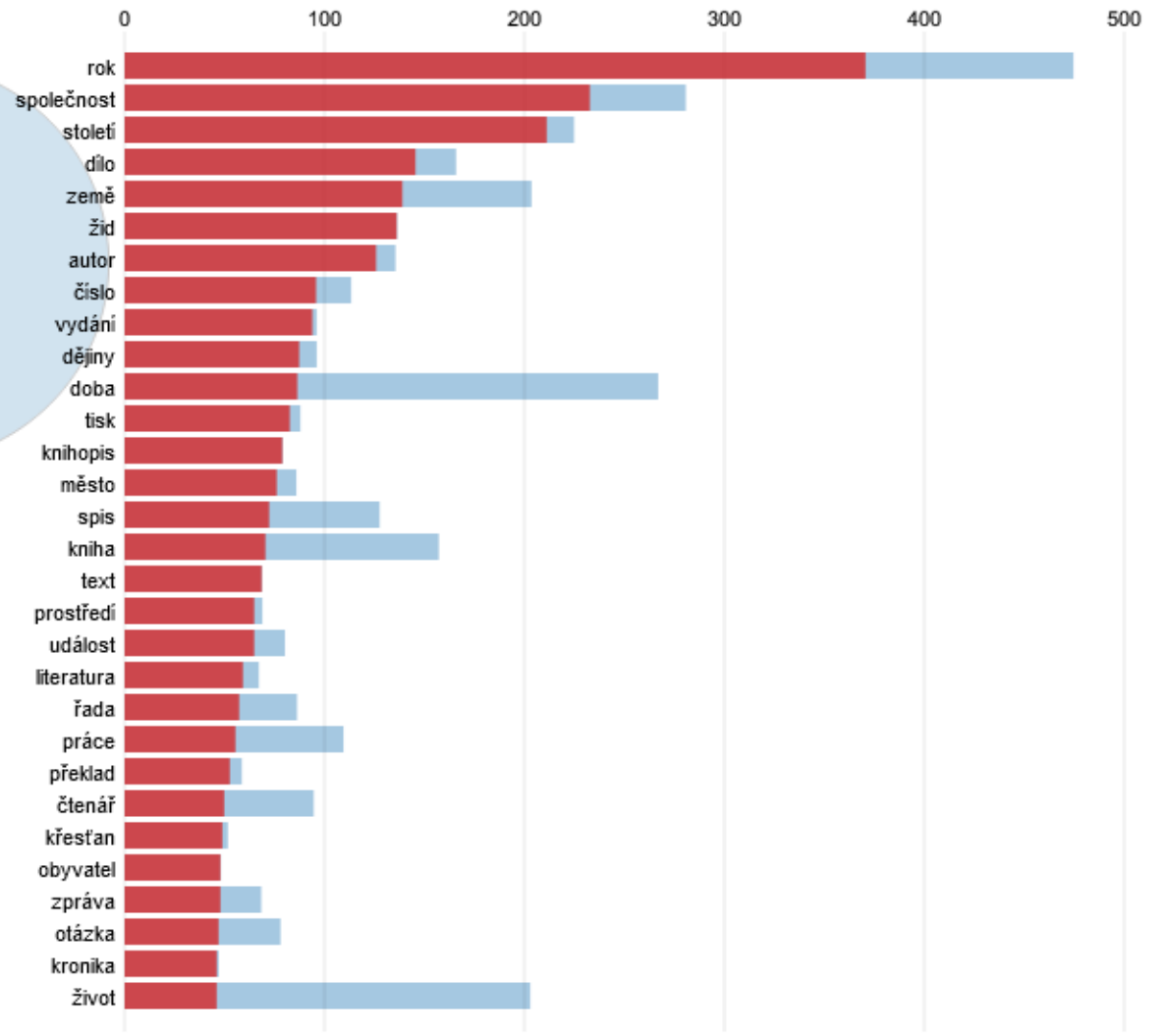
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



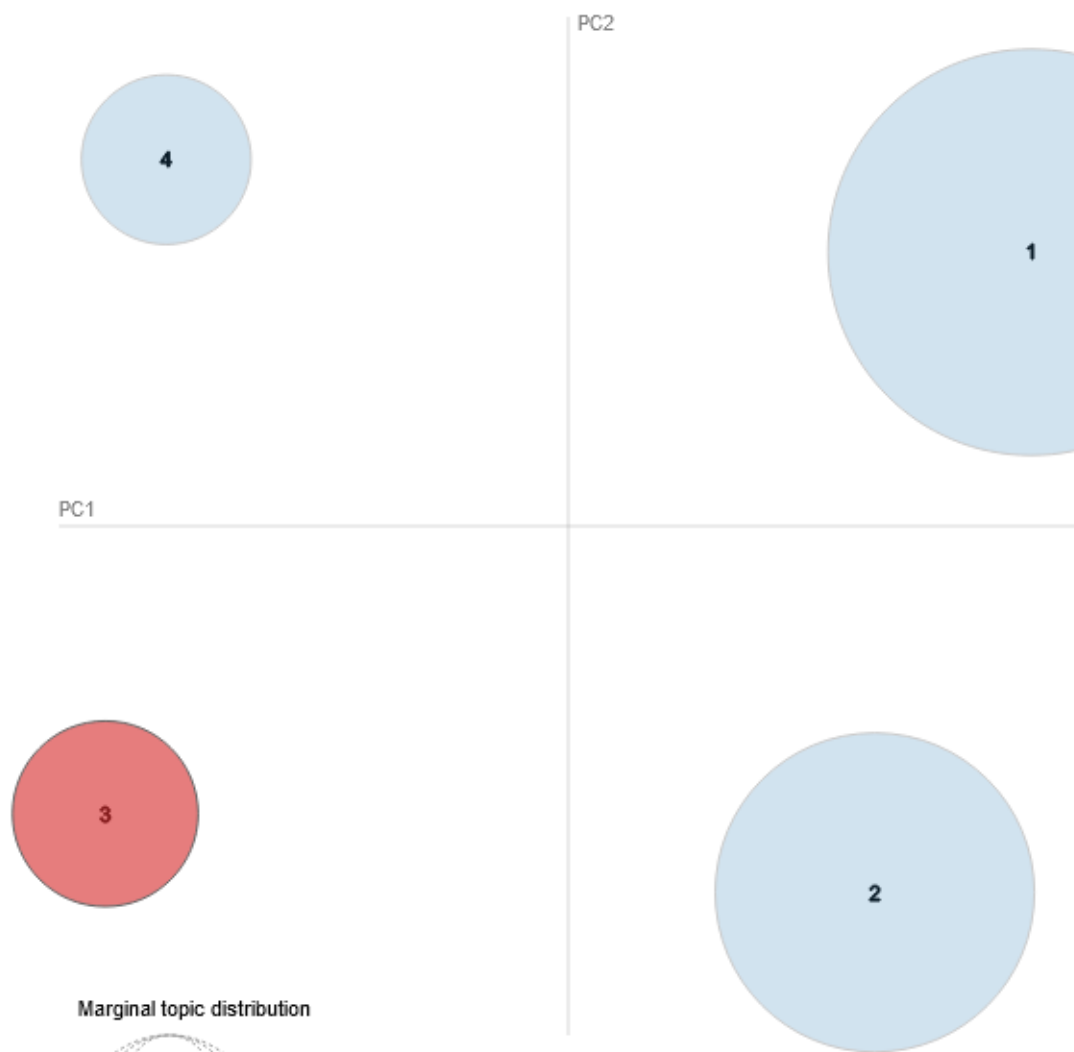
Top-30 Most Relevant Terms for Topic 2 (30.9% of tokens)



Overall term frequency
 Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
 2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)

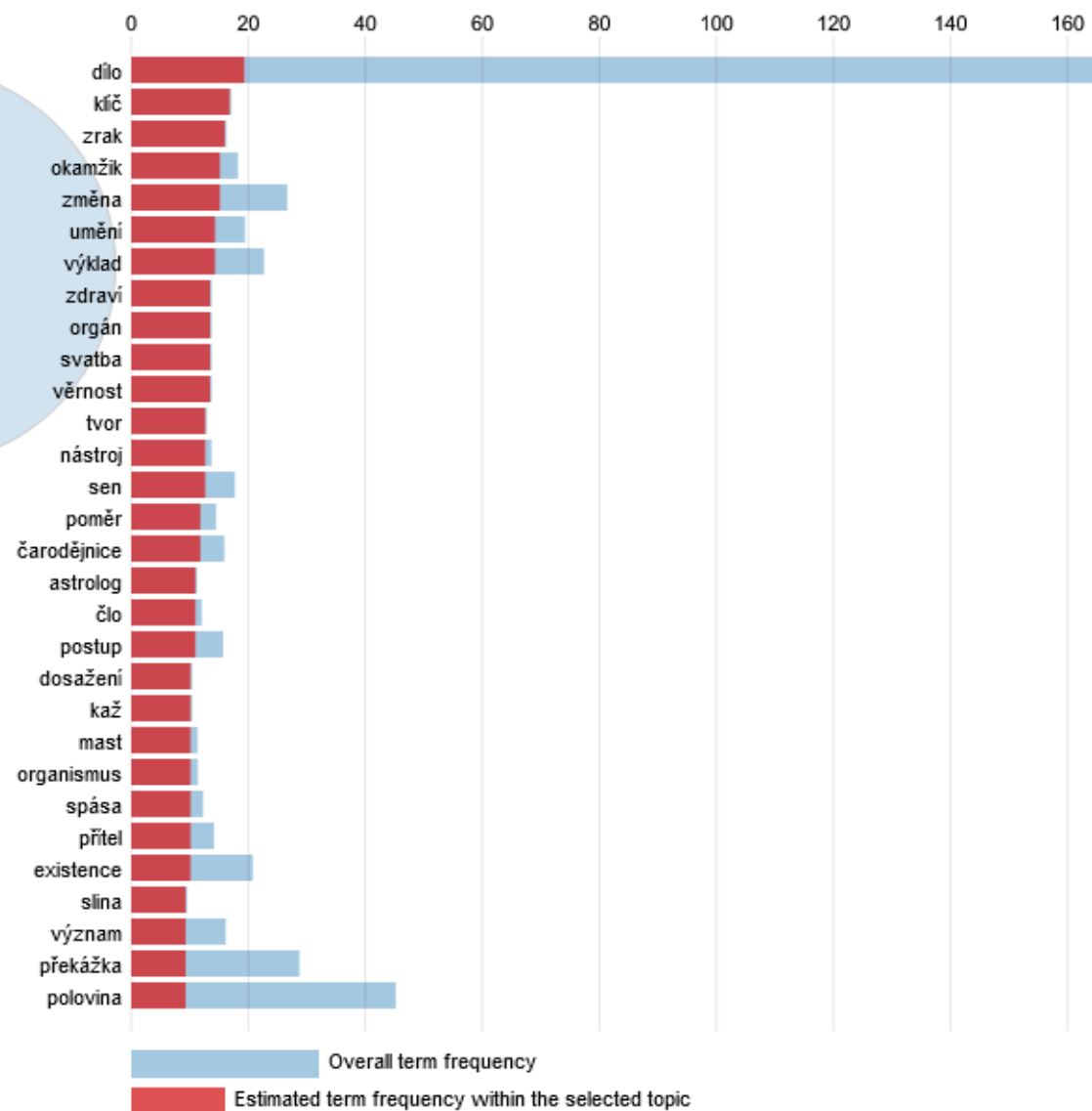
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution

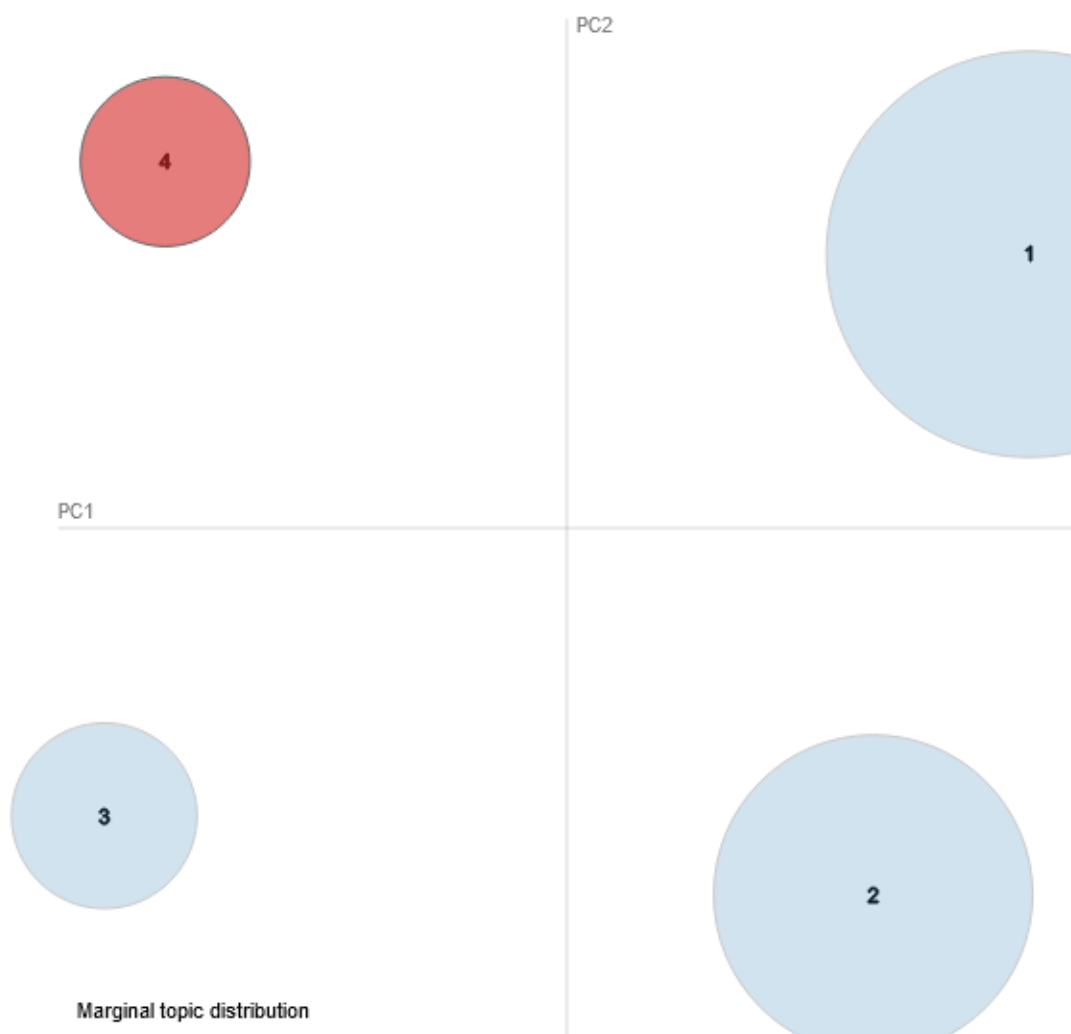


Top-30 Most Relevant Terms for Topic 3 (10.4% of tokens)



1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
 2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)

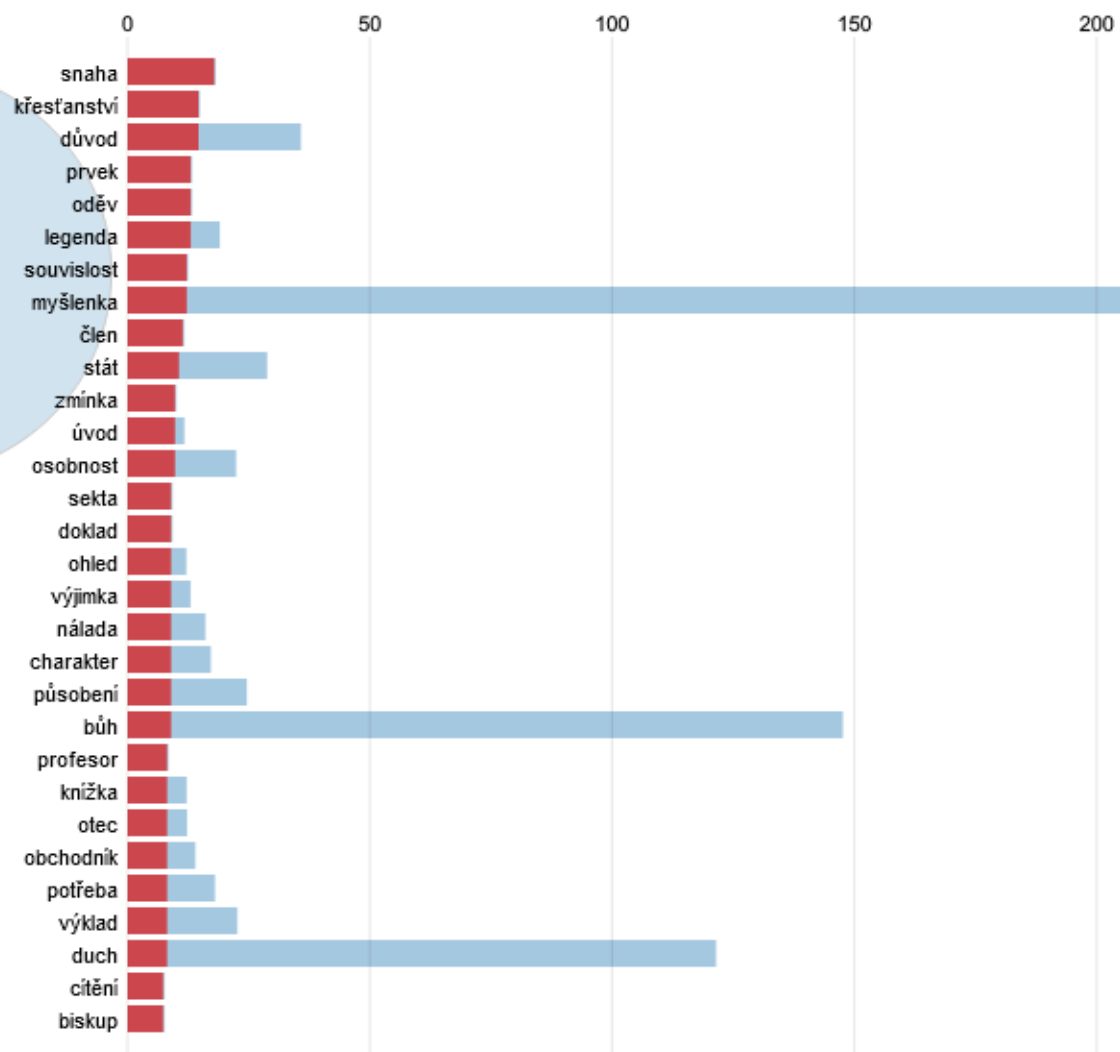
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 4 (8.7% of tokens)



Overall term frequency (blue bar)
 Estimated term frequency within the selected topic (red bar)

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
 2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)