# Text analysis 1

Lukáš Lehotský

"text analysis is just a fancy and convoluted way how to obtain independent or dependent variable"

Inaki Sagarzazu

# Concepts

# Bag of words

# Bag of words

- The quick brown fox jumps over the lazy dog

| Word | Occurrence |
|------|------------|
| brown | 1 |
| dog | 1 |
| fox | 1 |
| jumps | 1 |
| lazy | 1 |
| over | 1 |
| quick | 1 |
| the | 2 |

# Co-occurrence

# Co-occurrence

- The quick **brown fox** jumps over the lazy dog. **Brown dog** sleeps well.

| Word | Sentence 1 | Sentence 2 |
|------|------------|------------|
| brown | 1 | 1 |
| dog | 1 | 1 |
| fox | 1 | |
| jumps | 1 | |
| lazy | 1 | |
| over | 1 | |
| quick | 1 | |
| sleeps | | 1 |
| the | 2 | |
| well | | 1 |

# Co-locations and N-grams

# Co-locations/n-grams

- Established phrases – usually occur together and form a meaning
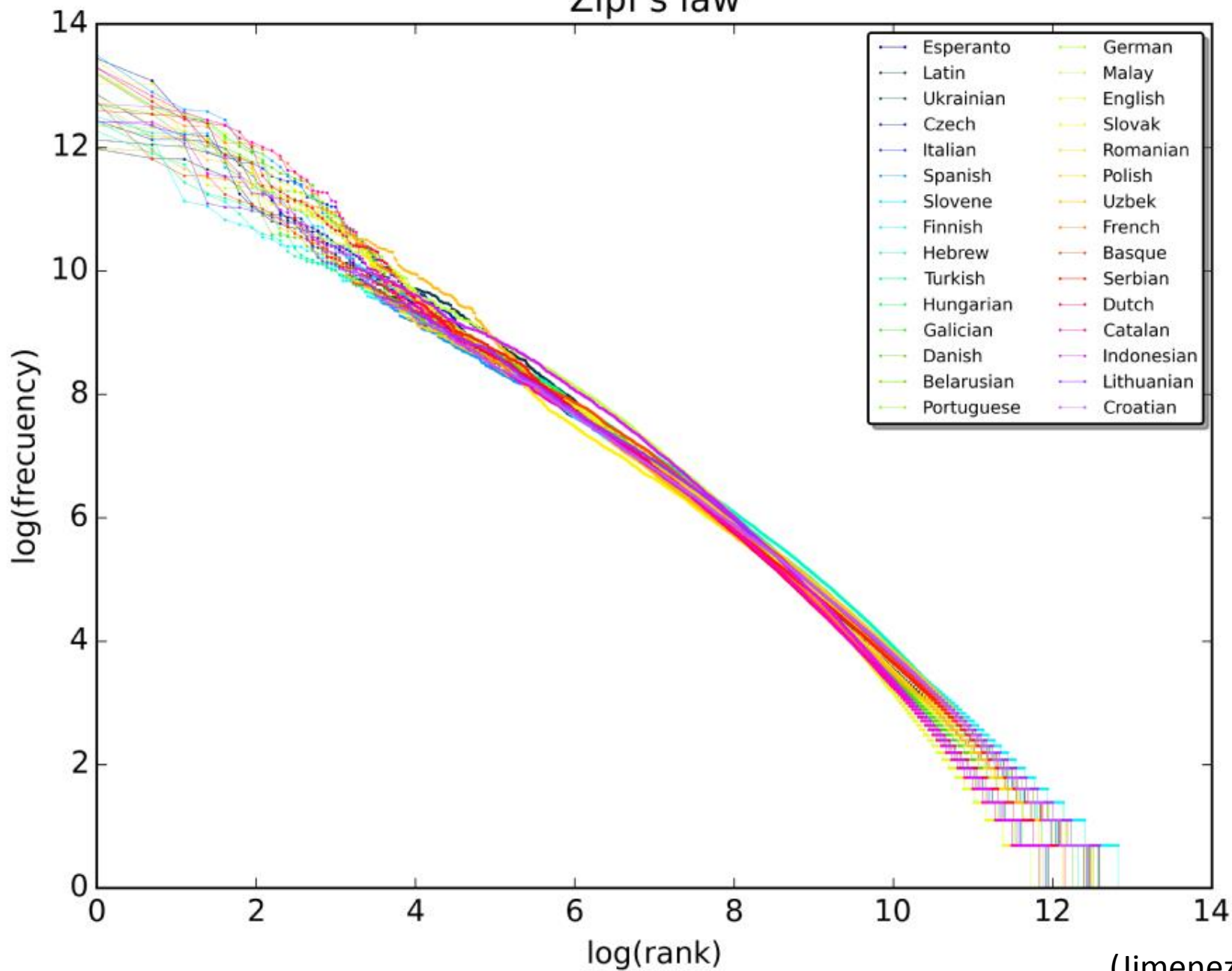
Ministry of the Environment

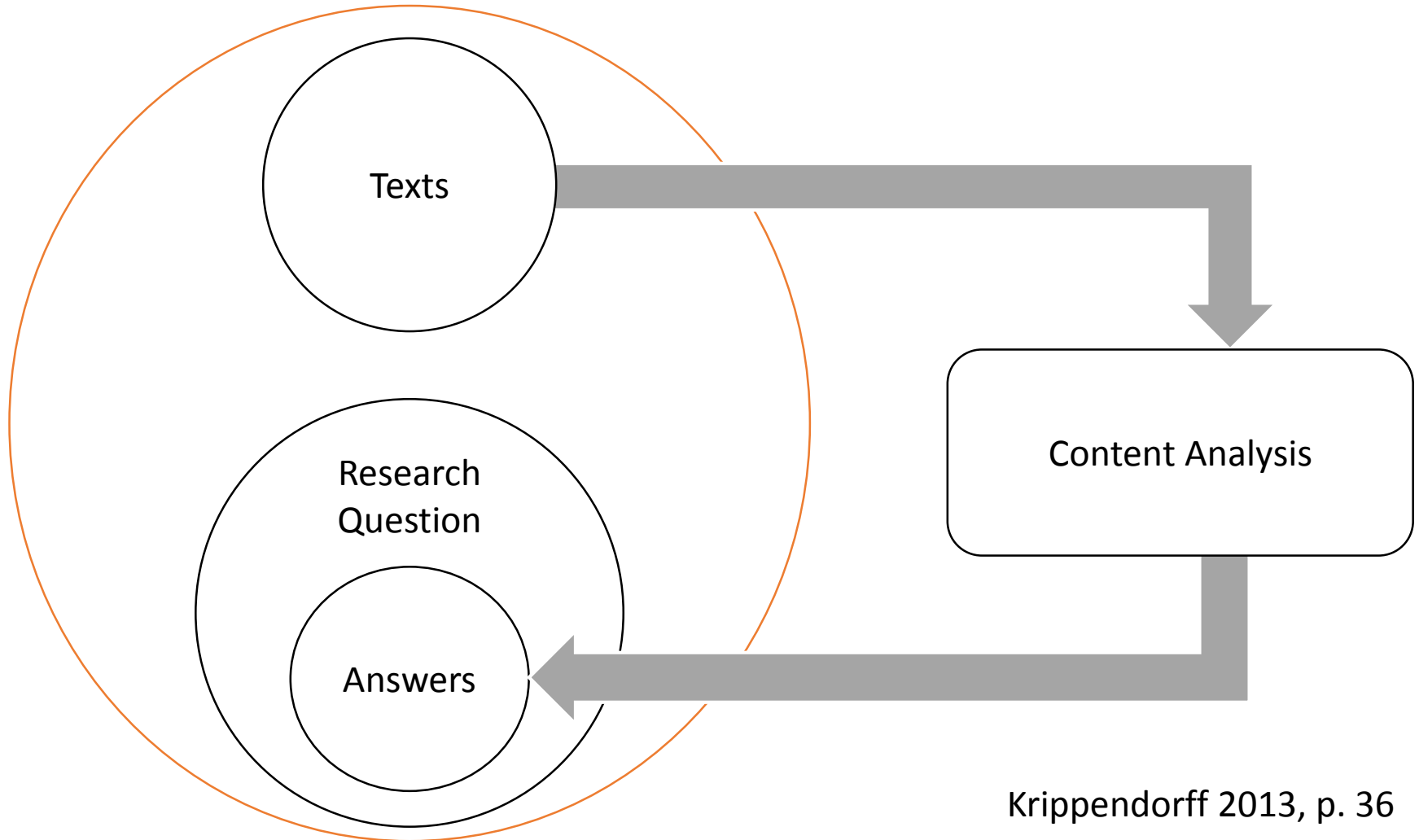European Union

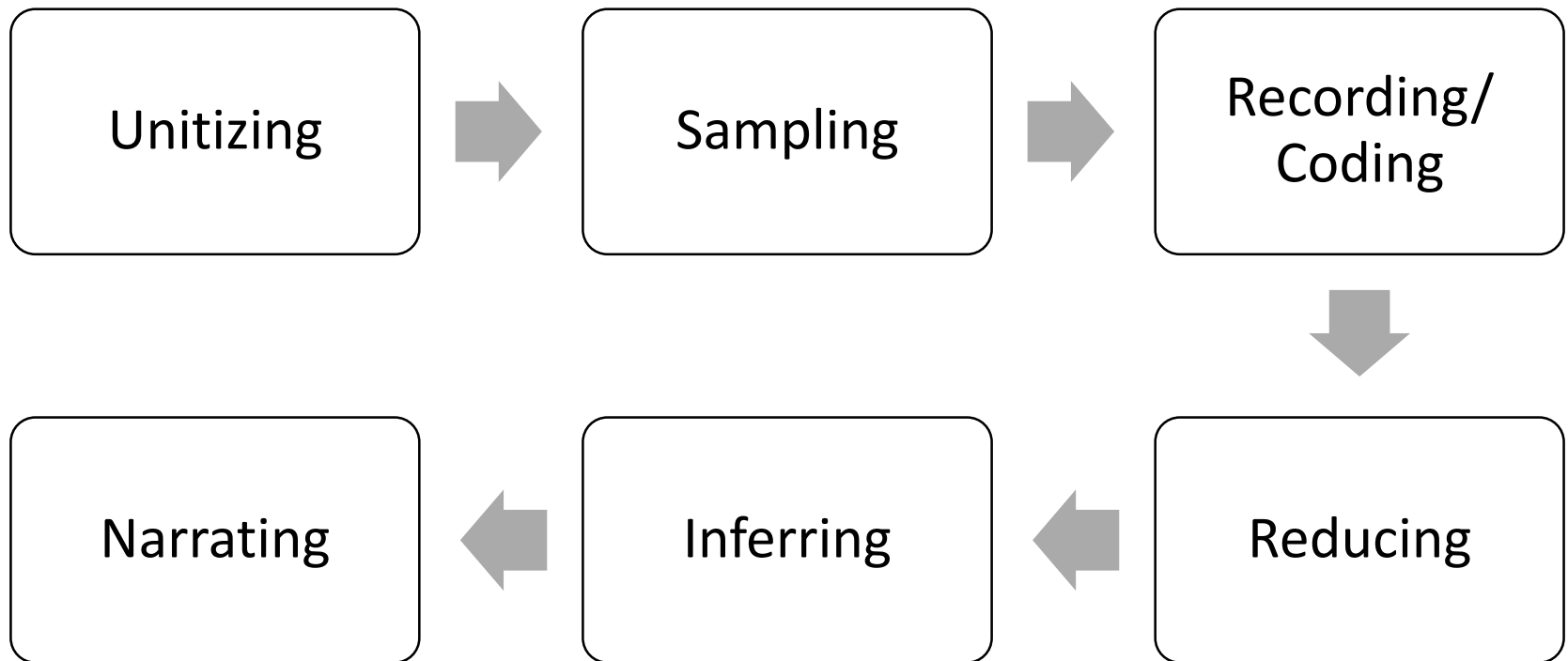prime minister

toilet paper

# Zipf law

Zipf's law

Legend:
- Esperanto
- Latin
- Ukrainian
- Czech
- Italian
- Spanish
- Slovene
- Finnish
- Hebrew
- Turkish
- Hungarian
- Galician
- Danish
- Belarusian
- Portuguese
- German
- Malay
- English
- Slovak
- Romanian
- Polish
- Uzbek
- French
- Basque
- Serbian
- Dutch
- Catalan
- Indonesian
- Lithuanian
- Croatian

x-axis: log(rank)
y-axis: log(frecuency)

(Jimenez, 2015)

# Manifest vs. latent content

# Design of CA research

Texts

Research Question

Answers

Content Analysis

Krippendorff 2013, p. 36

# Design of CA research
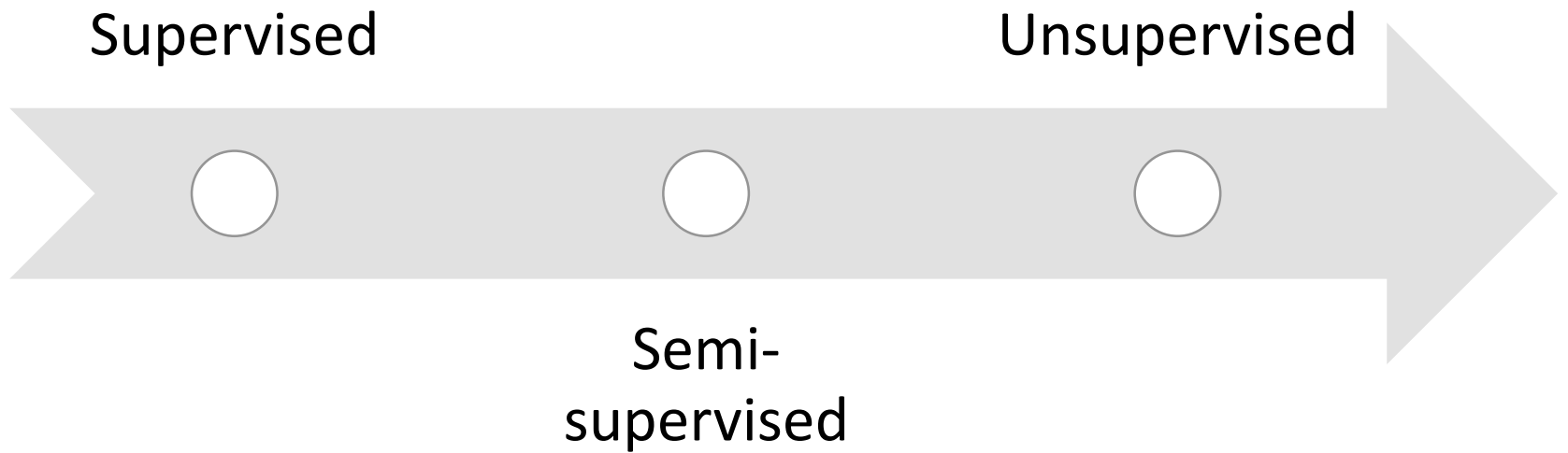


Krippendorff 2013, p. 86

# Basic terminology

- Corpus
  - Body of all text pieces available for the content analysis
- Term
  - Text token, usually word
- Term-document matrix
  - Matrix which records occurrence of terms in documents

# Methods

# Methods of TA

Supervised

Unsupervised

Semi-
supervised

# Methods of TA

- Supervised methods
  - Manual coding
- Semi-supervised
  - Dictionary-based methods
    - Deductively given dictionary
    - Dictionary obtained from data
      - Automatically
      - Manually
- Unsupervised
  - Frequencies
  - Topic modeling
  - …

# Fully supervised – manual coding

- Manual coding of text units
- Inductive vs. deductive coding
  - Inductive – data-driven
    - Categories not known
    - Open coding – categories emerge in iterative text reading
    - Axial coding – abstraction from open coding into categories
  - Deductive – theory-driven
    - Categories known a-priori
    - Existing code-book applied over data

# Fully supervised

- Coding is input for further analysis
  - Frequencies of codes
  - Temporal development
  - Standard statistical methods
  - Socio-semantic networks
- Discourse network analysis
  - Socio-semantic networks of actors and meanings (codes) they use

# Fully supervised - DNA



(Haunss et al. 2013)

# Fully supervised - DNA

# Issues with manual coding

- Questions of validity and reliability
- Reliability of human coders needs to be measured and accounted for
  - Intra-coder reliability (variation by same coder)
  - Inter-coded reliability (variation by different coder)
- Ways how to measure, e.g. Krippendorff $\alpha$
- Ways how to account for
  - Only overlap
  - Resolution of differences

# Semi-supervised

- Dictionary-based automated coding
  - Words in dictionary are discovered across the corpus
  - Coding process is done automatically
- Construction of dictionaries
  - Given pre-defined dictionary
    - WordStat, LIWC, …
  - Constructed from data
    - Theoretically-informed
    - Automatically generated
      - WordFish
      - WordScores

# Semi-supervised

- Existing dictionaries
  - WordStat (Laver & Garry 2000)
    - Estimation of policy positions from political texts
    - 415 words, 19 categories
  - LIWC (Linguistic Inquiry Word Count)
  - Sentiment dictionaries
  - General Inquirer
- Logic of this approach is to crawl over texts, discover tokens in dictionary and score texts
  - Scoring whole corpus
  - Scoring individual texts

# Semi-supervised

- Dictionary from data
  - Sample of texts with known properties
  - Other texts related – e.g. legal/conceptual documents

- Dictionaries built by researchers
  - Long process
  - High validity – researchers know texts
  - Lower reliability – same reasons as manual coding

# Semi-supervised

- Automated dictionary constructing
  - Laver, Benoit and Garry 2003
  - Two populations of texts
    - Texts with known properties – training set
    - Texts with unknown properties – target set
  - Logic of the process
    - Assign values of the category to known texts (training sample)
    - Let computer find words in the training sample and assign individual scores to words from texts
    - Code unknown texts with existing dictionary
  - High reliability, but questionable validity

# Unsupervised

- Most naïve – word frequencies
  - Just a crude exploratory hint of what is in text
- Clustering and multidimensional scaling of words
  - Based on co-occurrence of words
- Unsupervised categorization on term-document matrix
  - Topic modeling
- Co-occurrence term networks

prostředek ukázat včetně parlament vést potřeba
italie systém konec schengen festival existovat
pomáhat martin blízký vlna národní fungovat důležitý solidarita
jediný viktor bezpečnost práce příští úřad orbán takzvaný mimo dítě
znamenat mezinárodní souvislost hovořit
dohodnout diskuse kdyby rozhodnout předevšim často
fórum debata unijní něco uprchlický vztah myslit volba řada
významný nejen zákon konference přijmout obrana případ poslední zástupce poslanec
platit brzy občan trh některý zahraniční prostor návštěva
czech cesta rusko zahraniční vysoký rozdíl
brát čas vůči autor dohoda kvůli zeman právo sám jednat postoj téma shodnout
kontrola zaorálek pomoc skupina visegrádský patřit
vojenský visegrád republika summit maďarsko maďarský současný postup
hlava čekat procento čssd spolječný kvóta hranice návrh politik ministerstvo
divadlo spíše totiž nějaký řešení vědět sobotka politika ekonomický přístup
klíčový milión řešit ruský tlak jít uprchlík zahraničí cíl členský výsledek hl
ovšem chovanec miliarda čtyřka fond lze německo vláda evropa krize energetický bývalý
život peníze oba člověk usa stát premiér zájem doba přitom útok
prý důvod slovenský pan tam schengenský
pozice tady merkelová migrační polsko věc takový možnost názor říkat západní
avíc kolega spolupráce začít síla moc země chtít schůzka informace
kancléřka ukrajinský včera česko ukrajina například hlavní region
zahrada služba ctk program praha plyn evropský šéf otázka podpora teď udělat andre
putin americký jednání slovensko dodat velmi student jednotlivý
anizace stejně euro vidět ministr říci velký uvést zatím dojít
slovák zároveň svět oblast datetimestamp unie dobře nato povinný jasný cena
angela východní problém běženec zdůraznit německý akce řecko opatření
daleko snažit krok projekt tisíc dát prezident politický vnitro možný malý západ
policista odmítat měsíc tiskový nyní rámec polský migrant společnost většina afrika
podpořit část předseda bohuslav komise celý mluvit prohlásit turecko
británie miloš druhý kdy člen dobrý brusel změna základ
rakouský jan několik žádný ochrana jaký rád rakousko podporovat naopak francie
veřejný minulý zejména vnější válka
demokracie kyjev pozice odmítnout sociální zdroj nikdo centrum
mluvčí letošní boj řada právě týden dále robert

ewa kopaczová ochrana vnější předseda vláda český premiér situace kdy francie němec
ně západní david cameron tisíc běženec petr
polský vláda ministr finance viktor orbán migrační krize schengenský hranice příští týden z
vladimir putin stát eu nato eu poslední rok jean claude česko maďarsko
hlavní cíl islámský stát povinný kvóta ostatní země evropský fo
andrej babiš vláda čr sobotka čssd celý evropa ministr vnitro situace ukrajina minulý rok
ně střední ostatní štát čtk ap ministr zahraničí sociální demokr
řešení migrační polsko maďarsko slovensko polsko eu nato rad
čr polsko čtyřka v evropský komise prezident miloš český
ej duda východní evropa země eu český republika maďarský premiér brát vážn
společný postup ministr obrana summit eu politick
edsnictví čech slovak robert fico čr sr právo – doba kdy žadatel azyl průmysl ob
te některý země skupina v evropský unie český vláda tomáš prouza eu země před
a včera země unie milión euro uprchl
přerozdělování uprchlík ministerstvo zahraničí hran
tó mimořádný summit visegrádský čtyřka boj proti andrej kiska v
e jean milan chovanec velký británie pa
nacie lubomír stát v bohuslav sobotka lubomír zaorálek itálie řecko česko polsk
ka ministr vnitro milan slovenský premiér společný ochran
ský stát západní evropa země visegrádský říci in členský stát premiérka ewa
avedlnost – čtk sociální demokrat
ečný země britský premiér ochrana hranice evropský země návrh
miliarda koruna v – rámec
álení říci včera vnější hranice řešit příčina prezid
premiér sobotka příští rok premiér bohuslav země v několik rok
šet ruský plyn člen eu řecko itálie miliarda euro celý svět posle
er v ruský prezident země visegrádský skupina uprchlický krize šéf český společný evro
segrádský země česko slovensko letošní rok pracovní míst
kvóta odmítat obrana martin členský země slovensko maďarsko stát unie
ký kvóta schengenský prostor konec rok miloš zeman tisíc uprchlík václav klaus prezident vl
prchlík přijímání uprchlík uvnitř eu miliarda kč rok kdy halo noviny šéf evropský uvést pr
tický útok angela merkelová kancléřka angela velký množství v včera
ý prohlášení maďarsko slovensko střední evropa maďarsko polsko premiér země mutace mladý
áce rámec rámec eu donald tusk premiér viktor kvóta uprchlík azylový politik
kvóta přerozdělování zahraničí lubomír hlava stát republika země
at ministr miroslav lajčák evropský parlament migrační vlna příliv uprchlík některý stát
zahraniční věc mutace metro český ministr český diplomacie maďarský hranice
ánit vnější země chtít martin stropnický minulý týden spojený stát sr polsko
životní prostředí devadesátý rok ukrajinský krize proti kvóta předseda evropský zeman prez
německo francie rada eu energetický bezpečnost bílý dům polsko český
polský premiérka syrský uprchlík evropský rada lidský právo loňský rok plynovod nord parlamentní volba země původu
ministr zahraniční světový válka slovenský ministr polský prezident proti islámský zahraniční v
metro praha český prezident stát visegrádský vojtěch filip šéf maďarský
společný stanovisko jednotlivý stát německo rakousko plyn ukrajina
ministerstvo vnitro evropský záležitost

# Unsupervised – co-occurrence net

# Unsupervised – topic modeling

# How to get TDM?

# Data pre-processing

- Any text analysis must be preceded by data pre-processing
  - Dropping sparse terms – has a word which occurs in 1.5M corpus once, any value?
  - Dropping most frequent terms – does most profound word of interest any informative value?
  - Dropping "stopwords" – a, the, …
  - Dropping numerals, punctuation, …
  - Dropping time and place information
  - …
- No general rules on how to do that – rule of thumb

# Data pre-processing

- Stemming/lemmatization
  - Disposal of grammatical features of text
    - Dictionary-based
    - Rules-based
  - Both introduce some error into the corpus
- Lemmatization
  - Identification of lemmas (lexemes) of the words - transformation to lemmas
- Stemming
  - Stripping the word of prefixes or suffixes, leaving only word stems

# Lemmatization and stemming

"This was the most tranquil presidential address. President's approach was very relaxed."

- Lemmatization

"This be the most tranquil presidential address. President approach be very relax."

- Stemming

"This be the most tranquil presidenti address.

Presid approach be veri relax."

# Corpus generation

- Decision on document unitizing

- Decision over sampling
  - Does 5M texts provide more information than 15k?
  - Random vs. non-random sampling

- Inclusion of metadata – allow for filtering later
  - Author
  - Time and date
  - Source (e.g. media/newspaper)
  - …

# Term-document matrix

- Matrix – most methods based on this
  - $1^{st}$ dim – Tokens
  - $2^{nd}$ dim – Documents/units
  - Cells – frequency of tokens in documents
    - Boolean – Present vs. Not present (1/0)
    - Weighted
      - Absolute frequency (how many times word occur in document)
      - TF-IDF
- Grows large easily
  - 500 documents * easily 4k unique tokens = 2M cells
- At the same time, very sparse
  - Most of cells are empty – contain 0

# Term-document matrix

|  | 2003-2004-cz | 2004-2005-pl | 2005-2006-hu | 2006-2007-sk | 2007-2008-cz | Sum |
|---|---|---|---|---|---|---|
| **agriculture** | 3 | 6 | 2 | 5 | 3 | **19** |
| **aim** | 4 | 2 | 7 | 12 | 6 | **31** |
| **area** | 11 | 8 | 8 | 28 | 26 | **81** |
| **base** | 1 | 2 | 2 | 2 | 5 | **12** |
| **border** | 5 | 9 | 9 | 3 | 3 | **29** |
| **central** | 2 | 3 | 6 | 3 | 5 | **19** |
| **cohesion** | 3 | 1 | 7 | 4 | 4 | **19** |
| **commission** | 2 | 7 | 3 | 2 | 4 | **18** |
| **common** | 10 | 9 | 17 | 8 | 17 | **61** |
| **community** | 2 | 2 | 3 | 3 | 6 | **16** |
| **concern** | 9 | 13 | 12 | 18 | 6 | **58** |

# Programming in R

# Learning curves of popular stats programs

# R community / resources

- there is huge number of free resources
- R package / library manuals
- R site: http://cran.r-project.org
- community forums:
  - http://stackoverflow.com
  - http://www.statmethods.net
  - http://www.r-bloggers.com
- Youtube videos:
  https://www.youtube.com/watch?v=qHfSTRNg6jE
- googling (often fastest)

# R as language – focus on logic

- Any programming language is just very **condensed and formalized** speech
  - Just like mathematical notation
- Understand and formulate the **process**
- If you think about the **procedure** of what needs to be done, scripting becomes matter of knowing right expressions

# R studio layout

Scripting window

Environment (stored objects)
History

Console window

Plots
Packages
Help
Viewer

Untitled1

Source on Save | Run | Source

1

(Top Level)                                                    R Script

Environment | History

Global Environment

Environment is empty

Files Plots Packages Help Viewer

Zoom | Export

Console ~/

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

>

RStudio

File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Go to file/function          Addins ▾                                                    Project: (None) ▾

Untitled1 ✕                                                                    □  □

     Source on Save  🔍  ✏ ▾  ▤         ▾                    → Run   ⇥ ▾   → Source  ▾  ▤

1

**Scripting window**

1:1   (Top Level) ⬍                                                              R Script ⬍

Environment   History                                                          □  □

📂 ▾  💾   Import Dataset ▾  🧹                                       ☰ List ▾  ⟳

🗂 Global Environment ▾                                          🔍

Environment is empty

**Environment**
**History**

Files  Plots  Packages  Help  Viewer                                           □  □

◀  ▶   🔍 Zoom   📤 Export ▾  ⊗  🧹                                            ⟳

**Plots**
**Packages**
**Help**
**Viewer**

Console  ~/  ⬀                                                                 □  □

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> |

**Console**

# Terminology used

- Data, (data) element
  - **unit of information** (e.g. 1, 2, "word", TRUE, FALSE)
- Data class
  - describes **properties of data elements** (numeric, character, logical, etc.)
- Object
  - a **"container"** that stores and organizes data in the R environment
- Object type
  - describes **properties of objects** (vector, matrix, list, data frame, etc.)
- Function
  - transforms inputs into outputs based on certain rules (methods/procedures)
  - arguments of the function specify the inputs and applied rules

# Object

- object: instance of a certain data class that can be manipulated according set of procedures (methods)

```
one <- 1
```

# Object

File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Project: (None)

Untitled1* ×

Source on Save

Run | Source

```
1  one <- 1
2
```

1:9  (Top Level)  R Script

Console  ~/

>

Environment | History

Import Dataset | List

Global Environment

Environment is empty

Files | Plots | Packages | Help | Viewer

Install | Update

| Name | Description | Version |
|---|---|---|
| **User Library** | | |
| assertthat | Easy Pre and Post Assertions | 0.2.0 |
| audio | Audio Interface for R | 0.1-5 |
| beepr | Easily Play Notification Sounds on any Platform | 1.2 |
| BH | Boost C++ Header Files | 1.62.0-1 |
| bindr | Parametrized Active Bindings | 0.1 |
| bindrcpp | An 'Rcpp' Interface to Active Bindings | 0.2 |
| bitops | Bitwise Operations | 1.0-6 |
| Cairo | R graphics device using cairo graphics library for creating high-quality bitmap (PNG, JPEG, TIFF), vector (PDF, SVG, PostScript) and display (X11 and Win32) output | 1.5-9 |
| chron | Chronological Objects which can Handle Dates and Times | 2.3-50 |
| colorspace | Color Space Manipulation | 1.3-2 |
| curl | A Modern and Flexible Web Client for R | 2.8.1 |
| data.table | Extension of `data.frame` | 1.10.4 |
| dichromat | Color Schemes for Dichromats | 2.0-0 |

# Creating/storing objects



Lukas <-

# Creating/storing objects

Obj. name **<-** Object

# Object

- Once objects exist, operations over objects may be applied

```
one <- 1

one + one


> one <- 1
> one + one
[1] 2
```

RStudio

File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Go to file/function          Addins ▼                                          R Project: (None) ▼

Untitled1* ×

Source on Save                                                    Run        Source ▼

```
1  one <- 1
2
3  one + one
4
```

4:1    (Top Level) ≑                                                           R Script ≑

Console ~/

```
> one <- 1
> one + one
[1] 2
>
```

Environment   History

Import Dataset ▼                                                List ▼

Global Environment ▼

**Values**

one                    1

Files  Plots  Packages  Help  Viewer

Install    Update

| Name | Description | Version |
|---|---|---|
| **User Library** | | |
| assertthat | Easy Pre and Post Assertions | 0.2.0 |
| audio | Audio Interface for R | 0.1-5 |
| beepr | Easily Play Notification Sounds on any Platform | 1.2 |
| BH | Boost C++ Header Files | 1.62.0-1 |
| bindr | Parametrized Active Bindings | 0.1 |
| bindrcpp | An 'Rcpp' Interface to Active Bindings | 0.2 |
| bitops | Bitwise Operations | 1.0-6 |
| Cairo | R graphics device using cairo graphics library for creating high-quality bitmap (PNG, JPEG, TIFF), vector (PDF, SVG, PostScript) and display (X11 and Win32) output | 1.5-9 |
| chron | Chronological Objects which can Handle Dates and Times | 2.3-50 |
| colorspace | Color Space Manipulation | 1.3-2 |
| curl | A Modern and Flexible Web Client for R | 2.8.1 |
| data.table | Extension of `data.frame` | 1.10.4 |
| dichromat | Color Schemes for Dichromats | 2.0-0 |

# What is an object?

- **Anything** may become an object

- Temporary objects
  - Only **appear** in console
  - Their values must be stored in order to use them in operations

- Stored objects
  - Must be **defined** by user
  - **Remain the same** unless overwritten
  - Must be **removed** by user as well

RStudio

File   Edit   Code   View   Plots   Session   Build   Debug   Profile   Tools   Help

Go to file/function          Addins ▾                                    Project: (None) ▾

Untitled1* ×

Source on Save                                                    Run    Source ▾

```
1  one <- 1
2
3  one + one
4
5  two <- one + one
6  |
```

6:1   (Top Level) ⬍                                                    R Script ⬍

Console  ~/

```
> one <- 1
> one + one
[1] 2
> two <- one + one
>
```

Environment   History

Import Dataset ▾                                              List ▾

Global Environment ▾

**Values**

one          1
two          2

Files   Plots   Packages   Help   Viewer

Install    Update

| Name | Description | Version |
|------|-------------|---------|
| **User Library** | | |
| assertthat | Easy Pre and Post Assertions | 0.2.0 |
| audio | Audio Interface for R | 0.1-5 |
| beepr | Easily Play Notification Sounds on any Platform | 1.2 |
| BH | Boost C++ Header Files | 1.62.0-1 |
| bindr | Parametrized Active Bindings | 0.1 |
| bindrcpp | An 'Rcpp' Interface to Active Bindings | 0.2 |
| bitops | Bitwise Operations | 1.0-6 |
| Cairo | R graphics device using cairo graphics library for creating high-quality bitmap (PNG, JPEG, TIFF), vector (PDF, SVG, PostScript) and display (X11 and Win32) output | 1.5-9 |
| chron | Chronological Objects which can Handle Dates and Times | 2.3-50 |
| colorspace | Color Space Manipulation | 1.3-2 |
| curl | A Modern and Flexible Web Client for R | 2.8.1 |
| data.table | Extension of `data.frame` | 1.10.4 |
| dichromat | Color Schemes for Dichromats | 2.0-0 |

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function | Addins | Project: (None)

Untitled1* |

Source on Save | Run | Source |

```
1  one <- 1
2
3  one + one
4
5  two <- one + one
6
7  two
8
```

8:1 (Top Level) | R Script

Environment | History

Import Dataset | List

Global Environment

**Values**

| one | 1 |
| two | 2 |

Files | Plots | Packages | Help | Viewer

Install | Update

| Name | Description | Version |
|---|---|---|
| **User Library** | | |
| assertthat | Easy Pre and Post Assertions | 0.2.0 |
| audio | Audio Interface for R | 0.1-5 |
| beepr | Easily Play Notification Sounds on any Platform | 1.2 |
| BH | Boost C++ Header Files | 1.62.0-1 |
| bindr | Parametrized Active Bindings | 0.1 |
| bindrcpp | An 'Rcpp' Interface to Active Bindings | 0.2 |
| bitops | Bitwise Operations | 1.0-6 |
| Cairo | R graphics device using cairo graphics library for creating high-quality bitmap (PNG, JPEG, TIFF), vector (PDF, SVG, PostScript) and display (X11 and Win32) output | 1.5-9 |
| chron | Chronological Objects which can Handle Dates and Times | 2.3-50 |
| colorspace | Color Space Manipulation | 1.3-2 |
| curl | A Modern and Flexible Web Client for R | 2.8.1 |
| data.table | Extension of `data.frame` | 1.10.4 |
| dichromat | Color Schemes for Dichromats | 2.0-0 |

Console ~/

```
> one <- 1
> one + one
[1] 2
> two <- one + one
> two
[1] 2
>
```

# Data classes – prop. of elements

- Numeric
  - continuous numeric data
  - -1, 0.5, 10.49
- Integer
  - discrete numeric data
  - -1, 0, 1
- Character
  - string values
  - "anythingWithinQuotes"
- Logical
  - output of logical operation – TRUE/FALSE
  - 5 > 10

# Data classes

```
> as.numeric(10.49)
[1] 10.49
>
> as.integer(10.49)
[1] 10
>
> as.character(-1)
[1] "-1"
>
> as.numeric("anythingwithinquotes")
[1] NA Warning message: NAs introduced by coercion
>
> 5 > 10
[1] FALSE
>
> as.character(5 > 10)
[1] "FALSE"
```

# Object types – prop. of objects

# Object types – prop. of objects

- vector
  - sequence (1-dimensional) of elements of same data class

- matrix
  - 2-dimensional rectangular collection of elements of same data class
  - array: n-dimensional matrix

- list
  - vector that can contain elements of different data classes

- data frame
  - list of vectors of equal length
  - table data

# Vector

```
> c(2,3,5)
[1] 2 3 5
>
> c("aa", "bb", "cc", "dd", "ee")
[1] "aa" "bb" "cc" "dd" "ee"
>
> c(TRUE, FALSE, TRUE, FALSE, FALSE)
[1] TRUE FALSE TRUE FALSE FALSE
>
```

# Matrix

```
> m <- matrix(data = c(1,2,3,4,5,6,7,8,9,10,11,12),
+                nrow = 3,
+                ncol = 4)
>
> m
      [,1]   [,2]   [,3]   [,4]
[1,]  1      4      7      10
[2,]  2      5      8      11
[3,]  3      6      9      12
>
```

# List

```
> n <- c(2, 3, 5)
> s <- c("aa", "bb", "cc", "dd", "ee")
> x <- list(n, s, b, 3) # x contains copy of n, s
> x
[[1]]
[1] 2 3 5

[[2]]
[1] "aa" "bb" "cc" "dd" "ee"

[[3]]
[1] TRUE FALSE TRUE FALSE FALSE

[[4]]
[1] 3
```

# Data frame

```
> teams <- c("PHI","NYM","FLA","ATL","WSN")
> wins <- c(92,89,94,72,59)
> losses <- c(70,73,77,90,102)
>
> data <- data.frame(teams,wins,losses)
>
> data
      teams  wins   losses
1     PHI    92     70
2     NYM    89     73
3     FLA    94     77
4     ATL    72     90
5     WSN    59     102
>
```

# R functions

- `word()` indicates function

```
> sqrt(9)
[1] 3
```

- `function(argument_1, argument_2, …)`

```
> sample(x = 0:100, size = 10, rep = FALSE)
[1] 48 50 37 94 42 39 21 19 63 95
```

- basic functions (part of the basic R package)
- package functions (part of the particular package)
- user functions (user-defined functions)

# R libraries

- Libraries allow to load pre-defined functions according to problem at hand

- Load, install and unload either using R Studio or using functions in script

- Libraries download and install automatically

File Edit Code View Plots Session Build Debug Profile Tools Help

Project: (None)

Untitled1*

Source on Save | Run | Source

```
1  library(beepr)
2
```

2:1 (Top Level)

Environment | History

Import Dataset | List

Global Environment

Environment is empty

Console ~/

```
> library(beepr)
>
```

**Install Packages**

Install from: ? Configuring Repositories

Repository (CRAN, CRANextra)

Packages (separate multiple with space or comma):

Install to Library:

C:/Users/Lukas/Documents/R/win-library/3.4 [Default]

☑ Install dependencies

Install | Cancel

Plots | Packages | Help | Viewer

all | Update

| Name | Description | Version |
|------|-------------|---------|
| **brary** | | |
| sertthat | Easy Pre and Post Assertions | 0.2.0 |
| udio | Audio Interface for R | 0.1-5 |
| eepr | Easily Play Notification Sounds on any Platform | 1.2 |
| H | Boost C++ Header Files | 1.62.0-1 |
| ndr | Parametrized Active Bindings | 0.1 |
| bindrcpp | An 'Rcpp' Interface to Active Bindings | 0.2 |
| bitops | Bitwise Operations | 1.0-6 |
| Cairo | R graphics device using cairo graphics library for creating high-quality bitmap (PNG, JPEG, TIFF), vector (PDF, SVG, PostScript) and display (X11 and Win32) output | 1.5-9 |
| chron | Chronological Objects which can Handle Dates and Times | 2.3-50 |
| colorspace | Color Space Manipulation | 1.3-2 |
| curl | A Modern and Flexible Web Client for R | 2.8.1 |
| data.table | Extension of `data.frame` | 1.10.4 |
| dichromat | Color Schemes for Dichromats | 2.0-0 |

RStudio

File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Go to file/function          Addins ▾

Project: (None) ▾

Untitled1* ✕

☐ Source on Save        Run    ↦ Source ▾

```r
1  library(beepr)
2
```

2:1  (Top Level)  R Script

Environment    History

Import Dataset ▾        List ▾

Global Environment ▾

Environment is empty

Console  ~/

```
> library("network", lib.loc="~/R/win-library/3.4")
network: Classes for Relational Data
Version 1.13.0 created on 2015-08-31.
copyright (c) 2005, Carter T. Butts, University of California-Irvine
                    Mark S. Handcock, University of California -- Los An
geles
                    David R. Hunter, Penn State University
                    Martina Morris, University of Washington
                    Skye Bender-deMoll, University of Washington
 For citation information, type citation("network").
 Type help("network-package") to get started.

> detach("package:network", unload=TRUE)
>
```

Files    Plots    Packages    Help    Viewer

☐ Install    ↻ Update

| | Name | Description | Version |
|---|---|---|---|
| ☐ | ldatuning | Tuning of the Latent Dirichlet Allocation Models Parameters | 0.2.0 |
| ☐ | magrittr | A Forward-Pipe Operator for R | 1.5 |
| ☐ | maptools | Tools for Reading and Handling Spatial Objects | 0.9-2 |
| ☐ | mime | Map Filenames to MIME Types | 0.5 |
| ☐ | modeltools | Tools and Classes for Statistical Models | 0.2-21 |
| ☐ | munsell | Utilities for Using Munsell Colours | 0.4.3 |
| ☐ | network | Classes for Relational Data | 1.13.0 |
| ☐ | NLP | Natural Language Processing Infrastructure | 0.1-10 |
| ☐ | openNLP | Apache OpenNLP Tools Interface | 0.2-6 |
| ☐ | openNLPdata | Apache OpenNLP Jars and Basic English Language Models | 1.5.3-2 |
| ☐ | openssl | Toolkit for Encryption, Signatures and Certificates Based on OpenSSL | 0.9.6 |
| ☐ | PCIT | Partial Correlation Coefficient with Information Theory | 1.5-3 |
| ☐ | pkgconfig | Private Configuration for 'R' Packages | 2.0.1 |
| ☐ | plogr | The 'plog' C++ Logging Library | 0.1-1 |
| ☐ | plotrix | Various Plotting Functions | 3.6-5 |

# Basic R functions

```
c() # combine two or more elements into an object

class() # explore elements' data class
length() # explore number of first dim. of object
dim() # explore dimensions of two-dimensional obj.
nrow() # number of rows
ncol() # number of columns

head() # first few rows of data
tail() # last few rows of data
str() # explore structure of object

names() # names in the named vector - one dimension
rownames() # names of rows - two dimensions
colnames() # names of columns - two dimensions
```

# Working directory

- Folder, where all imports and exports are taking place – enough to set once

- Makes data import and export easier

- Functions `setwd()` and `getwd()`

- Does not accept single backslash in Win path
  - Replace backslash \ with forwardslash / or double backslash \\

```
setwd("C:\\Users\\Lukas\\Documents\\R intro")

setwd("C:/Users/Lukas/Documents/R intro")
```

# Data output

- Save entire workspace
  - Save all R objects you've created so far
  - Allows to return to work/backup current work
- Save particular object
  - Export data to tabular objects
  - CSV as most common format

# CSV - most common data format

- **C**omma-**S**eparated **V**alues
- Tabular data separated by commas (separator/delimiter) or other signs (tabulator, space, semicolon)
- CSV file (.csv), TSV file (.tsv) – always a **text file** (.txt)
- Must have **same number of columns** (separators)

```
cars,type,price,consumption,emissions,expensive
BMW,3,1200000,6.2,0,0
Audi,A4,1164000,5.9,0,0
VW,Passat,950500,6.2,NA,NA
```

# CSV – other examples

```
cars;type;price;consumption;emissions
BMW;3;1200000;6.2;0
Audi;A4;1164000;5.9;0
VW;Passat;950500;6.2;0
```

```
"cars" "type" "price" "consumption"  "emissions"
"BMW" "3" "1,200,000" "6.2" "0"
"Audi" "A4" "1,164,000" "5.9" "0"
"VW" "Passat" "950,500" "6.2" "0"
```

```
cars,type,price,consumption,emissions
BMW,3,1,200,000,6.2
Audi,A4,1,164,000,5.9
VW,Passat,950,500,6.2
```

Bad data – improper use of comma delimiter results in uneven # of rows

# Exporting object – tabular

- Function `write.table()`
- Name of file must be specified
- Easy to import to Excel or other software

```r
frequencies <- c(92,89,94,72,59)

write.table(frequencies,
            "frequencies.csv",
            sep = ",",
            row.names = FALSE,
            col.names = TRUE,
            fileEncoding = "UTF-8")
```

# Exporting object – unstructured

- Function `writeLines()`

- Has basically no arguments

- Saves the whole object as one text

```
frequencies <- c(92,89,94,72,59)

writeLines(frequencies,
           "frequencies.txt")
```

# Text analysis in R

# Text analysis in R

- Most prominent package for text analysis is "tm" (stands for text mining)
  - Provides tools corpus creation, text manipulation, term-document matrix creation
  - Easily allows to read text documents as corpus
- Competing packages – "quanteda"
  - Developed by Ken Benoit (WordScores)
  - Provides some TA methods
  - Overlaps with "tm" package – if both packages loaded, it will generate conflicts (feature, not bug)

# Corpus

- `getSources()` provides list of available sources
  - Files inside a directory – `DirSource()`
  - Text inside a vector – `VectorSource()`
  - Dataframe, XML, links to web-sites, …
- `Corpus()` creates a corpus object out of text sources

# Corpus

```
my.texts <-  "C:\\Users\\Lukas\\Desktop\\data\\"

directory.source <- DirSource(directory = my.texts)

text.corpus <- Corpus(directory.source)
```

# Corpus operations – functions

- Useful functions:
  - `removePunctuation()` – remove all punctuation
  - `removeWords()` – remove stopwords
  - `stripWhitespace()` – remove duplicate white space
  - `removeNumbers()` – remove all numbers
  - `stemDocument()` – stem document
  - `plainTextDocument()` – turn document into tm package's plain text format

# Corpus operations

- `tm_map()` function allows to apply manipulations over the corpus data

```
edited.corpus <- text.corpus

edited.corpus <- tm_map(edited.corpus, removeNumbers)

edited.corpus <- tm_map(edited.corpus, removePunctuation)

edited.corpus <- tm_map(edited.corpus, stripWhitespace)

edited.corpus <- tm_map(edited.corpus,
                    removeWords,
                    stopwords("english"))
```

# Term-document matrix

- **Function** `TermDocumentMatrix()`
  - Terms in rows
  - Documents in columns
- `DocumentTermMatrix()` creates inverse TDM
- Output is non-standard matrix object
  - If matrix operations are needed, it must be converted to basic matrix format with `as.matrix()` function

# Term-document matrix

```
tdm <- TermDocumentMatrix(edited.corpus)

dtm <- DocumentTermMatrix(edited.corpus)

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

# Useful functions in "tm"

- `removeSparseTerms()`
  - Removes terms to a defined sparsity of the TDM matrix – removes terms which are used sparsely across documents
- `findFreqTerms()`
  - Lists most frequent terms across the TDM matrix
  - Does **not provide frequencies**, though
- `findAssocs()`
  - Correlation of appearance of a term with other terms across TDM – returns **Pearson's r**

# Frequencies

- `findFreqTerms()` shows frequent terms
  - Has two attributes defining bounds – `lowfreq`, `highfreq`
- Easier to calculate frequencies separately
  - Convert TDM to matrix with `as.matrix()`
  - Calculate sums of rows with `rowSums()`
  - Sort the vector with `sort()` with `decreasing` attribute

```r
tdm.matrixed <- as.matrix(tdm)

frequencies <- rowSums(tdm.matrixed)

frequencies <- sort(frequencies, decreasing = T)
```

# Wordclouds

- Package "wordcloud"
- Function `wordcloud()`

| Attribute | Description |
|---|---|
| `words` | Terms |
| `freq` | Frequencies of terms |
| `scale` | Two values in `c()` function to bound upper and lower scale |
| `max.words` | Maximum number of words rendered |
| `random.order` | Binary - should terms be placed in random order? |
| `rot.per` | Percentage of terms placed vertically |
| `colors` | Color or color palette |
| `random.color` | Binary – should colors be assigned randomly or based on the word frequency? |

# Wordclouds

```r
tdm.matrixed <- as.matrix(tdm)

frequencies <- rowSums(tdm.matrixed)

frequencies <- sort(frequencies,decreasing = T)

terms <- names(frequencies)

library(wordcloud)

wordcloud(words = terms,
          freq = frequencies,
          scale = c(5,0.5),
          max.words = 150,
          random.order = F,
          rot.per = 0,
          colors = "red")
```

# Wordclouds