# Koronavírus

Granát Marcell és Mazzag Bálint

2021. március 1.
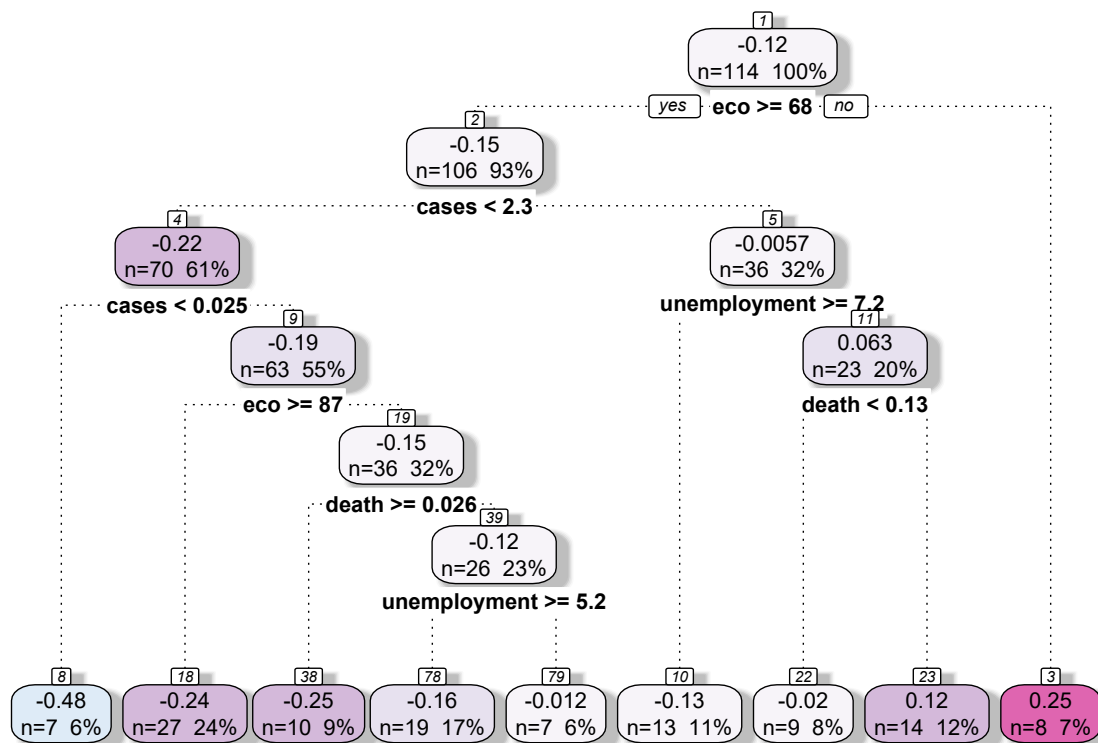
## Tartalomjegyzék
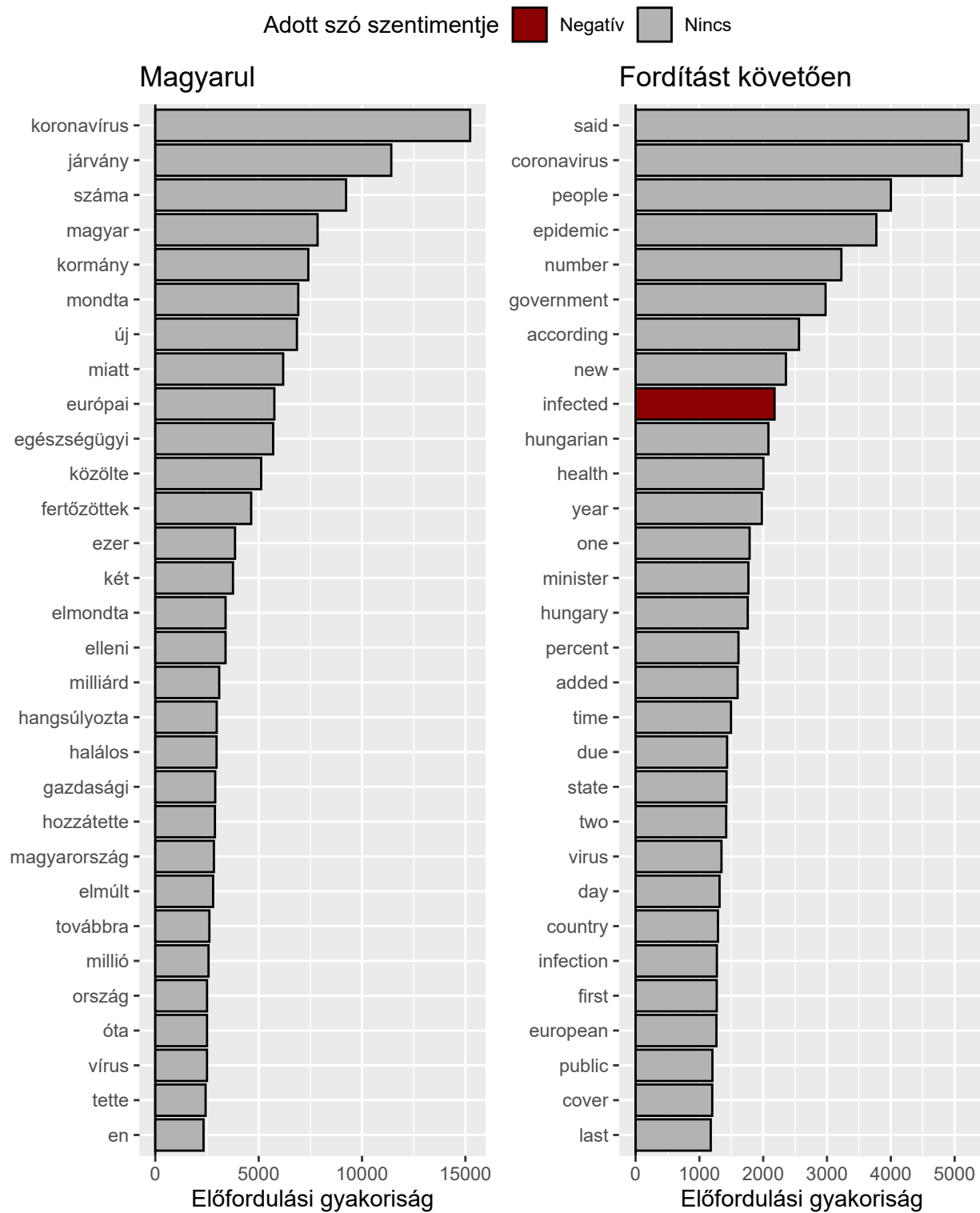
**Absztrakt**

Here is the abstract.
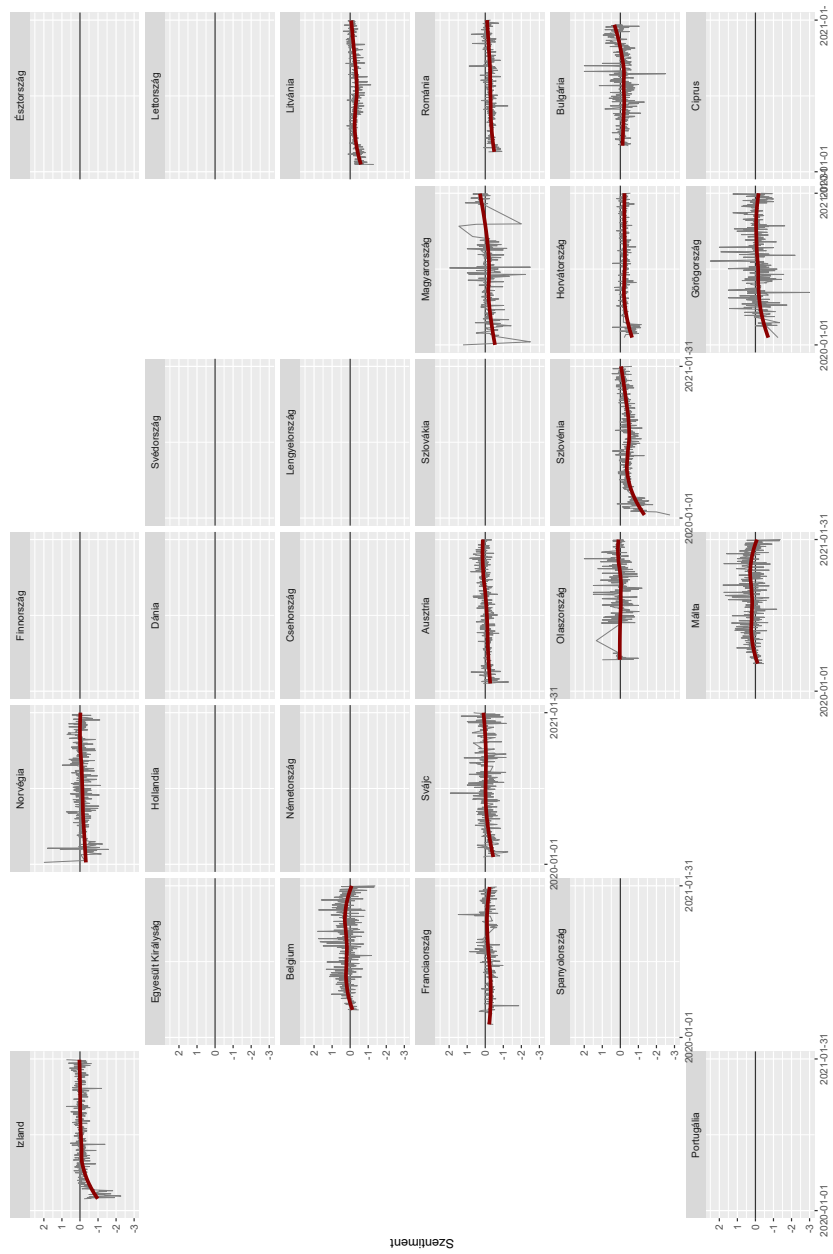
# Bevezetés

# Adatok

### Gépi fordítás

# Leíró statisztikák

1. ábra. Leggyakrabban előforduló szavak a magyar nylevű cikkekben a fordítást megelőzően és azt követően.

2. ábra. A szentiment alakulása országonként

3. ábra. Leggyakrabban előforduló pozitív és negatív szentimenttel rendelkező szavak

# Hivatkozások

Adsera, A. (2004), Changing fertility rates in developed countries. the impact of labor market institutions', *Journal of population economics* **17**(1), 17–43.

## Függelék: R kódok

```r
1   # Set up ---------------------------------------------------------------------------
2
3   ## Packages ========================================================================
4
5   library(tidyverse)
6   library(patchwork)
7   library(knitr)
8   library(broom)
9   library(geofacet)
10  library(tidytext)
11  library(tm)
12  library(wordcloud)
13
14  ## Gg theme ========================================================================
15
16  update_geom_defaults("point", list(fill = "cyan4",
17                                      shape = 21,
18                                      color = "black",
19                                      size = 1.4))
20  update_geom_defaults("line",
21                       list(color = "midnightblue", size = 1.4))
22
23  update_geom_defaults("smooth", list(color = "red4", size = 1.4))
24
25  update_geom_defaults("density",
26                       list(color = "midnightblue", fill =  "midnightblue",
27                            alpha = .3, size = 1.4))
28
29  extrafont::loadfonts(device="win")
30
31  theme_set(theme_grey() + theme(
32    legend.direction = "vertical",
33    plot.caption = element_text(family = "serif")
34  ))
35
36  # Data -----------------------------------------------------------------------------
37
38  # Articles =========================================================================
39
40  load("dat.RData")
41  # This RData contains the articles after the main cleaning process
42  # To ensure full reproducibility see the attached files at the corresponding
43  # GitHub Repo: -> https://github.com/MarcellGranat/CoronaSentiment <-
44
45  Hungary_rawtext <- readxl::read_excel("scrapping raw csv/Hungary_rawtext.xlsx") %>%
46    # Hungarian articles before translation
47    select(date, title, URL = links, text) %>%
48    mutate_all(function(x) str_remove_all(x, "\r")) %>%
49    mutate_all(function(x) str_remove_all(x, "\t")) %>%
50    mutate_all(function(x) str_remove_all(x, "\n")) %>%
51    mutate_at(-1, function(x) zoo::na.locf(x)) %>%
52    filter(!str_detect(date, '_x000') & date != '0') %>%
```

```r
53    filter(!str_detect(text, 'mtva_player')) %>% # TODO consider a better solution
54    mutate(
55      date = gsub(" -.*", "", date),
56      text = str_remove_all(text, "_x000D_"),
57      date = lubridate::ymd(date)
58    ) %>%
59    tidytext::unnest_tokens(words, text)
60
61  ### Add sentiment values to our data #################################################
62
63  dat_sentiment <- dat %>%
64    select(date, text, country) %>%
65    mutate(country = ifelse(str_detect(country, "BE"), "BE", country)) %>%
66    {left_join(tidytext::unnest_tokens(., words, text),
67               get_sentiments("afinn"), by=c("words"="word"))}
68  # TODO other packages
69
70  dat_sentiment_daily <- dat_sentiment %>%
71    group_by(date, country) %>%
72    summarise(value = mean(value, na.rm = T), n = n()) %>%
73    ungroup() %>%
74    na.omit() %>%
75    rename(code = country)
76
77  dat_sentiment_monthly <- dat_sentiment %>%
78    na.omit() %>%
79    mutate(
80      date = lubridate::ym(paste(lubridate::year(date), lubridate::month(date), sep = "-"))
81    ) %>%
82    group_by(date, country) %>%
83    summarise(value = mean(value, na.rm = T), n = n()) %>%
84    ungroup() %>%
85    na.omit() %>%
86    rename(code = country)
87
88  # COVID data ===================================================================
89
90  dat_covid <-
91    readr::read_csv("https://covid.ourworldindata.org/data/owid-covid-data.csv") %>%
92    transmute(code = countrycode::countrycode(iso_code, origin = 'iso3c',
93                                              destination = 'iso2c'),
94              date,
95              cases = new_cases_per_million/1000,
96              death = new_deaths_per_million/1000
97    )
98
99  dat_covid_monthly <- dat_covid %>%
100   mutate(
101     date = lubridate::ym(paste0(lubridate::year(date), '-', lubridate::month(date)))
102   ) %>%
103   group_by(date, code) %>%
104   summarise(cases = sum(cases, na.rm = T), death = sum(death, na.rm = T)) %>%
105   ungroup()
```

```r
106
107   # Data from Eurostat ============================================================
108
109   dat_eco_sent <- eurostat::get_eurostat('ei_bssi_m_r2')
110   # Economic sentiment indicator
111
112   dat_unemployment <- eurostat::get_eurostat("une_rt_m") %>%
113   # unemployment
114     filter(age == "TOTAL", sex == "T", s_adj == "NSA", unit == "PC_ACT") %>%
115     select(date = time, code = geo, unemployment = values)
116   # Grid to facet_geo ============================================================
117
118   mygrid <- data.frame(
119     row = c(5, 1, 1, 1, 1, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 4, 4, 4, 4, 4, 4, 5, 5, 5, 5, 6,
120             6, 6, 6),
121     col = c(7, 1, 3, 4, 7, 7, 5, 4, 2, 3, 7, 2, 3, 5, 4, 4, 7, 6, 2, 5, 3, 6, 4, 5, 2, 4
122             , 7, 1, 6),
123     code = c("BG", "IS", "NO", "FI", "EE", "LV", "SE", "DK", "UK", "NL", "LT", "BE", "DE",
124             "PL", "CZ", "AT", "RO", "HU", "FR", "SK", "CH", "HR", "IT", "SI", "ES", "MT",
125             "CY", "PT", "EL"),
126     name = c("Bulgária", "Izland", "Norvégia", "Finnország", "Észtország", "Lettország",
127             "Svédország", "Dánia", "Egyesült Királyság", "Hollandia", "Litvánia",
128             "Belgium", "Németország", "Lengyelország", "Csehország", "Ausztria",
129             "Románia", "Magyarország", "Franciaország", "Szlovákia", "Svájc",
130             "Horvátország", "Olaszország", "Szlovénia", "Spanyolország", "Málta", "Ciprus",
131             "Portugália", "Görögország"),
132     stringsAsFactors = FALSE
133   )
134
135   # Automatic translation ========================================================
136
137   st_hu <- c(stopwords::stopwords('hungarian'), "is", "ha", "hozzá", "címlapfotó",
138             "illusztráció") %>%
139     {ifelse(str_starts(., "új"), NA, .)} %>%
140     na.omit()
141
142   ggpubr::ggarrange(
143     Hungary_rawtext %>%
144       filter(!str_detect(words, '\\d')) %>%
145       anti_join(data.frame(words = st_hu)) %>%
146       count(words, sort = T) %>%
147       arrange(desc(n)) %>%
148       head(30) %>%
149       mutate(
150         words = fct_reorder(words, n)
151       ) %>%
152       ggplot() +
153       aes(n, words) +
154       geom_vline(xintercept = 0) +
155       geom_col(color = 'black', fill = "gray70") +
156       labs(title = 'Magyarul', x = 'Előfordulási gyakoriság', y = NULL),
157
158     dat_sentiment %>%
159       filter(country == 'HU') %>%
```

```r
160       filter(!str_detect(words, '\\d')) %>%
161       anti_join(data.frame(words = c(stopwords::stopwords(), "also", "can"))) %>%
162       count(words, value, sort = T) %>%
163       arrange(desc(n)) %>%
164       head(30) %>%
165       mutate(
166         value = case_when(
167           value < 0 ~ "Negatív",
168           value > 0 ~ "Pozitív",
169           T ~ "Nincs"
170         ),
171         words = fct_reorder(words, n)
172       ) %>%
173       ggplot() +
174       aes(n, words, fill = value) +
175       geom_vline(xintercept = 0) +
176       geom_col(color = "black") +
177       labs(title = 'Fordítást követően', x = 'Előfordulási gyakoriság', y = NULL,
178            fill = "Adott szó szentimentje") +
179       scale_fill_manual(values = c('red4', 'gray70', 'green')) +
180       theme(
181         legend.position = 'bottom',
182         legend.direction = 'horizontal'
183       ), common.legend = T
184   )
185
186   # Explore the data ------------------------------------------------------------
187
188   ggplot(dat_sentiment_daily, aes(date, value)) +
189     geom_hline(yintercept = 0, color = "grey20") +
190     geom_line(size = .3, color = 'grey50') +
191     geom_smooth(size = 1.5, se = F) +
192     facet_geo(~ code, grid = mygrid, label = 'name') +
193     scale_x_date(limits = c(min(dat_sentiment_daily$date), max(dat_sentiment_daily$date)),
194                  breaks = c(min(dat_sentiment_daily$date), max(dat_sentiment_daily$date))) +
195     labs(y = "Szentiment", x = NULL)
196
197   library(reshape2)
198
199   dat_sentiment %>%
200     na.omit() %>%
201     mutate(
202       sentiment = ifelse(value > 0, "Pozitív", "Negatív")
203     ) %>%
204     count(words, sentiment, sort = TRUE) %>%
205     acast(words ~ sentiment, value.var = "n", fill = 0) %>%
206     comparison.cloud(colors = c("red4", "cyan4"),
207                      max.words = 100)
208
209   dat_plm <- dat_eco_sent %>%
210     filter(indic == "BS-ESI-I") %>%
211     select(date = time, code = geo, eco = values) %>%
212     merge(dat_sentiment_monthly) %>%
```

```r
    merge(dat_unemployment) %>%
    merge(dat_covid_monthly) %>%
    mutate(
      t = lubridate::interval(lubridate::ymd('2020-01-01'), date),
      t = lubridate::as.period(t) %/% months(1)
    )

# Regression tree ---------------------------------------------------------------

m_tree <- rpart::rpart(data = dat_plm, formula = value ~ .-date-code-n,
                       cp = .01)

rattle::fancyRpartPlot(m_tree, palettes = 'PuRd', sub = NULL)
```