



# Airbnb

Marcell P. Granát & Zója M. Szabó

March 6, 2021

## Contents

|                           |   |
|---------------------------|---|
| Introduction              | 3 |
| Theoretical consideration | 6 |
| Data                      | 6 |
| Explore data              | 6 |
| Model building            | 6 |
| Consequences              | 6 |
| Summary                   | 6 |
| Appendix: R codes         | 8 |

### **Abstract**

Here is the abstract.

## Introduction

[leíró statok: eu bizottság felhívása]

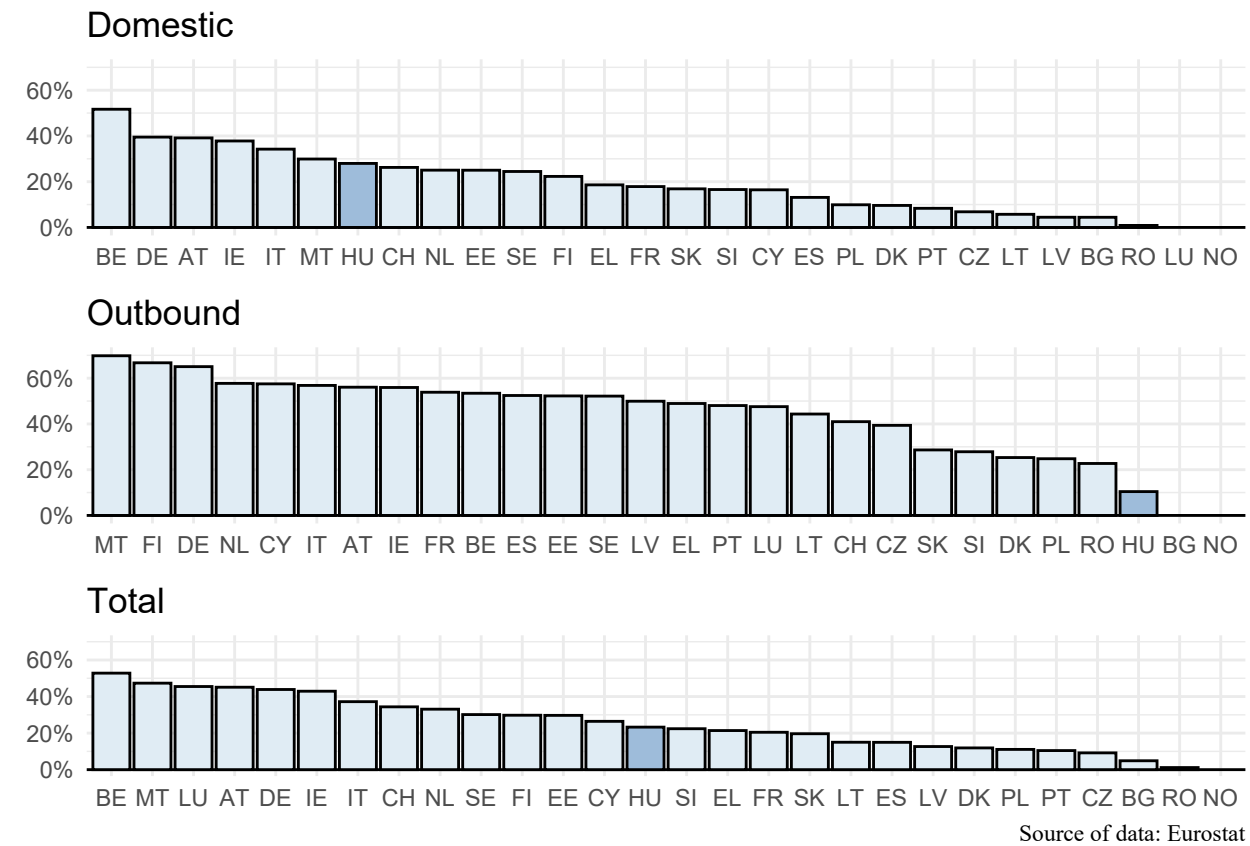


Figure 1: Proportion of internet bookings of the main means of accommodation by countries and the partner type

```
[1] "sf"          "data.frame"
```

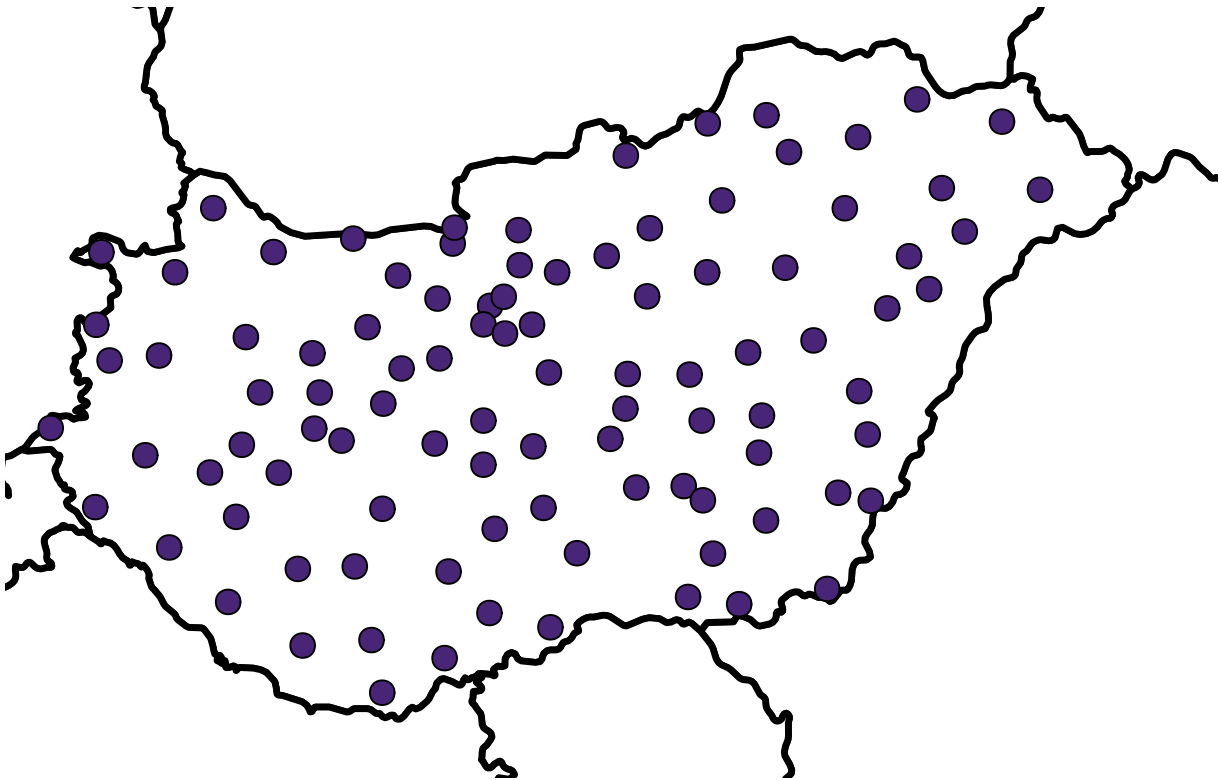


Figure 2: Starting points to our scrapping algorithm

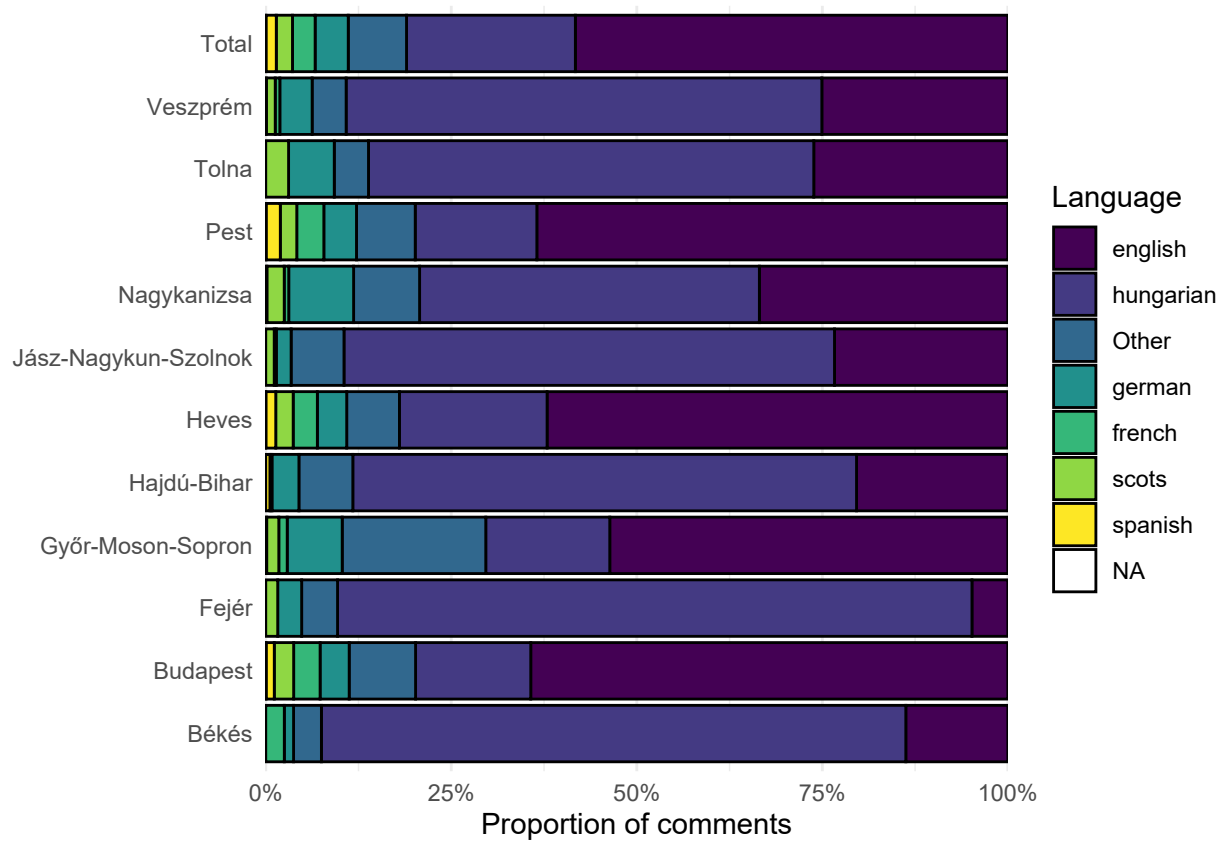


Figure 3: Most common languages found in the comments by counties



Low

## Data

## Model building

## Summary

## References

Silge, J. & Robinson, D. (2017), *Text mining with R: A tidy approach*, O'Reilly Media, Sebastopol, CA.

Zervas, G., Proserpio, D. & Byers, J. W. (2017), ‘The rise of the sharing economy: Estimating the impact of airbnb on the hotel industry’, *Journal of Marketing Research* **54**(5), 687–705.

## Appendix: R codes

```

1  # Packages -----
2
3  library(tidyverse)
4  library(knitr)
5  library(tidytext)
6
7  load('dat.RData')
8  ## Gg theme =====
9
10 update_geom_defaults("point", list(fill = "#B1339E",
11                                     shape = 21,
12                                     color = "black",
13                                     size = 1.4))
14 update_geom_defaults("line",
15                       list(color = "midnightblue", size = 1.4))
16
17 update_geom_defaults("smooth", list(color = "red4", size = 1.4))
18
19 update_geom_defaults("density",
20                       list(color = "midnightblue", fill = "midnightblue",
21                             alpha = .3, size = 1.4))
22
23 extrafont::loadfonts(device="win")
24
25 theme_set(theme_minimal() + theme(
26   legend.direction = "vertical",
27   plot.caption = element_text(family = "serif")
28 ))
29
30 f.plot_eurostat <- function(x){
31   eurostat::get_eurostat('tour_dem_ttorg', time_format = 'num') %>%
32     filter(trip_arr %in% c('ACC_WEB', 'TOTAL') & duration == 'N1-3' & time == 2017 & purpose == 'TOTAL')
33     pivot_wider(names_from = trip_arr, values_from = values) %>%
34     mutate(
35       value = ACC_WEB/TOTAL
36     ) %>%
37     filter(partner == x) %>%
38     mutate(geo = fct_reorder(geo, -value)) %>%
39     ggplot() +
40     aes(geo, value, fill = geo == 'HU') +
41     geom_hline(yintercept = 0) +
42     geom_col(color = 'black') +
43     scale_fill_brewer(palette = 3, guide = F) +
44     scale_y_continuous(labels = scales::percent, limits = c(0, .7)) +
45     labs(
46       x = NULL, y = NULL, title = case_when(
47         x == 'DOM' ~ 'Domestic',
48         x == 'OUT' ~ 'Outbound',
49         T ~ 'Total'
50       )
51     )
52 }

```



```

53
54 ggpubr::ggarrange(
55   f.plot_eurostat('DOM'),
56   f.plot_eurostat('OUT'),
57   f.plot_eurostat('WORLD') +
58     labs(caption = 'Source of data: Eurostat'),
59   ncol = 1
60 )
61
62 hun_cities <- read_csv("worldcities.csv") %>%
63   filter(country == "Hungary") %>%
64   select(city, lat, lng, admin_name, population)
65
66 cities <- readxl::read_excel("cities.xlsx")
67
68 library("rnaturalearth")
69 library("rnaturalearthdata")
70
71 world <- ne_countries(scale = "large", returnclass = "sf")
72 class(world)
73
74 merge(cities, hun_cities, by = 'city') %>%
75   tibble() %>%
76   ggplot() +
77   geom_sf(data = world, size = 1.2, fill = 'white', color = 'black') +
78   coord_sf(xlim = c(16, 23.4), ylim = c(45.5, 48.7), expand = FALSE) +
79   geom_point(aes(x = lng, y = lat), size = 4,
80             shape = 21, fill = viridis::viridis(1, begin = .1)) +
81   theme(
82     text = element_blank()
83   )
84
85 lapply(dat, function(x) {
86   tibble(city = x[["source"]][["city"]], comments = x$comments) %>%
87     mutate(language = textcat::textcat(comments))
88 }) %>%
89   reduce(rbind) %>%
90   mutate(
91     language = fct_lump(language, n = 6) %>%
92     fct_infreq()
93   ) %>%
94   merge(hun_cities) %>%
95   {rbind(., mutate(., admin_name = 'Total'))} %>%
96   mutate(admin_name = fct_reorder(admin_name, admin_name == 'Total')) %>%
97   ggplot() +
98   aes(y = admin_name, fill = language) +
99   scale_x_continuous(labels = scales::percent, limits = c(0,1), expand = c(0,0)) +
100  geom_bar(color = "black", position = position_fill()) +
101  scale_fill_viridis_d() +
102  labs(x = 'Proportion of comments', y = NULL, fill = "Language")
103
104 dat_words <- lapply(dat, function(x) {
105   tryCatch({
106     tibble(comments = x$comments) %>%

```

```

107   mutate(language = textcat::textcat(comments)) %>%
108   filter(language == "english") %>%
109   tail(-2) %>%
110   tidytext::unnest_tokens(words, comments, to_lower = T) %>%
111   mutate(assessment = x[["source"]][["assessment"]], n_reviews = x[["source"]][["n_reviews"]])
112 }, error = function(e) NULL)
113 }
114 ) %>%
115 {reduce(Filter(f = Negate(is.null), .), rbind)}
116
117 dat_words <- dat_words %>%
118   filter(n_reviews > 10 & words != 'bövebben') %>%
119   mutate(type = case_when(
120     assesment >= quantile(assessment, .99) ~ 'High',
121     assesment <= quantile(assessment, .01) ~ 'Low',
122     T ~ 'Middle'
123   ))
124
125 dat_words_High <- dat_words %>%
126   filter(type == 'High') %>%
127   group_by(words) %>%
128   summarise(n = n()) %>%
129   ungroup() %>%
130   merge(dat_words %>%
131     group_by(words) %>%
132     summarise(total = n()) %>%
133     ungroup()) %>%
134   filter(n >= 10) %>%
135   mutate(
136     rtf = n/total, type = 'High'
137   ) %>%
138   arrange(desc(rtf)) %>%
139   anti_join(data.frame(words = c(stopwords::stopwords(), "also", "can"))) %>%
140   head(50)
141
142 dat_words_Low <- dat_words %>%
143   filter(type == 'Low') %>%
144   group_by(words) %>%
145   summarise(n = n()) %>%
146   ungroup() %>%
147   merge(dat_words %>%
148     group_by(words) %>%
149     summarise(total = n()) %>%
150     ungroup()) %>%
151   filter(n >= 10) %>%
152   mutate(
153     rtf = n/total, type = 'Low'
154   ) %>%
155   arrange(desc(rtf)) %>%
156   anti_join(data.frame(words = c(stopwords::stopwords(), "also", "can"))) %>%
157   head(50)
158 rbind(dat_words_Low, dat_words_High) %>%
159   reshape2::acast(words ~ type, value.var = "rtf", fill = 0) %>%

```

```
160 wordcloud::comparison.cloud(colors = viridis::viridis(2, direction = -1 , end = .7),
161                             max.words = 100)
```