

## 3. Feladat

StatWars

2021. november 16.

## Tartalomjegyzék

Függelék: R kódok

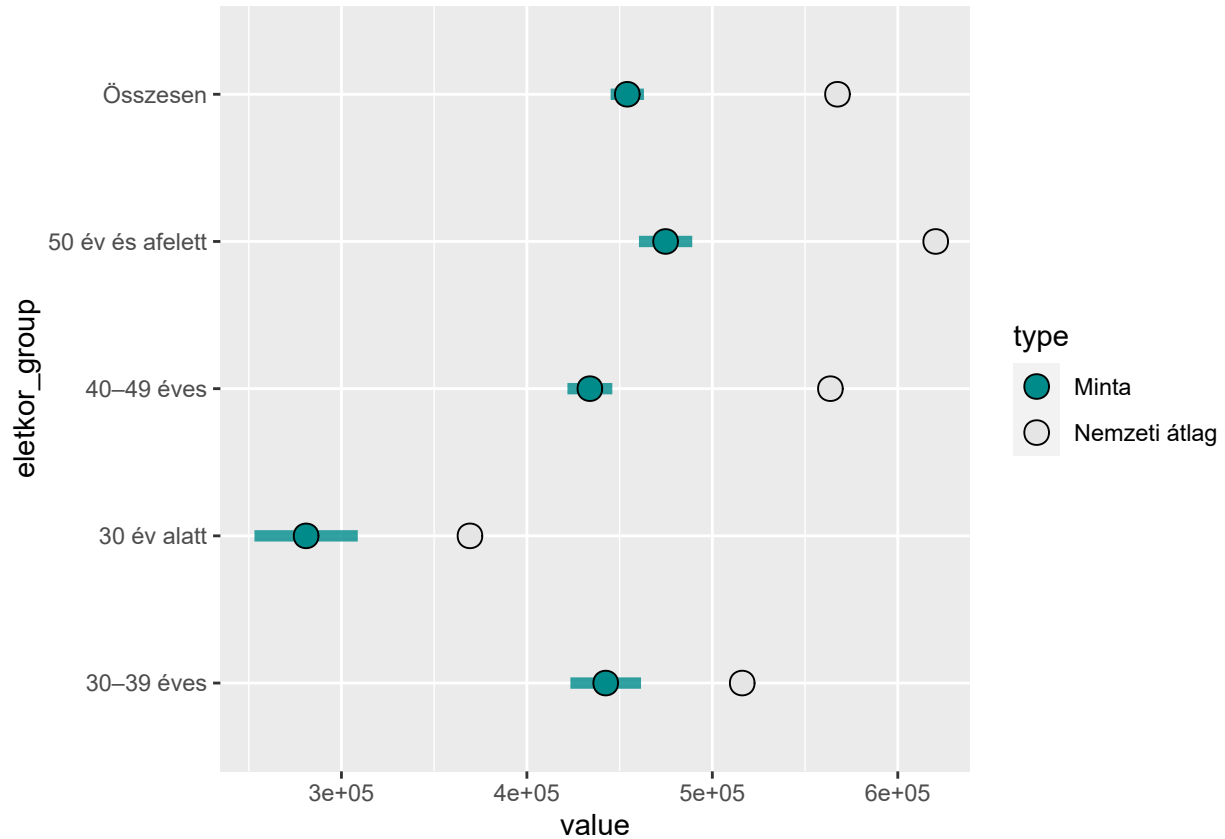
10

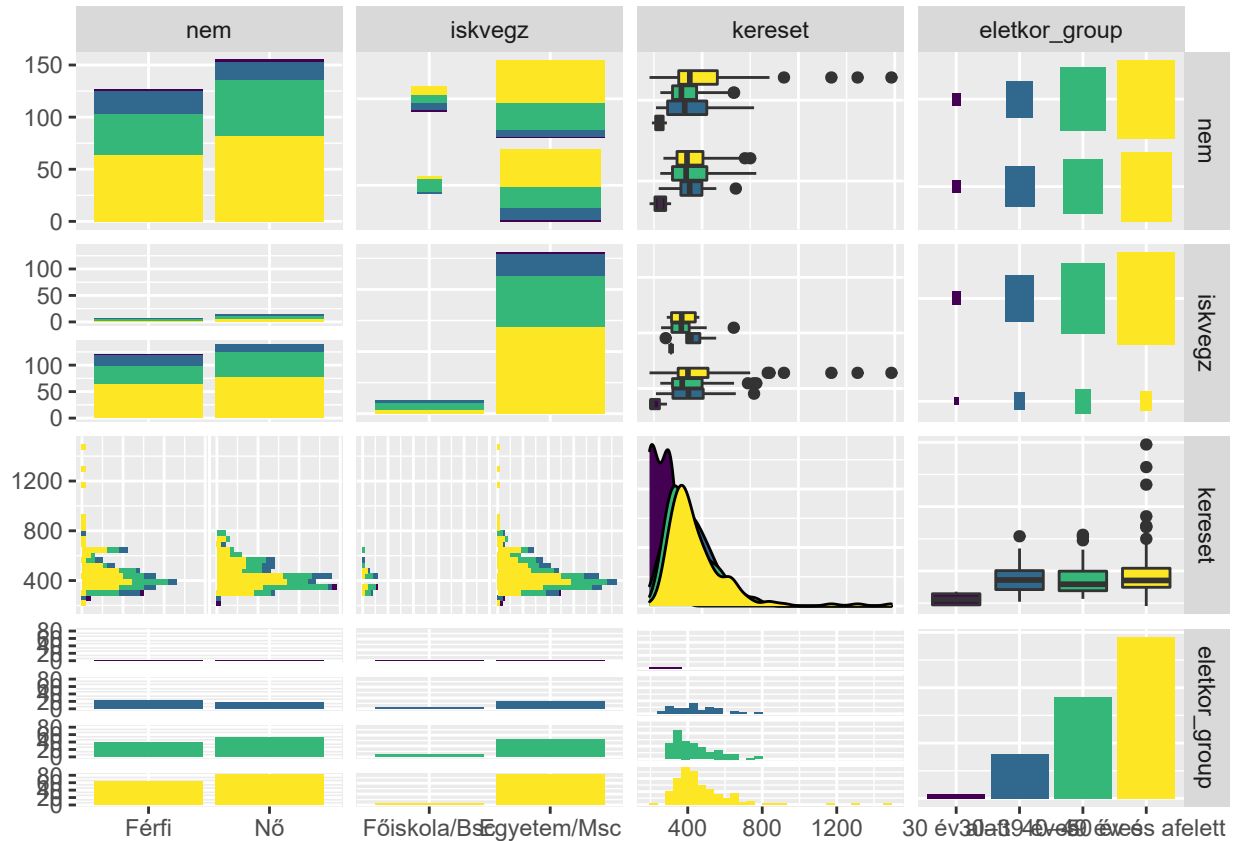
1. Az egyetemi főiskolai oktatók/tanárok fizetése életkori csoportonként hogyan különbözik egymástól? Hasonlítsák össze valamilyen benchmark adattal, ezen az egyetemen mennyivel keresnek jobban/rosszabbul az egyetemi/főiskolai oktatók/tanárok életkori csoportonkénti bontásban, mint az országos átlag?

☒ ábra a ksh összehasonlításról

☐ kereset vs. életkor boxplot

## Warning: Removed 5 rows containing missing values (geom\_segment).

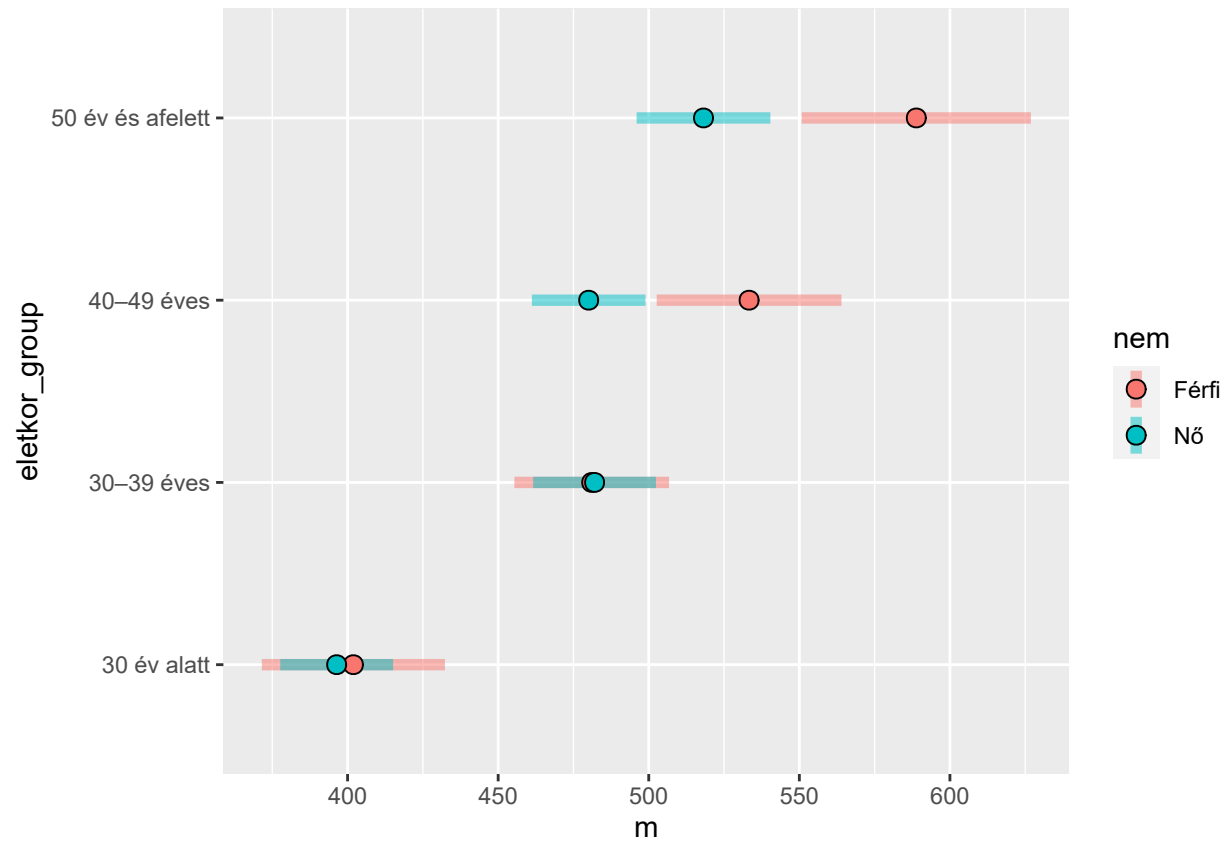




```
## # A tibble: 5 x 7
##   életkor_group   mean median    sd alpha3 kurtosis    n
##   <chr>          <dbl> <dbl> <dbl> <dbl>    <dbl> <int>
## 1 30 év alatt    398.   372.  96.9  0.636    2.50    37
## 2 30-39 éves    482.   443. 171.   1.28     4.24   116
## 3 40-49 éves    501.   419. 242.   2.14     8.46   209
## 4 50 év és afelett 544.   426. 335.   2.90    13.0   285
## 5 Összesen     511.   425. 274.   3.02    15.6   647
```

2. Hasonlítsák össze az oktatók (4-es csoport) és az ügyintézők (5 és 6-os csoport együtt) keresetek szerinti eloszlását a lehető legteljesebben!

Az eredményeket foglalják össze, ahol annak helye van érzékeltessék ábrákkal! Igyekezzenek tömören, lényegretörően végezni a számításokat! Kérjük, egy word vagy pdf fájlban legyenek az eredmények, elemzések! Excelt, vagy más szoftvert természetesen használhatnak, de azok outputja ha feltétlenül kell, függeléként lehet az elemzésükben.

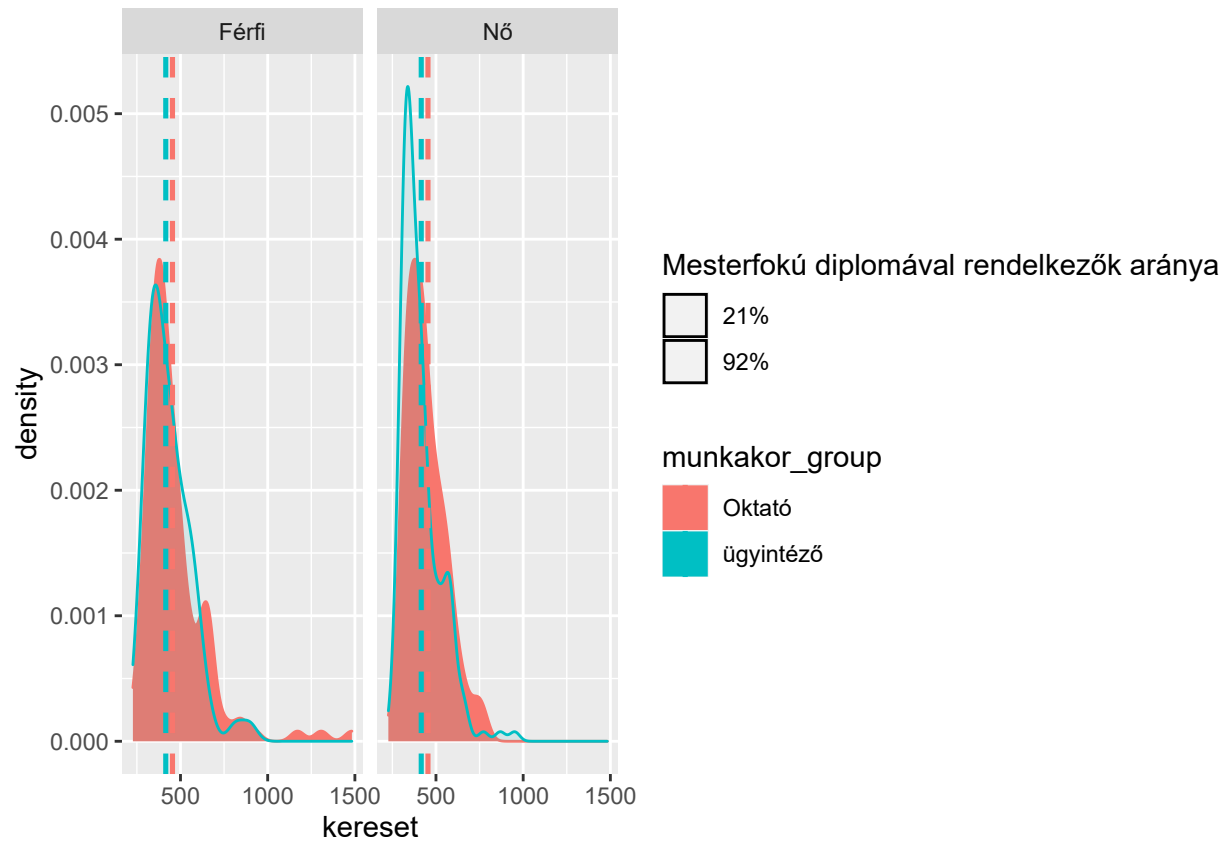




Életkor alapján való eloszlása a tanári fizetéseknek

Oktatúk és ügyintézők kereseti eloszlása

## Warning: Using alpha for a discrete variable is not advised.

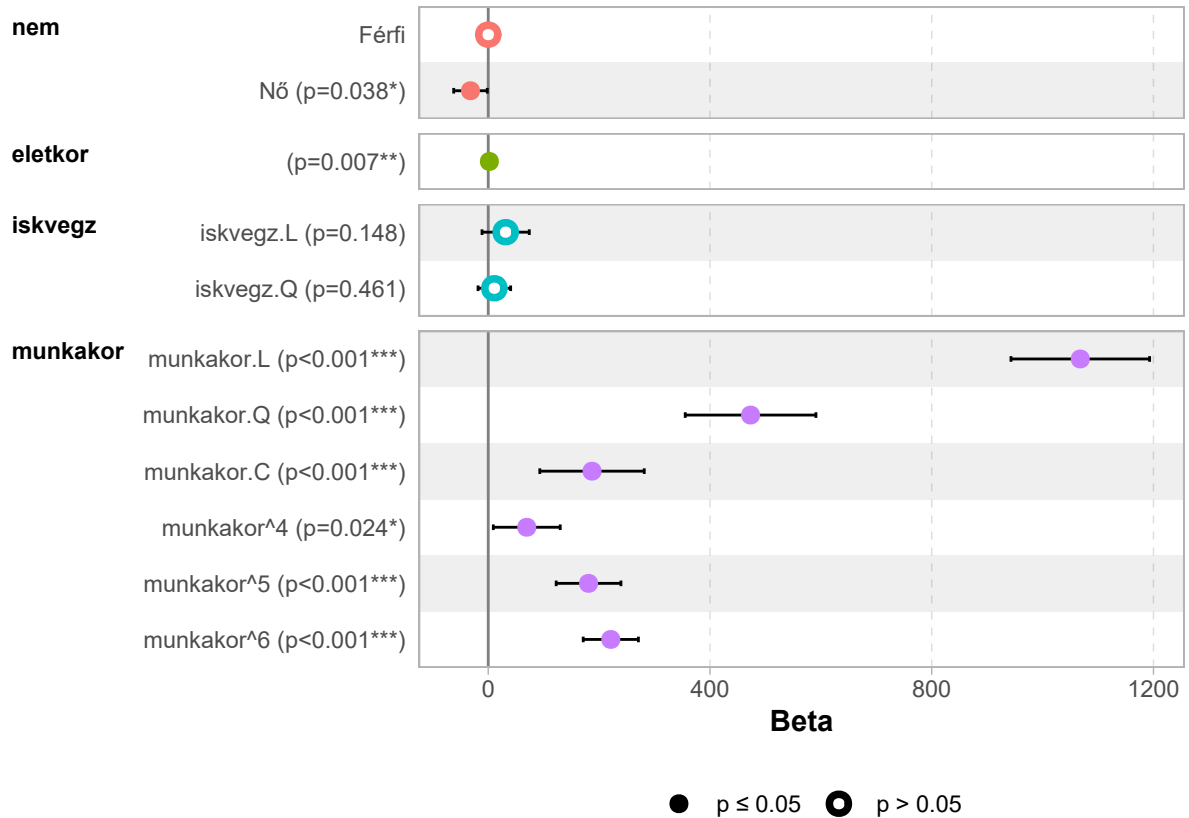




3. Készítsenek elemzést arról, hogy nem (férfi-nő) szerint a havi átlagos bruttó keresetekben mekkora az átlagos különbség összességében és az egyéb ismérvek hatását kiszűrve, illetve azokkal összekapcsolódva! Használjanak az elemzéshez kétféle módszertant/modellt és hasonlítsák össze a kétféle módszerrel kapott eredmény(ek)e)t! Írjanak egy összefoglalást is az elemzések tapasztalatairól!

p-score, ols, fa, ..

```
## # A tibble: 2 x 2
##   nem   `mean(kereset)`
##   <chr>         <dbl>
## 1 Férfi         539.
## 2 Nő           493.
```

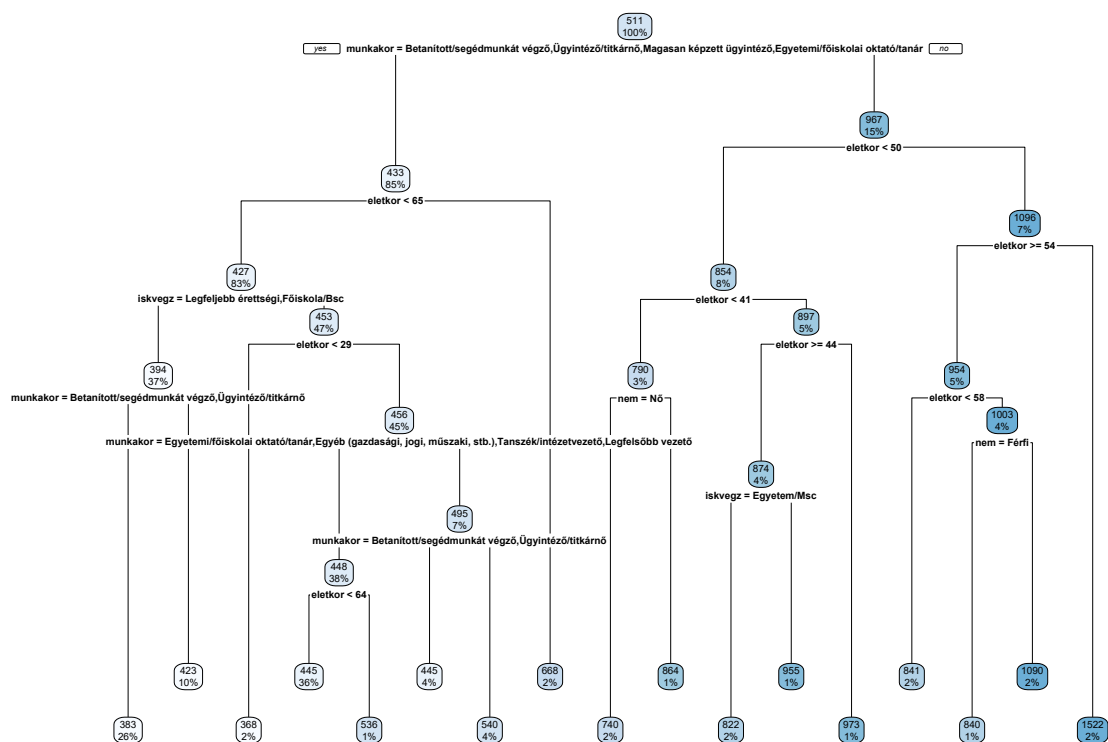


```
## # A tibble: 647 x 8
##       z nem   életkor iskveg  munkakor kereset életkor_group id
##   <dbl> <chr>   <int> <ord>   <ord>      <dbl> <ord>      <int>
## 1 0.521 Férfi    50 Egyetem/Msc Legfelsőbb vez~ 2411. 50 év és afel~ 1
## 2 0.516 Nő      51 Egyetem/Msc Legfelsőbb vez~ 1073. 50 év és afel~ 2
## 3 0.507 Férfi    53 Egyetem/Msc Legfelsőbb vez~ 1990. 50 év és afel~ 3
## 4 0.455 Nő      64 Egyetem/Msc Legfelsőbb vez~ 1609. 50 év és afel~ 4
## 5 0.393 Férfi    32 Főiskola/Bsc Tanszék/intéze~ 706. 30-39 éves 5
## 6 0.477 Nő      33 Egyetem/Msc Tanszék/intéze~ 994. 30-39 éves 6
## 7 0.384 Nő      34 Főiskola/Bsc Tanszék/intéze~ 632. 30-39 éves 7
## 8 0.380 Férfi    35 Főiskola/Bsc Tanszék/intéze~ 512. 30-39 éves 8
## 9 0.468 Férfi    35 Egyetem/Msc Tanszék/intéze~ 987. 30-39 éves 9
## 10 0.468 Nő      35 Egyetem/Msc Tanszék/intéze~ 880. 30-39 éves 10
## # ... with 637 more rows
```

```
## # A tibble: 1 x 3
##   ate atet atet_no
##   <dbl> <dbl> <dbl>
## 1 31.3 33.9 29.7
```

```
## Warning: Cannot retrieve the data used to build the model (so cannot determine roundint and is.binary)
## To silence this warning:
##   Call rpart.plot with roundint=FALSE,
##   or rebuild the rpart model with model=TRUE.
```





## Függelék: R kódok

```
1  # setup -----
2
3  library(tidyverse)
4  # data -----
5
6  teacher_df <- readxl::read_excel("3. forduló STAT WARS UNI.xlsx", sheet = 2) %>%
7    mutate(
8      nem = case_when(
9        nem == 1 ~ "Férfi",
10       nem == 2 ~ "Nő"
11     ),
12     életkor = as.integer(életkor),
13     iskveg = factor(iskveg, levels = 1:3, ordered = TRUE),
14     iskveg = fct_relabel(iskveg, function(l) {
15       case_when(
16         l == 1 ~ "Legfeljebb érettségi",
17         l == 2 ~ "Főiskola/Bsc",
18         l == 3 ~ "Egyetem/Msc"
19       )}),
20     munkakor = factor(munkakor, levels = 7:1, ordered = TRUE),
21     munkakor = fct_relabel(munkakor, function(l) {
22       case_when(
23         l == 1 ~ "Legfelsőbb vezető",
24         l == 2 ~ "Tanszék/intézetvezető",
25         l == 3 ~ "Egyéb (gazdasági, jogi, műszaki, stb.)",
26         l == 4 ~ "Egyetemi/főiskolai oktató/tanár",
27         l == 5 ~ "Magasan képzett ügyintéző",
28         l == 6 ~ "Ügyintéző/titkárnő",
29         l == 7 ~ "Betanított/segédmunkát végző"
30       )})
31   )
32  # utils -----
33
34  total_summarise <- function(x, g, ...) {
35    # original summarise function from tidyverse, but contains TOTAL row
36    bind_rows(
37      x %>%
38        group_by({{ g }}) %>%
39        summarise(...) %>%
40        ungroup(),
41      x %>%
42        summarise(...) %>%
43        mutate(g = "Összesen") %>%
44        select(g, everything()) %>%
45        rename("{ g }" := 1)
46    )
47  }
48  national_avg <- rio::import("https://www.ksh.hu/stadat_files/mun/hu/mun0059.csv") %>%
49    # download data from KSH website: https://www.ksh.hu/stadat_files/mun/hu/mun0059.html
50    tibble() %>%
51    janitor::row_to_names(2) %>%
52    select(2, starts_with("2020")) %>%
```

```

53 rename_all(str_remove_all, "2020 Korcsoport ") %>%
54 rename_all(str_remove_all, "2020 ") %>%
55 rename(profession = 1, Összesen = Együtt) %>%
56 filter(str_detect(profession, "Egyetemi")) %>%
57 mutate_at(-1, str_remove, " ") %>%
58 mutate_at(-1, as.numeric) %>%
59 pivot_longer(-1, names_to = "eletkor_group") %>%
60 select(-profession)
61 teacher_df <- teacher_df %>%
62 mutate(
63   eletkor_group = cut(eletkor, breaks = c(c(0, 3, 4, 5)*10, Inf), right = FALSE,
64     labels = FALSE),
65   eletkor_group = factor(eletkor_group, levels = 1:4, ordered = TRUE),
66   eletkor_group = fct_relabel(eletkor_group, function(l) {
67     case_when(
68       l == 1 ~ "30 év alatt",
69       l == 2 ~ "30-39 éves",
70       l == 3 ~ "40-49 éves",
71       l == 4 ~ "50 év és afelett"
72     )
73   })
74 )
75 compare_df <- bind_rows(
76 teacher_df %>%
77   filter(munkakor == "Egyetemi/főiskolai oktató/tanár") %>%
78   total_summarise(eletkor_group,
79     value = mean(kereset)*1e3,
80     s = sd(kereset*1e3),
81     n = n()
82   ) %>%
83   mutate(type = "Minta"), # TODO név
84 national_avg %>%
85   mutate(type = "Nemzeti átlag", s = NA, n = NA) # TODO teljes munkaidő hipotézise
86 )
87 compare_df %>%
88   mutate(
89     lb = value - s/(n^.5),
90     ub = value + s/(n^.5),
91   ) %>%
92   ggplot() +
93   geom_linerange(aes(xmin = lb, xmax = ub, y = eletkor_group),
94     color = "cyan4", size = 2, alpha = .8) +
95   geom_point(aes(value, eletkor_group, fill = type), shape = 21, size = 4) +
96   scale_fill_manual(values = c("cyan4", "grey90"))
97 teacher_df %>%
98   filter(munkakor == "Egyetemi/főiskolai oktató/tanár") %>%
99   select(-eletkor, -munkakor) %>%
100   GGally::ggpairs(aes(color = eletkor_group))
101 total_summarise(teacher_df, eletkor_group,
102   mean = mean(kereset),
103   median = median(kereset),
104   sd = sd(kereset),
105   alpha3 = moments::skewness(kereset),

```

```

106         kurtosis = moments::kurtosis(kereset),
107         n = n()
108     )
109     teacher_df %>%
110       group_by(eletkor_group, nem) %>%
111       summarise(m = mean(kereset), s = sd(kereset), n = n()) %>%
112       mutate(
113         cl = m - s/(n^.5),
114         ch = m + s/(n^.5)
115       ) %>%
116       ggplot() +
117       aes(m, életkor_group) +
118       geom_linerange(aes(xmin = cl, xmax = ch, color = nem), size = 2, alpha = .5) +
119       geom_point(aes(fill = nem), shape = 21, size = 3)
120
121     GGally::ggpairs(teacher_df, ggplot2::aes(colour=iskvegz))
122
123     teacher_df %>%
124       filter(
125         munkakor %in% c("Ügyintéző/titkárnő", "Magasan képzett ügyintéző", "Egyetemi/főiskolai oktató/tanár")
126       ) %>%
127       mutate(munkakor_group = ifelse(
128         munkakor == "Egyetemi/főiskolai oktató/tanár", "Oktató", "ügyintéző")
129       ) %>%
130       group_by(munkakor_group) %>%
131       mutate(iskvegz_ratio = round(sum(iskvegz == "Egyetem/Msc")/n(), 2),
132             iskvegz_ratio = scales::percent(iskvegz_ratio),
133             mean = mean(kereset)) %>%
134       ungroup() %>%
135       ggplot(aes(x = kereset, group = munkakor_group, color = munkakor_group,
136                 fill = munkakor_group, alpha = iskvegz_ratio)) +
137       geom_density() +
138       geom_vline(aes(xintercept = mean, color = munkakor_group), linetype = "dashed", size = 0.9) +
139       facet_wrap(~nem) +
140       labs(alpha = "Mesterfokú diplomával rendelkezők aránya")
141     teacher_df %>%
142       filter(
143         munkakor %in% c("Ügyintéző/titkárnő", "Magasan képzett ügyintéző", "Egyetemi/főiskolai oktató/tanár")
144       ) %>%
145       mutate(munkakor_group = ifelse(
146         munkakor == "Egyetemi/főiskolai oktató/tanár", "Oktató", "ügyintéző")
147       ) %>%
148       GGally::ggpairs(aes(color = munkakor_group))
149     teacher_df %>%
150       group_by(nem) %>%
151       summarise(mean(kereset)) # TODO
152     teacher_df %>%
153       lm(formula = kereset ~ .-életkor_group) %>%
154       GGally::ggcoef_model()
155     teacher_df %>%
156       select(- életkor_group) %>%
157       mutate(nem = nem == "Férfi") %>%
158       glm(formula = nem ~ életkor + iskvegz + munkakor, family = "binomial") %>%

```

```
159   predict( type = "response") %>%
160   cbind(teacher_df) %>%
161   rename(z = 1) %>%
162   tibble() %>%
163   mutate(id = row_number())
164 teacher_df %>%
165   group_by(nem, iskvezg, munkakor, eletkor_group) %>%
166   summarise(kereset = mean(kereset), n = n()) %>%
167   pivot_wider(names_from = nem, values_from = c(kereset, n)) %>%
168   mutate(
169     d = kereset_Férfi - kereset_Nő,
170     n = n_Férfi + n_Nő
171   ) %>%
172   ungroup() %>%
173   summarise(ate = weighted.mean(d, n, na.rm = T),
174             atet = weighted.mean(d, n_Férfi, na.rm = TRUE),
175             atet_no = weighted.mean(d, n_Nő, na.rm = TRUE))
176 teacher_df %>%
177   select(- eletkor_group) %>%
178   rpart::rpart(formula = kereset ~ ., cp = .001) %>%
179   rpart.plot::rpart.plot()
```