

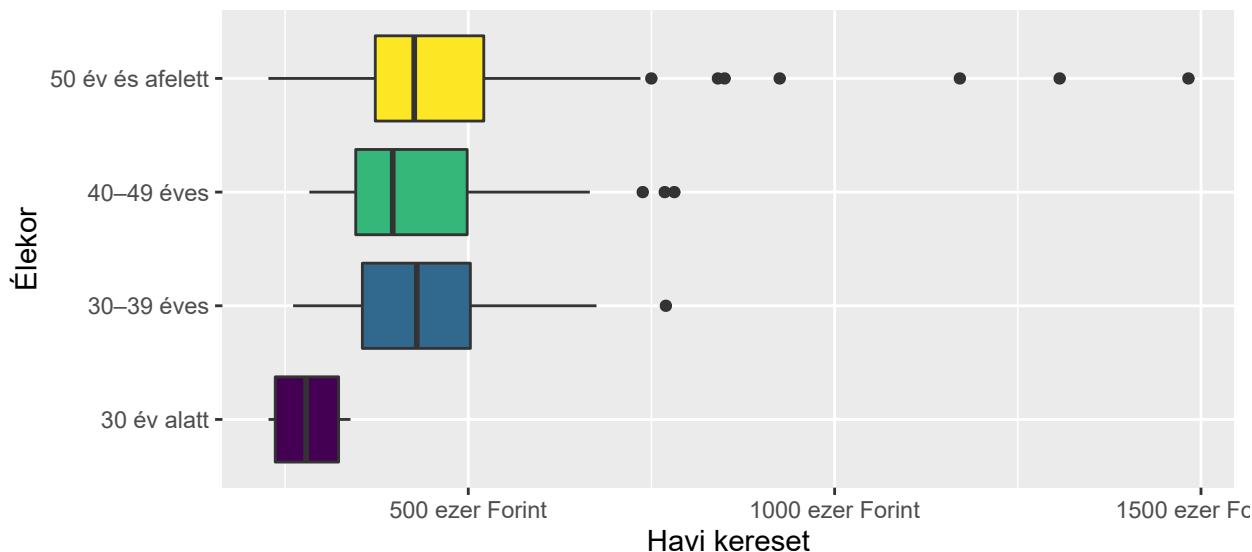
3. Feladat

StatWars

2021. november 18.

1. Feladat

A korosztályok felbontásakor figyelembe vettük a KSH módszertanát, így 4 korcsoportot alkottunk a megfigyelésekből, a 30 év alatti, 30-39 év közötti, 40-49 év közötti, és 50 év feletti csoportját. Ez alapján elmondhatjuk, hogy a kereset átlagos értéke a legalacsonyabb a 30 év alatti korosztálynál, míg korcsoportonként fokozatosan növekszik. Azonban fontos kiemelni, hogy a fizetések mediánértéke a korosztályokon belül a 30-39 évesek között a legmagasabb, így a másik két korosztálynál a kiugró értékek jobbra ferde eloszlást implikálnak. Láthatjuk, hogy a legtöbb kiugró értéket az 50 év felettiéknél találjuk (közülük is a férfiaknál), ahol akár 1 milliós bruttó fizetéssel is rendelkező oktatókat találhatunk, így a relatív szórás értéke ebben a csoportban a 60%-ot is meghaladta, míg a többiben nem érte el az 50%-ot.

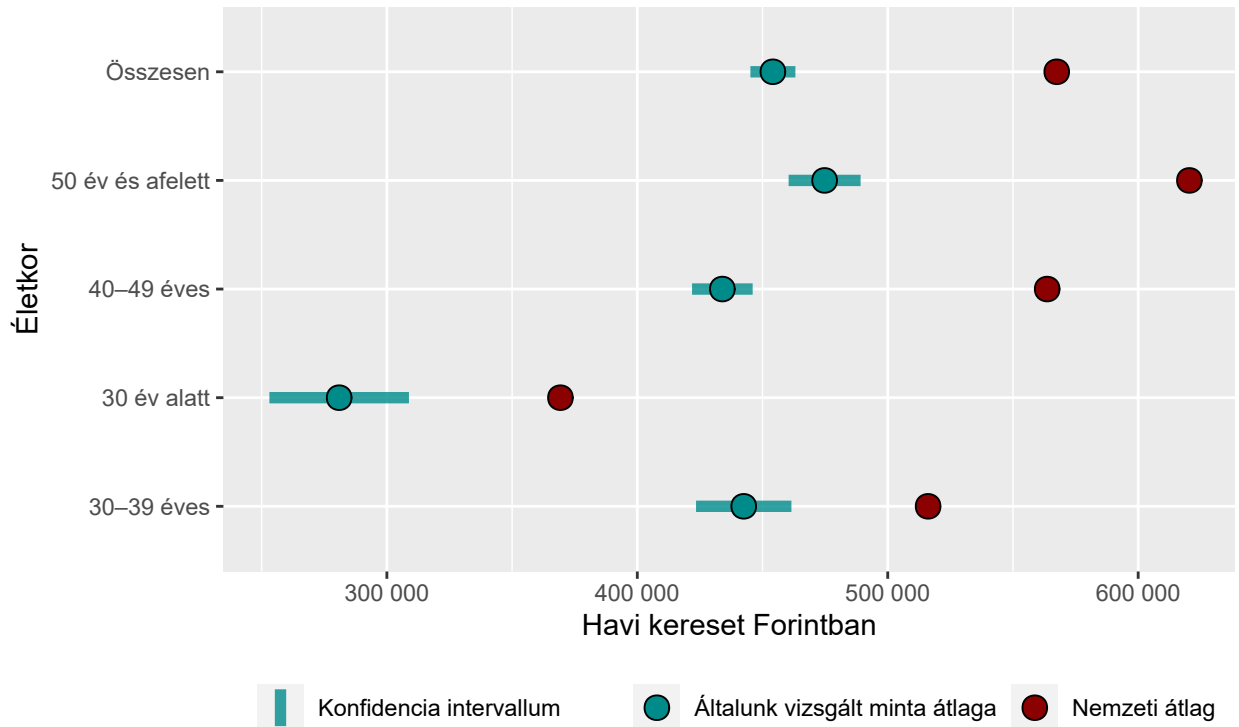


1. ábra. Egyetemi/főiskolai oktatók havi keresetének dobozábrája életkor csoportok szerinti bontásban.

Referenciaértékként a KSH 2020-as 2410-es FEOR '08 kódja (Egyetemi, főiskolai oktató, tanár) alá tartozó értékeket vizsgáltuk¹. Mivel a nemzeti bruttó átlagbér a teljes munkaidőben dolgozó oktatókra vonatkozik, így ezt az összehasonlítást megtehetjük, mivel a mintában szereplő munkavállalók is főállású alkalmazottak voltak. Összehasonlítva a mintában szereplő életkori csoportokat a KSH módszertanában megadott referenciacsoportokkal azt láthatjuk, hogy az 5%-os szignifikancia szinten vizsgált kétoldalas t-próba alapján a 2020-as országos bruttó átlagfizetések mind az 5 korcsoportban meghaladják az általunk vizsgált egyetem oktatóinak fizetéseit. Az összesen vizsgált bruttó átlagfizetés nagyjából 120 ezer forinttal volt alacsonyabb az intézményben, a legnagyobb különbséget azonban az 50 év feletti korosztályban tapasztalhattuk, nagyjából 150 ezer forintos átlagos eltéréssel.

¹https://www.ksh.hu/stadat_files/mun/hu/mun0059.html

2. FELADAT



2. ábra. Általunk vizsgált minta és az országos havi átlag keresetek összehasonlítása életkor szerinti bontásban

1. táblázat: Leíró statisztikák a életkor szerinti bontásban

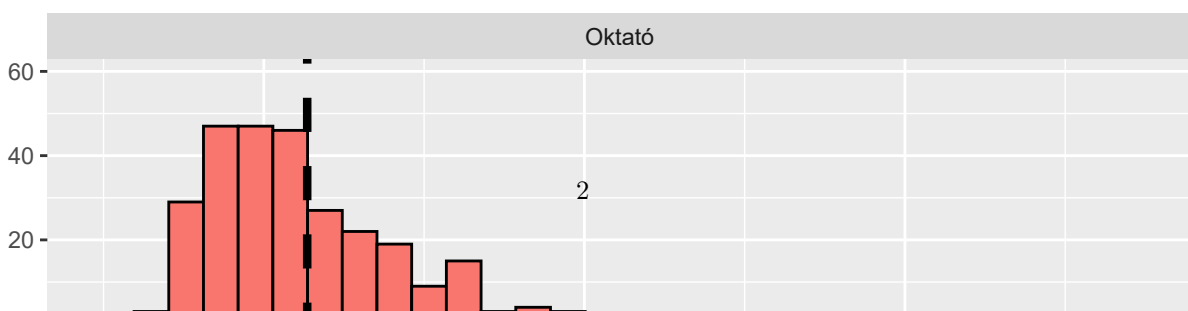
Életkor	Átlag	Medián	Szórás	Relatív szórás	Ferdeség	Csúcsosság	Elemsszám
30 év alatt	398,32	371,7	96,91	0,24	0,64	2,50	37
30–39 éves	481,62	443,2	171,49	0,36	1,28	4,24	116
40–49 éves	501,20	418,9	241,66	0,48	2,14	8,46	209
50 év és afelett	543,71	426,5	334,52	0,62	2,90	13,05	285
Összesen	510,53	424,9	274,13	0,54	3,02	15,56	647

2. Feladat

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.

- ☐ t-próba
- ☐ > 1m
- ☐ ezer HUF



kancia szinten elvetésre kerül az a nullhipotézis, miszerint a két foglalkoztatási csoportban megegyezne a sokassági átlag.

2. táblázat: Fizetések eloszlásának jellemzői munkakör jellege szerinti bontásban

Munkakör jellege	Átlag	Medián	Szórás	Relatív szórás	Ferdeség	Csúcsosság
Oktató	454,14	422,30	150,90	0,33	2,57	14,72
Ügyintéző	415,54	383,30	115,38	0,28	1,50	6,14
Összesen	435,77	404,15	136,42	0,31	2,36	13,84

3. Feladat

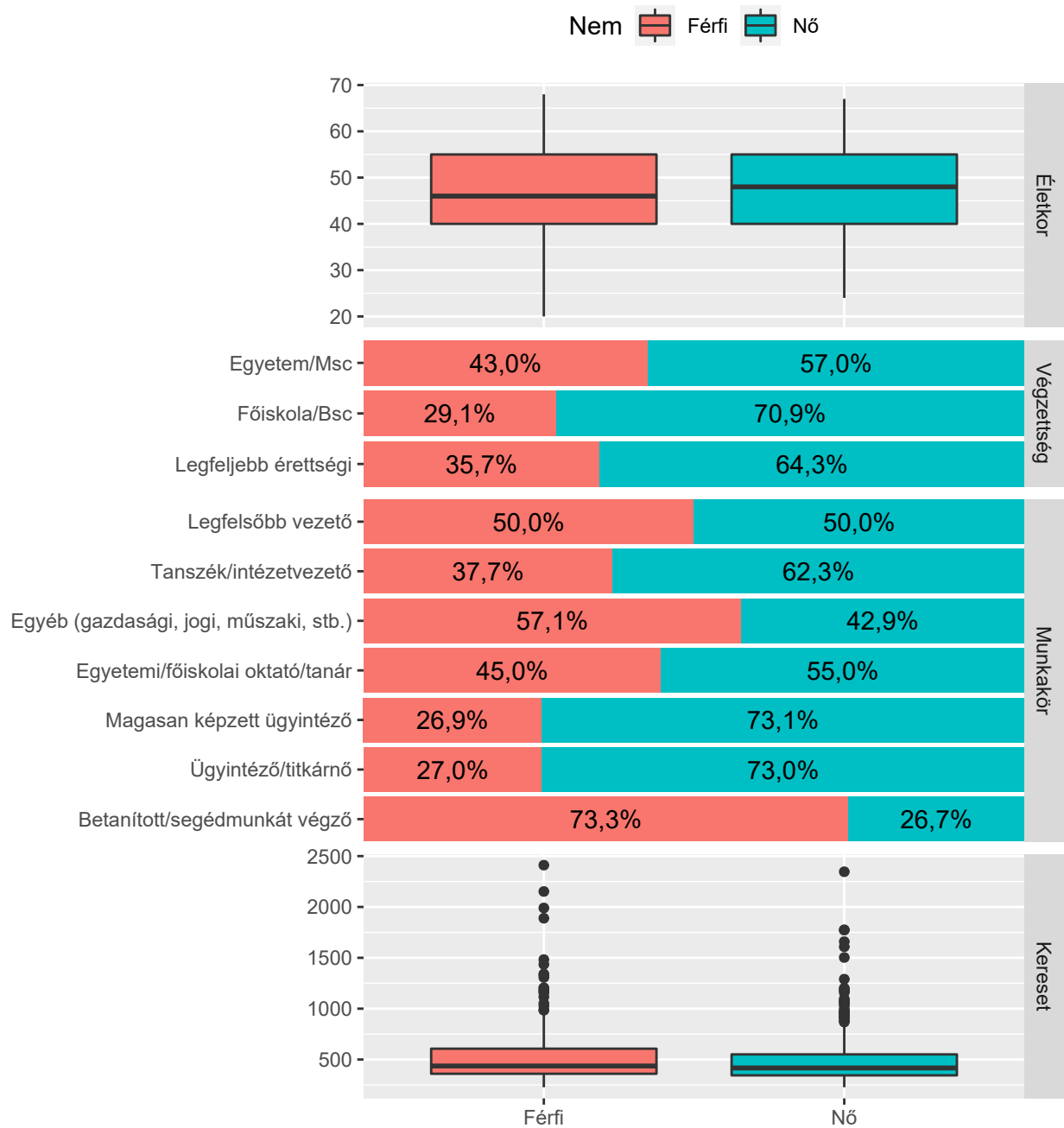
Készítsenek elemzést arról, hogy nem (férfi-nő) szerint a havi átlagos bruttó keresetekben mekkora az átlagos különbség összességében és az egyéb ismérvek hatását kiszűrve, illetve azokkal összekapcsolódva! Használjanak az elemzéshez kétféle módszertant/modellt és hasonlítsák össze a kétféle módszerrel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

p-score, ols, fa, ..

Kétoldalú alternatív hipotézis mellett a kétmintás t-próba teszt-statisztikájának értéke 2,2425 (p-érték = 0,0257), ami alapján 5%-os szignifikancia szinten elutasíthatjuk, hogy a férfi és női fizetések sokassági átlaga megegyezne. Egyoldalú nullhipotézis mellett (H_0 : Férfiak fizetése \leq Nők fizetése), a p-érték 0,0128, ami mellett 1%-os szignifikancia szinten továbbra sem tudjuk elvetni a nullhipotézist, miszerint a Férfiak fizetése szignifikánsan magasabb, mint a nőké.

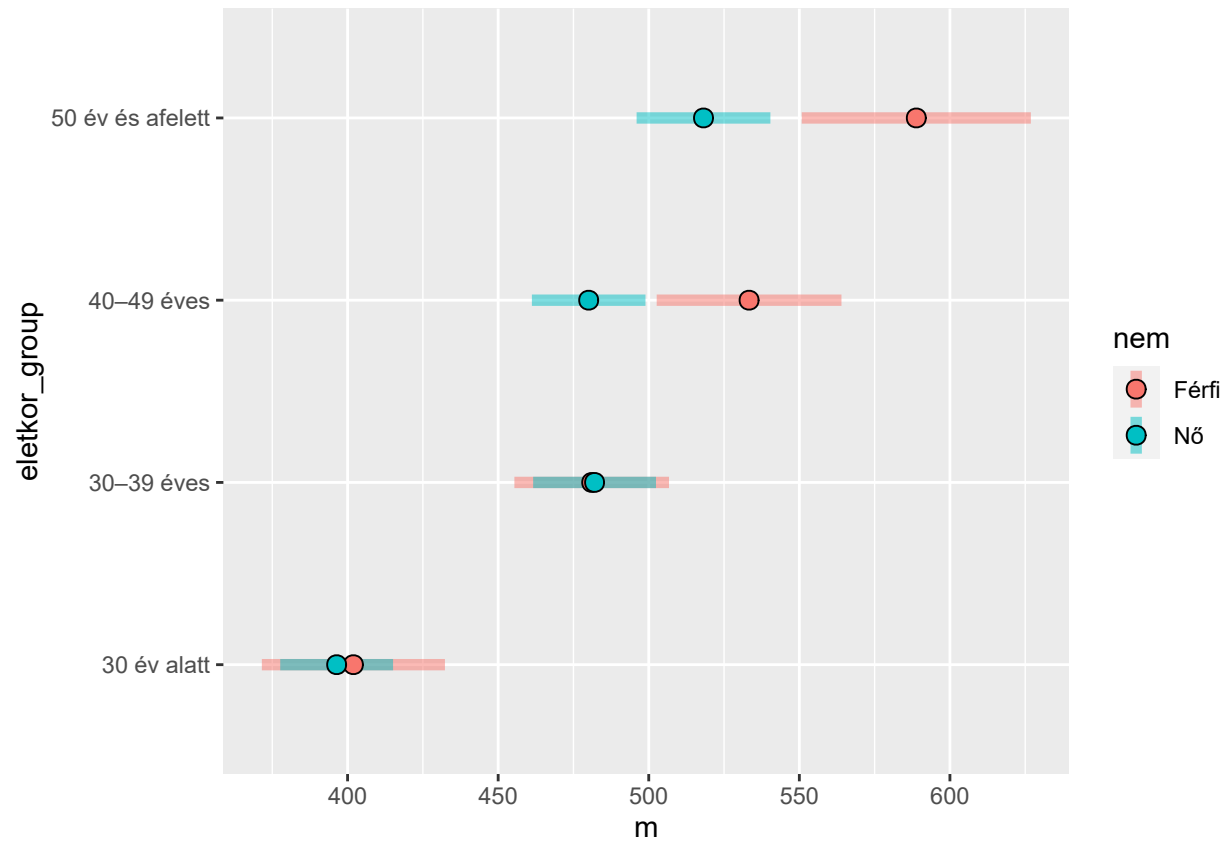
3. táblázat: Fizetések eloszlásának jellemzői nemek szerinti bontásban

Munkakör jellege	Átlag	Medián	Szórás	Relatív szórás	Ferdeség	Csúcsosság
Férfi	539,40	436,4	311,71	0,58	2,82	13,40
Nő	492,71	417,2	246,83	0,50	3,08	16,51
Összesen	510,53	424,9	274,13	0,54	3,02	15,56

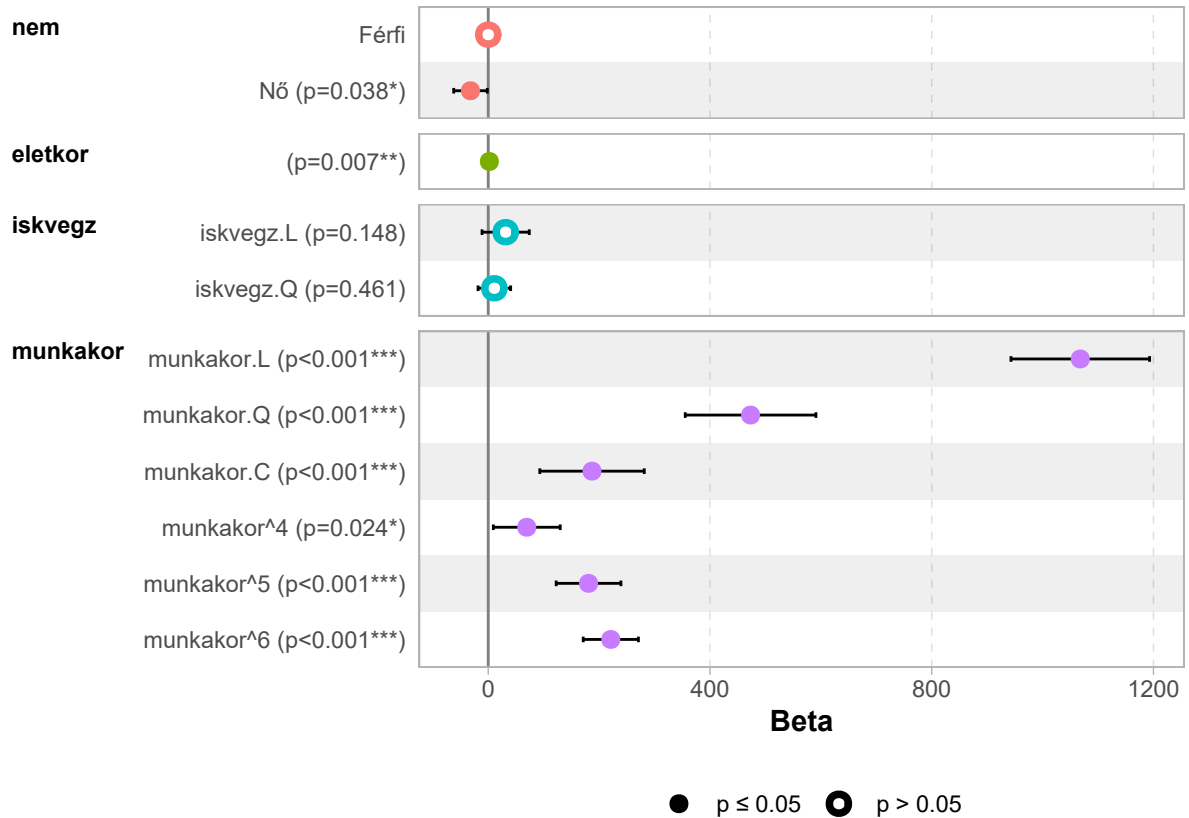


4. ábra. A vizsgált adattábla változóinak nemenkénti megoszlása

3. FELADAT



3. FELADAT



```
## # 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
```

Függelék: R kódok

```
1  # setup -----
2
3  library(tidyverse)
4  library(GGally)
5  options(scipen = 999)
6
7  # data -----
8
9  teacher_df <- readxl::read_excel("3. forduló STAT WARS UNI.xlsx", sheet = 2) %>%
10   mutate(
11     nem = case_when(
12       nem == 1 ~ "Férfi",
13       nem == 2 ~ "Nő"
14     ),
15     életkor = as.integer(életkor),
16     iskveg = factor(iskveg, levels = 1:3, ordered = TRUE),
17     iskveg = fct_relabel(iskveg, function(l) {
18       case_when(
19         l == 1 ~ "Legfeljebb érettségi",
20         l == 2 ~ "Főiskola/Bsc",
21         l == 3 ~ "Egyetem/Msc"
22       )}),
23     munkakor = factor(munkakor, levels = 7:1, ordered = TRUE),
24     munkakor = fct_relabel(munkakor, function(l) {
25       case_when(
26         l == 1 ~ "Legfelsőbb vezető",
27         l == 2 ~ "Tanszék/intézetvezető",
28         l == 3 ~ "Egyéb (gazdasági, jogi, műszaki, stb.)",
29         l == 4 ~ "Egyetemi/főiskolai oktató/tanár",
30         l == 5 ~ "Magasan képzett ügyintéző",
31         l == 6 ~ "Ügyintéző/titkárnő",
32         l == 7 ~ "Betanított/segédmunkát végző"
33       )})
34   )
35
36 teacher_df <- teacher_df %>%
37   mutate(
38     életkor_group = cut(életkor, breaks = c(c(0, 3, 4, 5)*10, Inf), right = FALSE,
39       labels = FALSE),
40     életkor_group = factor(életkor_group, levels = 1:4, ordered = TRUE),
41     életkor_group = fct_relabel(életkor_group, function(l) {
42       case_when(
43         l == 1 ~ "30 év alatt",
44         l == 2 ~ "30-39 éves",
45         l == 3 ~ "40-49 éves",
46         l == 4 ~ "50 év és afelett"
47       )
48     })
49   )
50
51 # utils -----
52
```

```

53 total_summarise <- function(x, g, ...) {
54   # original summarise function from tidyverse, but contains TOTAL row
55   bind_rows(
56     x %>%
57       group_by({{ g }}) %>%
58       summarise(...) %>%
59       ungroup(),
60     x %>%
61       summarise(...) %>%
62       mutate(g = "Összesen") %>%
63       select(g, everything()) %>%
64       rename("{ g }" := 1)
65   )
66 }
67
68 teacher_df %>%
69   filter(munkakor == "Egyetemi/főiskolai oktató/tanár") %>%
70   ggplot(aes(kereset, életkor_group, fill = életkor_group)) +
71   geom_boxplot(show.legend = FALSE) +
72   scale_x_continuous(labels = ~ str_c(., " ezer Forint")) +
73   labs(x = "Havi kereset", y = "Élekor")
74
75 national_avg <- rio::import("https://www.ksh.hu/stadat_files/mun/hu/mun0059.csv") %>%
76   # download data from KSH website: https://www.ksh.hu/stadat_files/mun/hu/mun0059.html
77   tibble() %>%
78   janitor::row_to_names(2) %>%
79   select(2, starts_with("2020")) %>%
80   rename_all(str_remove_all, "2020 Korcsoport ") %>%
81   rename_all(str_remove_all, "2020 ") %>%
82   rename(profession = 1, Összesen = Együtt) %>%
83   filter(str_detect(profession, "Egyetemi")) %>%
84   mutate_at(-1, str_remove, " ") %>%
85   mutate_at(-1, as.numeric) %>%
86   pivot_longer(-1, names_to = "életkor_group") %>%
87   select(-profession)
88
89 compare_df <- bind_rows(
90   teacher_df %>%
91     filter(munkakor == "Egyetemi/főiskolai oktató/tanár") %>%
92     total_summarise(életkor_group,
93       value = mean(kereset)*1e3,
94       s = sd(kereset*1e3),
95       n = n()
96     ) %>%
97     mutate(type = "Általunk vizsgált minta átlaga"),
98   national_avg %>%
99     mutate(type = "Nemzeti átlag", s = NA, n = NA)
100 )
101
102 compare_df %>%
103   mutate(
104     lb = value - s/(n^.5),
105     ub = value + s/(n^.5),

```



```

106 ) %>%
107 ggplot() +
108 geom_linerange(aes(xmin = lb, xmax = ub, y = életkor_group,
109                   color = "Konfidencia intervallum"), size = 2, alpha = .8) +
110 geom_point(aes(value, életkor_group, fill = type), shape = 21, size = 4) +
111 scale_fill_manual(values = c("cyan4", "red4")) +
112 scale_color_manual(values = "cyan4") +
113 scale_x_continuous(labels = ~ format(., big.mark = " ")) +
114 labs(x = "Havi kereset Forintban", y = "Életkor", color = NULL, fill = NULL) +
115 theme(
116   legend.position = "bottom"
117 )
118
119 teacher_df %>%
120   filter(munkakor == "Egyetemi/főiskolai oktató/tanár") %>%
121   select(-életkor, -munkakor) %>%
122   GGally::ggpairs(aes(color = életkor_group))
123 total_summarise(teacher_df, életkor_group,
124   `Átlag` = mean(kereset),
125   `Medián` = median(kereset),
126   `Szórás` = sd(kereset),
127   `Relatív szórás` = sd(kereset) / mean(kereset),
128   `Ferdesség` = moments::skewness(kereset),
129   `Csúcsosság` = moments::kurtosis(kereset),
130   `Elemszám` = n())
131 ) %>%
132   mutate_at(-1, ~ format(round(., 2), decimal.mark = ",")) %>%
133   rename(Életkor = 1) %>%
134   knitr::kable(caption = "Leíró statisztikák a életkor szerinti bontásban",
135     align = c("l", rep("c", 7)))
136 profession_df <- teacher_df %>%
137   filter(
138     munkakor %in% c("Ügyintéző/titkárnő", "Magasan képzett ügyintéző",
139       "Egyetemi/főiskolai oktató/tanár")
140   ) %>%
141   mutate(munkakor_group = ifelse(
142     munkakor == "Egyetemi/főiskolai oktató/tanár", "Oktató", "Ügyintéző"
143   ))
144
145 profession_df %>%
146   group_by(munkakor_group) %>%
147   mutate(
148     m = mean(kereset),
149     m = ifelse(!duplicated(m), m, NA)
150   ) %>%
151   ggplot(aes(kereset, fill = munkakor_group)) +
152   geom_histogram(color = "black", show.legend = FALSE) +
153   geom_hline(yintercept = 0) +
154   geom_vline(aes(xintercept = m, lty = "Átlag érték"), size = 1.5) +
155   facet_wrap(~ munkakor_group, ncol = 1) +
156   scale_linetype_manual(values = 2, name = NULL) +
157   scale_x_continuous(labels = ~ format(.*1e3, big.mark = " ")) +
158   theme(

```

```

159     legend.position = "bottom"
160   ) +
161     labs(x = "Havi kereset Forintban", y = "Darab")
162
163   t.test(kereset ~ munkakor_group, data = profession_df)
164
165   profession_df %>%
166     GGally::ggpairs(aes(color = munkakor_group))
167
168   profession_df %>%
169     total_summarise(g = munkakor_group,
170                     `Átlag` = mean(kereset),
171                     `Medián` = median(kereset),
172                     `Szórás` = sd(kereset),
173                     `Relatív szórás` = sd(kereset) / mean(kereset),
174                     `Ferdeség` = moments::skewness(kereset),
175                     `Csúcsosság` = moments::kurtosis(kereset),
176     ) %>%
177     mutate_at(-1, ~ format(round(., 2), decimal.mark = ",")) %>%
178     rename(`Munkakör jellege` = 1) %>%
179     knitr::kable(caption =
180       "Fizetések eloszlásának jellemzői munkakör jellege szerinti bontásban",
181       align = c("l", rep("c", )))
182
183   t.test(kereset ~ nem, data = profession_df, alternative = "two.sided")
184
185   t.test(kereset ~ nem, data = profession_df, alternative = "greater")
186
187   teacher_df %>%
188     total_summarise(g = nem,
189                     `Átlag` = mean(kereset),
190                     `Medián` = median(kereset),
191                     `Szórás` = sd(kereset),
192                     `Relatív szórás` = sd(kereset) / mean(kereset),
193                     `Ferdeség` = moments::skewness(kereset),
194                     `Csúcsosság` = moments::kurtosis(kereset),
195     ) %>%
196     mutate_at(-1, ~ format(round(., 2), decimal.mark = ",")) %>%
197     rename(`Munkakör jellege` = 1) %>%
198     knitr::kable(caption =
199       "Fizetések eloszlásának jellemzői nemek szerinti bontásban",
200       align = c("l", rep("c", )))
201   teacher_df %>%
202     select(-eletkor_group) %>%
203     rename(
204       `Életkor` = életkor,
205       `Végzettség` = iskvegz,
206       `Munkakör` = munkakor,
207       `Kereset` = kereset,
208       `Nem` = nem
209     ) %>%
210     GGally::ggbivariate(outcome = "Nem",
211       rowbar_args = list(

```

```

212         label_format = scales::label_percent(decimal.mark = ",", accuracy = .1)
213     )
214
215 )
216
217 teacher_df %>%
218   group_by(eletkor_group, nem) %>%
219   summarise(m = mean(kereset), s = sd(kereset), n = n()) %>%
220   mutate(
221     cl = m - s/(n^.5),
222     ch = m + s/(n^.5)
223   ) %>%
224   ggplot() +
225     aes(m, eletkor_group) +
226     geom_linerange(aes(xmin = cl, xmax = ch, color = nem), size = 2, alpha = .5) +
227     geom_point(aes(fill = nem), shape = 21, size = 3)
228 teacher_df %>%
229   lm(formula = kereset ~ .-eletkor_group) %>%
230   GGally::ggcoef_model()
231 teacher_df %>%
232   select(- eletkor_group) %>%
233   mutate(nem == "Férfi") %>%
234   glm(formula = nem ~ eletkor + iskveg + munkakor, family = "binomial") %>%
235   predict(type = "response") %>%
236   cbind(teacher_df) %>%
237   rename(z = 1) %>%
238   tibble() %>%
239   mutate(id = row_number())
240 matching_df <- teacher_df %>%
241   group_by(nem, iskveg, munkakor, eletkor_group) %>%
242   summarise(kereset = mean(kereset), n = n()) %>%
243   pivot_wider(names_from = nem, values_from = c(kereset, n)) %>%
244   janitor::clean_names() %>%
245   mutate(
246     d = kereset_ferfi - kereset_no,
247     n = n_ferfi + n_no
248   )
249 matching_df %>%
250   ungroup() %>%
251   summarise(ate = weighted.mean(d, n, na.rm = T),
252             atet = weighted.mean(d, n_ferfi, na.rm = TRUE),
253             atet_no = weighted.mean(d, n_no, na.rm = TRUE))

```