

Ökonometria

1. házi feladat

Granát Marcell

2020. december 1.

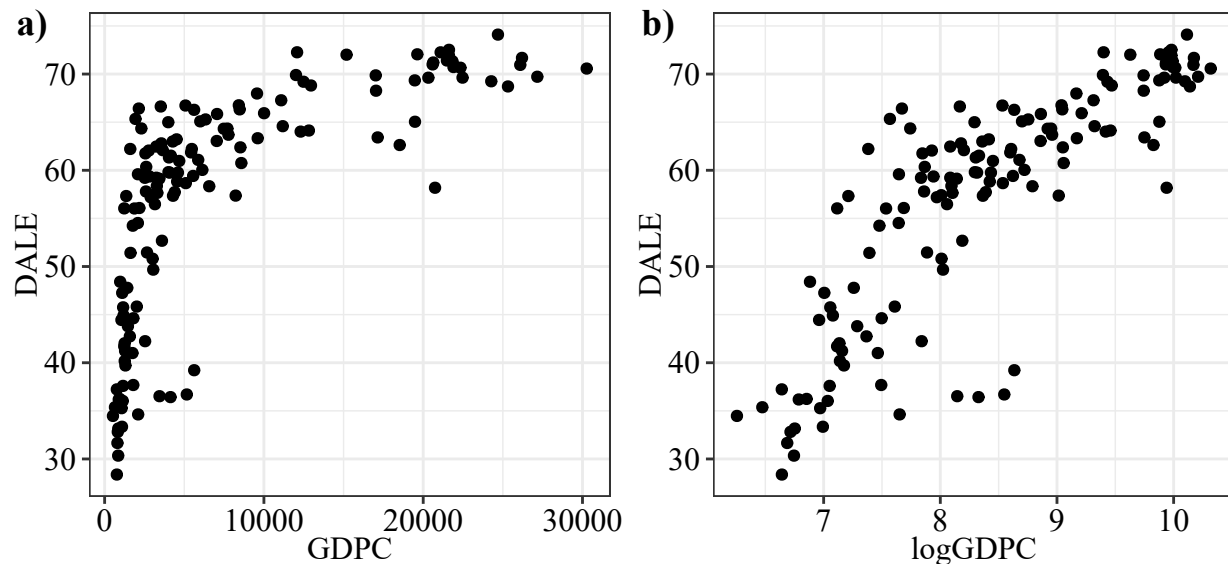
Tartalomjegyzék

1. feladat	2
a)	2
b)	2
c)	2
d)	3
e)	3
f)	3
g)	3
h)	3
2. feladat	5
a)	5
b)	5
c)	5
3. feladat	6
a)	6
b)	6
c)	6
d)	6
e)	7
f)	7
4. feladat	7
5. feladat	7
a)	7
b)	7
c)	8
d)	8
R kódok	9

1. feladat

a)

Ábrázoljuk a *DALE* változót a *GDPC* függvényében, illetve (külön ábrán) *DALE*-t a *GDPC* logaritmusára függvényében. Értelmezzük az ábrákat!



1. ábra. Az egészségkárosodással korrigált várható élettartam az egy főre eső GDP függvényében

Az 1. ábrának *a* panele ismerteti az egészségkárosodással korrigált várható élettartamot az egy főre eső bruttó kibocsátás függvényében, míg a *b* panelen ugyanezt láthatjuk, de utóbbi változó logaritmikusan átskálázott értéke szerepel a vízszintes tengelyen. Kivehető, hogy **lineáris modell jobban fog illeszkedni, amennyiben regresszorként a logaritmizált GDP/fő értéket használjuk fel (lin-log modell).**

b)

Becsüljük meg *OLS* módszerrel azt a modellt, amelyben *DALE*-t magyarázzuk a *GDPC* logaritmusával és *GINI* szintjével!

c)

Értelmezzük a *GDPC* logaritmusának együtthatóját és ellenőrizzük annak statisztikai szignifikanciáját! Gyakorlati (közgazdasági) értelemben jelentős az együttható nagysága?

1. táblázat: Az első modell paraméterei

Változó	Koefficiens	Standard hiba	T-statisztika	P-érték
konstans	-2,9	5,98	-0,49	62,83%
logGDPC	8,5	0,55	15,55	0,00%
GINI	-30,0	6,59	-4,55	0,00%

Statisztikailag szignifikánsnak bizonyul a GDP/fő logaritmus, mivel minden gyakorlatban bevett szignifikanciaszinten elutasításra kerül a H_0 , mely szerint nem különbözik a becült paraméter értéke szignifikánsan 0-tól. Amennyiben 1%-kal megnő az egy főre eső GDP értéke - minden más változatlansága mellett -, úgy

átlagosan 0,085 egységgel nő meg az egészségkárosodással korrigált várható élettartam.

Gyakorlatilag is jelentős, mivel a logaritmizált értéke az egy főre eső GDP-nek még mindig 4 egységnyi terjedelemben mozog, így a teljes intervallumon várhatóan 34,51 évnyi változást okoz, amely több mint az átlagos egészségkárosodással korrigált várható élettartam érték fele. Az előzetesen megfogalmazott elméleti megfontolásnak - magasabb életszínvonalon nagyobb a várható élettartam - megfelelő előjelet kaptunk.

d)

Értelmezzük a GINI változó együttthatóját! Gyakorlati (közgazdasági) értelemben jelentős az együtttható nagysága?

Statisztikailag szignifikánsnak bizonyul a GINI mutató, mivel minden gyakorlatban bevett szignifikanciaszinten elutasításra kerül a H_0 , mely szerint nem különbözik a becült paraméter értéke szignifikánsan 0-tól. Amennyiben 1 egységgel megnő a GINI értéke - minden más változatlansága mellett -, úgy várhatóan 30 egységgel csökken az egészségkárosodással korrigált várható élettartam.

Gyakorlatilag is jelentős, mivel a GINI mutató 0,4 egységnyi terjedelemben mozog, így a teljes intervallumon várhatóan 12,063 évnyi változást okoz az átlagos egészségkárosodással korrigált várható élettartamban. Előzetesen szakmai ismeretek hiányában nem tudok megfogalmazni feltevést a becült paraméter előjelére, de a kapott eredmény hihetőnek tűnik.

e)

Határozzuk meg a GINI paraméterbecslésének 95%-os konfidenciaintervallumát, teszteljük szignifikanciáját 1%-os szinten, és határozzuk meg a tesztstatisztika p-értékét!

A GINI becült paraméterének t-statisztikája -4,54, ami kisebb, mint az 1%-os szignifikanciaszinthez tartozó kritikus alsó érték (-2,61), tehát a nullhipotézist - miszerint nem szignifikáns a GINI hatása - elutasítjuk. A tesztthez tartozó p-érték (empirikus szignifikanciaszint) - az 1. táblázatból kiolvasható - 0,00%, tehát az előbb említett H_0 -t minden gyakorlatban bevett szignifikanciaszinten elutasítjuk.

2. táblázat: GINI paraméterbecslésének 95%-os konfidenciaintervalluma

Alsó határ	Felső határ
-43,03	-16,96

f)

Következtethetünk-e az eredmények alapján arra, hogy a magasabb Gini-együtttható alacsonyabb egészségkárosodással korrigált várható élettartamot okoz? Miért vagy miért nem?

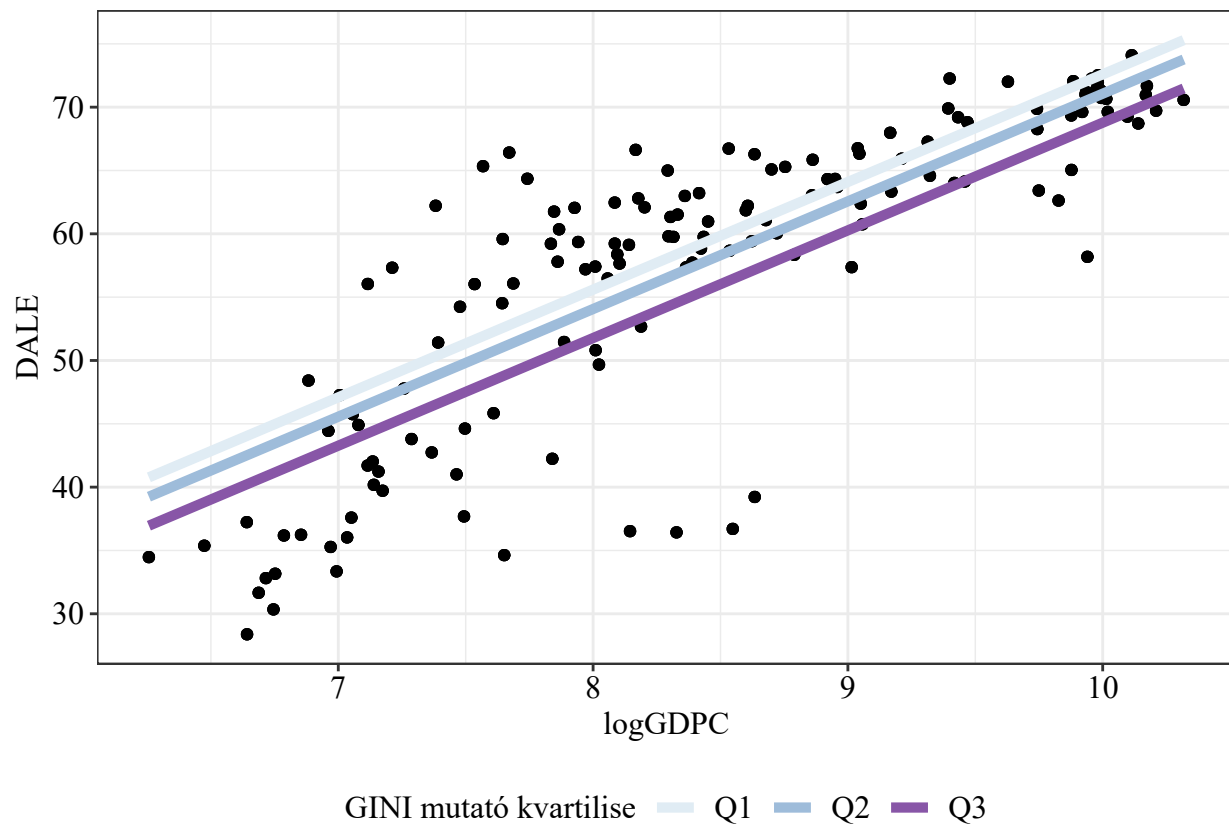
A modellben szereplő parciális hatás egyértelműen ezt az eredményt sugallja. Mindazontál szükséges lenne még megvizsgálni a teljes hatást (kiszámítani a logGDPC-n keresztüli közvetett hatást), illetve gondolni kell a kihagyott változók okozta torzításra is.

g)

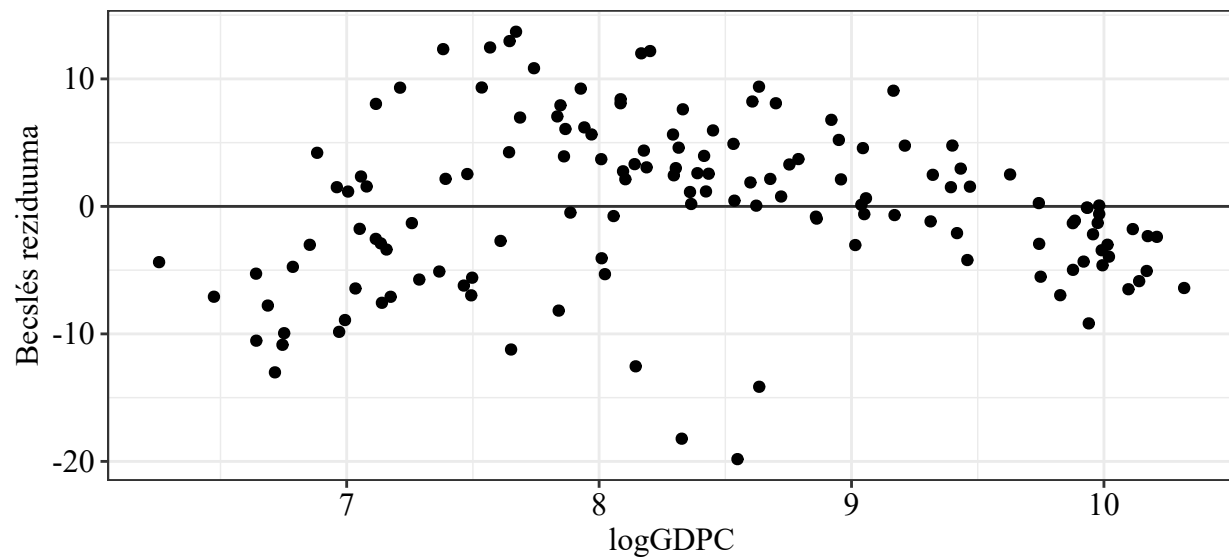
Az a) rész logaritmikus ábráján rajzoljuk be a DALE becült regressziós függvényét a logaritmikus GDPC függvényében, ha GINI az eloszlásának alsó kvartilisé, mediánját illetve felső kvartilisé veszi fel!

h)

Ábrázoljuk a reziduálisokat a logaritmikus GDPC függvényében!



2. ábra. Az első modellben a GDP/fő hatása különböző GINI értékek mellett



3. ábra. Az első modell becsléséből származó reziduumok

2. feladat

a)

A legutolsó ábra sugallata alapján bővítsük ki a modellt a logaritmusos GDPC négyzetével!

b)

Az 1. vagy a 2. rész modelljét választanánk? Miért?

3. táblázat: Az 1. és 2. modell jellemzői

	R négyzet	Korrigált R négyzet	AIC	BIC
1. modell	0,72	0,72	929,49	941,28
2. modell	0,78	0,78	897,96	912,71

A 2. modellt választanánk, ugyanis mind az R^2 , mind a korrigált R^2 értéke magasabb a 2. modellben, illetve a közölt információs kritériumok¹ (AIC, BIC) értékei alacsonyabbak. Így minden illeszkedés jóságát jellemző mutató alapján arra a döntésre jutunk, hogy a 2. modell jobban írja le a regressziós kapcsolatot.

c)

Értelmezzük a modellből származó előrejelzés és reziduális értékét egy tetszőlegesen választott országra!

4. táblázat: A 2. modell Magyarországra készült becslése

Valós érték	Becsülés	Reziduum
63,04	67,08	-4,04

A 3. táblázatból kiolvasható, hogy a modell Magyarország 1993-as egészségkárosodással korrigált várható élettartamának 67,08 évet becsül, ami 4,04 évvel magasabb, mint a valós érték, amely 63,04 év.

¹Hibára alapuló mutatók, így értéküket minimalizálni kell.

3. feladat

a)

Becsüljük meg DALE modelljét a GDPC (nem pedig a logaritmikus GDPC) négyzetes függvényét és a GINI szintjét használva magyarázó változóként!

b)

Teszteljük 1%-os szinten, hogy a GDPC és négyzete együttesen szignifikáns-e ebben a regresszióban!

A teszthez Wald-féle F-próbát hajtok végre, melynek nullhipotézise, hogy $\beta_{GDPC} = \beta_{GDPC^2} = 0$. Az F-próba értéke 84,341, amely minden gyakorlatban bevett szignifikanciaszinthez tartozó felső kritikus értéket meghalad (a p-érték 0,00%). Mivel a H_0 -t elutasítjuk, így kijelenthetjük, hogy a két tárgyal regresszor együttesen szignifikáns a modellben.

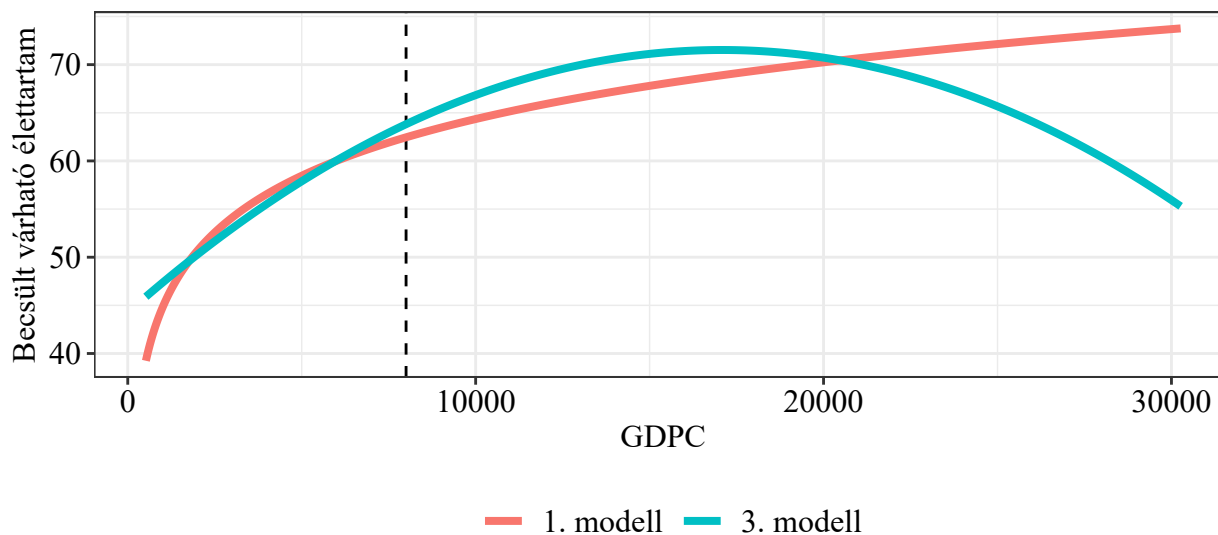
c)

A fenti modell alapján számítsuk ki GDPC parciális hatását, ha GDPC egy, a hallgató vezetéknévének kezdőbetűje által meghatározott értéket vesz fel. (GDPC=8000 USD A-F kezdőbetű esetén, GDPC=10000 USD G-P kezdőbetű esetén és GDPC=12000 USD Q-Z kezdőbetű esetén.)

Amennyiben 8000-ről 8001-re nőne a GDP/fő egy országban, úgy minden más változatlanlansága mellett várhatóan **19,56** évvel nőne meg ott az egészségkárosodással korrigált várható élettartam.

d)

Hasonlítsuk össze az eredményt az 1. részben kapott parciális hatással!



4. ábra. Az 1. és 3. modellel készített becslés a várható élettartamra a GDP/fő függvényében

Mivel a 3. modellben a $GDPC^2$ -hez tartozó becslés paraméter előjele negatív, így az első modellel konzisztens módon ellaposodó hatása van a GDP/fő-nek a várható élettartamra. Az előbbieken alkalmazott 8000 USD pontban - medián GINI érték mellett - a 3. modell parciális hatása nagyobb (meredekebb az egyenes).

e)

Számítsuk ki a parciális hatás standard hibáját!

A parciális hatás standard hibája **0,00031**.

f)

Az 1. vagy a 3. rész modelljét választanánk? Miért?

5. táblázat: Az 1. és a 3. modell jellemzői

	R négyzet	Korrigált R négyzet	AIC	BIC
1. modell	0,72	0,72	929,49	941,28
3. modell	0,66	0,65	961,09	975,84

Az 1. modellt választanánk, ugyanis mind az R^2 , mind a korrigált R^2 értéke magasabb az 1. modellben, illetve a közölt információs kritériumok (AIC, BIC) értékei alacsonyabbak. Így minden illeszkedés jóságát jellemző mutató alapján arra a döntésre jutunk, hogy az 1. modell jobban írja le a regressziós kapcsolatot.

4. feladat

Összességében, melyik modellt választanánk azon modellek közül, amelyek a GINI szintje mellett a GDPC vagy a logaritmikus GDPC tetszőleges polinomját használják magyarázó változóként (azaz a modellhalmaz az 1., 2. és 3. rész modelljeit is tartalmazza speciális esetként)?

A legjobban illeszkedő változók körének meghatározásához kiindulási modellnek vettem azt, amelyikben csak a GINI szerepel, mint magyarázóváltozó. Ezt követően AIC információs kritérium minimalizálási céllal bővítettem a modellt. Minden egyes lépésnél a GDP/fő és a GDP/fő logaritmusának polinómjai (maximum 6 rendű) közül azt emeltem be a modellbe, amellyel az új modell AIC információs kritériuma a legalacsonyabb volt. Akkor áltam meg a bővítéssel, mikor bármely változó bevonásával csökkent volna az AIC. Így legjobb modellnek mutatkozik az, melyben a GDP/fő logaritmusa és a GDP/fő logaritmusának négyzete szerepel.

5. feladat

A $\log(\text{bér})$ -t mint függő változót modellezzük az IQ-val mint magyarázó változóval egyváltozós regresszió segítségével, egy elég nagy mintán. A $\log(\text{bér})$ átlaga 12 és szórása 0,5, míg IQ átlaga 100 és szórása 15. A két változó mintabeli korrelációja 0,4.

a)

Számítsuk ki a regresszió R-négyzet értékét!

$$R^2 = (r)^2 = 0,4^2 = 0,16$$

b)

Számítsuk ki a hibateg varianciáját!

$$MSE = \sigma_y^2 \times (1 - R^2) = 0,21$$

c)*Mi a meredekségi paraméter OLS becslése?*

$$\hat{\beta}_1 = \hat{\rho} \frac{\hat{\sigma}_y}{\hat{\sigma}_x} = 0,4 \frac{0,5}{15} = 0,013\bar{3}$$

d)*Mi a tengelymetszet OLS becslése?*

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 12 - 0,4 \frac{0,5}{15} \times 100 = 10,6\bar{6}$$

R kódok

```

1  # setup -----
2  library(tidyverse)
3  library(granatlib) # my personal package: https://github.com/MarcellGranat/granatlib
4  theme_set(theme_granat())
5  dat <- rio::import("health_small.xlsx") %>% filter(YEAR == 1993)
6          # data import, NEPTUN: AYCOPF
7
8  # 1 #####
9
10 # a -----
11 dat <- dat %>% mutate(logGDPC = log(GDPC)) # add the log of GDPC
12
13 ggpubr::ggarrange(
14   ggplot(dat, aes(GDPC, DALE)) + geom_point(),
15   ggplot(dat, aes(logGDPC, DALE)) + geom_point(), labels = c("a)", "b)")
16 )
17
18 # b -----
19 modell <- lm(data = dat, formula = DALE ~ logGDPC + GINI)
20
21 # c, d -----
22 modell %>% broom::tidy() %>% prtbl("Az első modell paraméterei", ufc = F)
23
24 # e -----
25 confint(modell, 'GINI', level = .95) %>% data.frame() %>%
26   set_names("Alsó határ", "Felső határ") %>%
27   knitr::kable(format.args = list(decimal.mark = ","), digits = 2,
28               caption = "GINI paraméterbecslésének 95%-os konfidenciaintervalluma",
29               row.names = F, align = c("c", "c"))
30
31 # f -----
32 data.frame(
33   logGDPC = dat$logGDPC, DALE = dat$DALE,
34   Q1 = dat %>% select(logGDPC, GINI) %>%
35     mutate(GINI = quantile(GINI, .25)) %>% predict.lm(object = modell),
36   Q2 = dat %>% select(logGDPC, GINI) %>%
37     mutate(GINI = quantile(GINI, .5)) %>% predict.lm(object = modell),
38   Q3 = dat %>% select(logGDPC, GINI) %>%
39     mutate(GINI = quantile(GINI, .75)) %>% predict.lm(object = modell)
40 ) %>% pivot_longer(3:5) %>%
41   ggplot() + geom_point(aes(logGDPC, DALE)) +
42   geom_line(aes(logGDPC, value, color = name), size = 1.8) +
43   scale_color_brewer(palette = "BuPu") +
44   labs(color = "GINI mutató kvartilise")
45
46 modell %>% broom::augment() %>% ggplot() +
47   geom_hline(yintercept = 0, color = "grey20") +
48   geom_point(aes(x = logGDPC, y = .resid)) +
49   labs(y = "Becslés reziduuma")
50
51 # 2 #####
52

```

```

53 # a -----
54 dat <- dat %>% mutate(logGDPC2 = logGDPC^2)
55 model2 <- lm(data = dat, formula = DALE ~ logGDPC + logGDPC2 + GINI)
56
57 # b -----
58
59 rbind(broom::glance(model1), broom::glance(model2)) %>% mutate(
60   model = c("1. modell", "2. modell")
61 ) %>% column_to_rownames(var = 'model') %>% select(
62   r.squared, adj.r.squared, AIC, BIC
63 ) %>% rename(c("R négyzet" = r.squared, "Korrigált R négyzet" = adj.r.squared)) %>%
64   knitr::kable(digits = 2, format.args = list(decimal.mark = ","),
65               caption = "Az 1. és 2. modell jellemzői", align = rep("c", ncol(.)))
66
67 # c -----
68 model2 %>% broom::augment() %>% cbind(dat$COUNTRYNAME) %>%
69   filter(dat$COUNTRYNAME == "Hungary") %>%
70   select(DALE, .fitted, .resid) %>%
71   set_names("Valós érték", "Becsülés", "Reziduum") %>% knitr::kable(
72     caption = "A 2. modell Magyarországra készült becslése",
73     align = c("c", "c", "c"), format.args = list(decimal.mark = ","), digits = 2
74   )
75
76 # 3 #####
77
78 # a -----
79 dat <- dat %>% mutate(GDPC2 = GDPC^2)
80 model3 <- lm(data = dat, formula = DALE ~ GINI + GDPC + GDPC2)
81
82 # b -----
83 car::linearHypothesis(model3, c("GDPC = 0", "GDPC2 = 0"), test="F")
84
85 # c -----
86 # Neptun: AYCOFF -> GDPC = 8000
87 GDPC.partialeffect <- format(sum(model3$coefficients[3:4]*
88                               c(8000, 8000^2)), digits = 4, decimal.mark = ",")
89
90 # d -----
91 data.frame(GINI = rep(median(dat$GINI), 1000),
92            GDPC = seq(from = min(dat$GDPC), to = max(dat$GDPC), length.out = 1000)) %>%
93   mutate(logGDPC = log(GDPC), GDPC2 = GDPC^2) %>%
94   cbind(data.frame(model1 = predict.lm(object = model1, newdata = .))) %>%
95   cbind(data.frame(model3 = predict.lm(object = model3, newdata = .))) %>%
96   select(GDPC, model1, model3) %>% set_names("GDPC", "1. modell", "3. modell") %>%
97   pivot_longer(-1) %>%
98   ggplot(aes(x = GDPC, y = value, color = name)) +
99   geom_vline(xintercept = 8000, linetype = "dashed") +
100  geom_line(size = 1.4) +
101  labs(y = "Becsült várható élettartam", color = "")
102
103 # e -----
104 answer_3e <- vcov(model3) %>%
105   {.[ "GDPC", "GDPC" ] + .[ "GDPC2", "GDPC2" ] + .[ "GDPC", "GDPC2" ]} %>%
106   sqrt()

```

```

106
107 # f -----
108
109 rbind(broom::glance(model1), broom::glance(model3)) %>% mutate(
110   model = c("1. modell", "3. modell")
111 ) %>% column_to_rownames(var = 'model') %>% select(
112   r.squared, adj.r.squared, AIC, BIC
113 ) %>% rename(c("R négyzet" = r.squared, "Korrigált R négyzet" = adj.r.squared)) %>%
114   knitr::kable(digits = 2, format.args = list(decimal.mark = ","),
115               caption = "Az 1. és a 3. modell jellemzői")
116
117 # 4 #####
118 dat2 <- dat %>% select(DALE, GDPC, GDPC2, GINI, logGDPC, logGDPC2) %>%
119   mutate(
120     GDPC3 = GDPC^3,
121     logGDPC3 = logGDPC^3,
122     GDPC4 = GDPC^4,
123     logGDPC4 = logGDPC^4,
124     GDPC5 = GDPC^5,
125     logGDPC5 = logGDPC^5,
126     GDPC6 = GDPC^6,
127     logGDPC6 = logGDPC^6
128   )
129
130 model4 <- lm(data = dat2, formula = DALE ~ GINI)
131
132 MASS::stepAIC(model4, scope = list(lower = DALE ~ GINI,
133   upper = DALE ~ GINI + GDPC + GDPC2 + GDPC3 + GDPC4 + GDPC5 + GDPC6 +
134   logGDPC + logGDPC2 + logGDPC3 + logGDPC4 + logGDPC5 + logGDPC6),
135   direction = "forward", trace = F)

```