

Seminar_5

Zoltan Kekecs

April 3, 2023

Ket változó kapcsolatának vizsgálata, statisztikai inferencia

Az óra célja

Az óra célja hogy megismerkedjünk a **statisztikai inferencia alapjaival** két változó kapcsolatának elemzésén keresztül.

Package-ek betöltése

```
if (!require("tidyverse")) install.packages("tidyverse")  
library(tidyverse) # for dplyr and ggplot2
```

Hipotezisteszteles

A statisztikai inferencia, és hipotézis teszteles során az a célunk, hogy megállapítsuk, letezik-e egy bizonyos hatás vagy kapcsolat. De ezt a **null-hipotézis szignifikancia teszteles (NHST)** során egy fordított logikával tesszük: azt állapítjuk meg, hogy **mekkora a valószínűsége hogy az általunk megfigyelt adatot/trendet figyeljük meg (vagy annál meg extremer trendet), amennyiben a null-hipotézis igaz.**

Egy egyszerű példa: az a sejtés, hogy **egy penzermé cinkelt** (vagyis ki van sulyozva hogy az egyik oldalára nagyobb eséllyel essen mint a másik oldalára), megpedig úgy hogy nagy valószínűséggel **fej** legyen az eredmény amikor feldobjuk. Ebben az esetben a **null-hipotézisem** az, hogy az **erme nem cinkelt**. Vagyis a null-hipotézis szerint ugyanakkor a valószínűsége fejre és írást kapni eredményként.

- H1: cinkelt érme (fej fele)
- H0: nem cinkelt érme

Tegyük fel hogy 10-szer feldobjuk az ermet, és 9-szer fejet dobunk. Mekkora a valószínűsége, hogy az érme cinkelt? Ezt nem tudjuk megmondani. Többek között azért sem mert nem tudjuk, mennyire lehet cinkelve. Viszont azt meg tudjuk mondani, hogy mekkora a valószínűsége, hogy ezt az eredményt kapnánk, ha az érme **NINCS** cinkelve.

Annak a valószínűsége, hogy **legalább 9-szer** (vagy többször) fejet dobok **10 dobásból** egy nem cinkelt érmevel, $p = 0.0107$ (**nagyjából 1%**). (Ezt a kódot részben nem fontos megérteni, a lényege hogy a `pbinom()` funkcióval kiszámoltuk a valószínűséget, hogy 10 feldobásból legalább 9 fej lesz).

```

probability_of_heads_if_H0_is_true <- 0.5

heads <- 9
total_flips <- 10
probability_of_result = 1-pbinom(heads-1, total_flips, probability_of_heads_if_H0_is_true)

probability_of_result

## [1] 0.01074219

```

Ez a valószínűség **maskepp mondva** azt jelenti, hogy ha ugyan ezt a kísérletet 100-szor megismételjük (mindegyikben 10 feldobással), akkor a 100 kísérletből csak átlagosan nagyjából 1-szer várható, hogy 9 vagy több fejet kapjunk.

Ezt le is ellenőrizhetjük, ha véletlenszerűen **generálunk 10.000 hasonló kísérletet** az `rbinom()` funkcióval. Az ábrán látható, hogy csak a kísérletek igen kis százalékában kaptunk 9 vagy több "sikert". (Ezt a kódreszt sem fontos megérteni, a lényege, hogy az `rbinom()` funkcióval 10000-szer azt szimuláltuk, hogy egymás után 10-szer feldobtuk egy érmét (vagyis hogy véletlenszerűen választottunk egy számot 0 és 1 között), ez után ennek a 10000 kísérletnek az eredményét ábrázoltuk a `ggplot`-tal, az oszlopok magassága azt jelzi, hogy hány kísérletben jött ki az adott számú siker).

```

successes = rbinom(n = 10000, size = 10, prob = 0.5)
random_flips = data.frame(successes)

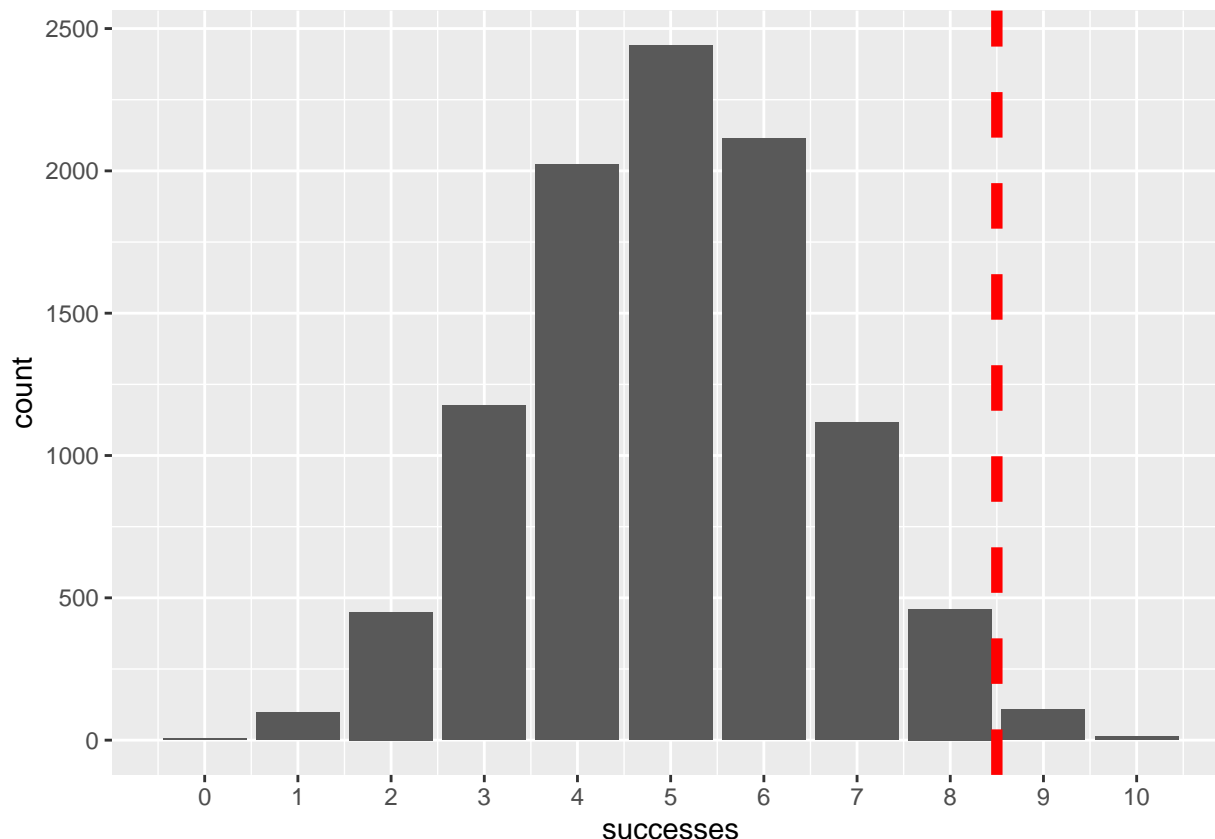
ggplot(data = random_flips) +
  aes(x = successes) +
  geom_bar() +
  scale_x_continuous(breaks = 0:10) +
  geom_vline(xintercept = 8.5, col = "red", linetype = "dashed", size = 2)

```

```

## Warning: Using 'size' aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use 'linewidth' instead.

```



Vagyis a 9 fej 10 feldobasbol egy **eleg meglepo** (nagyon ritka) eredmény, hiszen ez csak az esetek nagyjából 1%-ában fordul elő ha az érme nem cinkelt. De mit mond ez nekünk arról hogy **az érme valóban cinkelve van-e** vagy sem? Mekkora ennek az esélye? Ezt sajnos nem tudjuk meg. Amit megtudunk ebből a számításból, az az, hogy **milyen ritka ez az eredmény amit kaptunk ha azt feltételezzük hogy az érme nincs cinkelve**. Ezt a fajta fordított logikát kell megérteni ahhoz, amikor az NHST-vel dolgozunk.

Tegyük fel hogy egy **“igen-vagy-nem” dontest** kell hoznunk arról, hogy cinkelt-e az érme vagy sem. Mondjuk minket biztattak meg hogy ellenőrizzük az érmet egy fontos pénzfeldobás előtt, és el kell döntünk, hogy megbízunk-e ennek az érmenek a hitelességében, vagy kerjünk egy új érmet a pénzfeldobáshoz, mert ezt cinkeltnek itéljük. Itt jön az NHST **teszt** része. Ezt a dontest az NHST-ben egy előre meghatározott valószínűségi küszöbérték, **dontesi küszöbérték**, figyelembevételével hozzuk meg. Ha az általunk megfigyelt eredmény **kelloen meglepo, kelloen ritka** a null hipotézis helyessége feltételezve, akkor elvetjük azt a feltételezést, hogy a null-hipotézis helyes. Ilyenkor kizárásos alapon az alternatív hipotézis helyessége fogadjuk el.

A pszichológia tudományában a dontesi küszöbérték tradicionalisan 5%, vagyis ha annak a valószínűsége hogy az általunk megfigyelt eredményt (vagy annál extrémabb eredményt) kapjunk a null hipotézis helyessége esetén **kisebb mint 5%** ($p < 0.05$), akkor **elvetjük a null-hipotézist**.

Fontos azonban hangsúlyozni, hogy egy-egy NHTS teszt során nem tudjuk meg a null hipotézis helyességének, vagy az alternatív hipotézis helyességének a valódi valószínűségeit. Csak azt tudjuk, hogy mennyire valószínű vagy valószínűtlen hogy az általunk megfigyelt eredményt látjuk “egy olyan világban” ahol a null hipotézis helyes. Se többet, se kevesebbet. És ez alapján hozzuk meg a döntésünket a null-hipotézis elvetéséről, vagy megtartásáról.

Az NHST módszer fő előnye, hogy ha **konzisztensen használjuk a fent említett dontesi küszöböt** a kutatásainkban, akkor **elegge biztosak** lehetünk abban, hogy a statisztikai döntéseink során **csak 5%-ában vetjük el hibásan a null hipotézist**. Vagyis a statisztikai döntéseknek csak 5%-a lesz hibás, ha

a null hipotézis valóban igaz, így tehát az elsőfajú hiba (alpha-error) valószínűsége 5%. (Masszóval a teszteknek csak 5%-ában állítjuk hibásan, hogy van hatás, amikor valóban nincs hatás.)

Két fontos kitért érdemes megfigyelni a fenti állításban. Egyrészt hogy azt irtam hogy “elege biztosak” lehetünk. Azért csak “elege biztosak” lehetünk ebben, és nem teljesen biztosak, mert ahhoz hogy ez az állítás helyes legyen, az általunk használt statisztikai tesztek **előfeltevéseinek teljesülnie kell**, és ebben nem lehetünk teljesen biztosak a populáció szintjén. A másik, hogy **“ha a null hipotézis valóban igaz”**. Arról az NHST-ben nem kapunk garanciát, hogy a statisztikai döntéseinknek hány százaléka hibás ha az alternatív hipotézis az igaz. Azt is fontos megérteni, hogy ez nem jelenti azt, hogy az összes publikált null hipotézis-tesztelésben csak 5%-nyi lenne az elsőfajú hiba, mert nem minden statisztikai döntést publikálnak.

Statisztikai tesztek

Nem kell jonak lennünk valószínűségszámításból hogy jó statisztikai döntéseket tudjunk hozni. A megfigyeles valószínűsége a null-hipotézis helyessége feltételezve általában egy **statisztikai teszt** mondja meg nekünk. Ezen az órán 5 statisztikai tesztet fogunk megismerni.

- binomialis teszt
- khi-négyszet teszt
- t-teszt
- egyszempontos ANOVA
- korrelációs teszt

binomialis teszt

A hipotézist, hogy az érme cinkelt, tesztelhetjük a **binomialis teszttel**, aminek R-ben `binom.test()` a funkciója. Az `x` helyére a megfigyelt “celmegfigyelesek” vagy “sikerek” számát (a mi esetünkben a fejek számát, $x = 9$), az `n` helyére az összes megfigyeles számát ($n = 10$), a `p` helyére pedig a **null-hipotézis** helyessége feltételezve a “celmegfigyelesek” elérésének valószínűsége kell beírni (mivel a hipotézisünk az hogy az érme cinkelt, az null hipotézisünk az, hogy az érme “nem cinkelt”). Ezt valószínűségként kell megadni, amit egy 0 és 1 közötti számmal jellemezhetünk (ahol a 0 azt jelenti hogy a megfigyelesek 0%-a lesz “siker”, az 1 pedig azt hogy a megfigyelesek 100%-a lesz “siker”, vagyis a 0.6 jelentése hogy a megfigyelesek 60%-a lesz “siker”). A mi esetünkben a null hipotézis helyessége esetén a fej valószínűsége 50% ($p = 0.5$).

```
binom.test(x = 9, n = 10, p = 0.5, alternative = "greater")
```

```
##
## Exact binomial test
##
## data: 9 and 10
## number of successes = 9, number of trials = 10, p-value = 0.01074
## alternative hypothesis: true probability of success is greater than 0.5
## 95 percent confidence interval:
## 0.6058367 1.0000000
## sample estimates:
## probability of success
## 0.9
```

Ennek a tesztnek az eredménye a következőt mutatja:

- p-value: p-érték, annak a valószínűsége, hogy az általunk megfigyelt, vagy extremer eredményt kapunk, feltételezve hogy a null-hipotézis helyes. Általában ha ez az érték 0.05 alatti, akkor elvetjük a null-hipotézist.

- alternative hypothesis: Itt írja le, hogy mi volt a H_1 , ami a mi esetünkben az volt, hogy a fej valószínűsége nagyobb mint 0.5 (50%). Ez egyben azt is jelenti, hogy a null-hipotézisünk az volt, hogy a fej valószínűsége 0.5.
- 95 percent confidence interval (vagy röviden 95% CI): a 95%-os konfidencia intervallum. Ez azt jelenti, hogy ha a kísérletet sokszor megismételjük és ugyan így kiszámoljuk a konfidencia intervallumot minden kísérletnél, az így kapott konfidencia intervallumok 95%-a tartalmazni fogja a valós hatásmeretet (ami a mi esetünkben a “siker”/fej valószínűsége). Fontos, hogy nem tudjuk, hogy a mi konkrét kísérletünkben a konfidencia intervallum tartalmazza-e a valós hatásmeretet.
- sample estimates: A “siker” (“celmegfigyeles”, a mi esetünkben a fej) valószínűségének becsült merteke a populációban a megfigyelt valószínűség alapján. Ez egy pontbecslés, ami mindig megegyezik a megfigyelt valószínűséggel.

Az eredményt így írhatjuk le:

“A kutatásunkban 9 fejet figyeltünk meg 10 penzfeldobasbol (90%). Ez alapján úgy iteljük, hogy annak a valószínűsége, hogy fejet dobunk az érmevel szignifikansan több mint 50%. A fej dobás valószínűsége 0.9 volt a mintában (95% CI = 0.61, 1).”

Adatgeneralas az orahoz

Az alábbi kód **adatokat general** a számunkra. Az adatgeneralashoz használt kód megértése ezen a szinten meg nem szükséges.

```
n_per_group = 40

base_height_mean = 164
base_height_sd = 10
base_anxiety_mean = 18
base_anxiety_sd = 2
resilience_mean = 7
resilience_sd = 2

treatment_effect = - 3
resilience_effect = - 0.8

gender_bias = 0.7
gender_effect = - 1
gender_effect_on_height = 12

treatment <- rep(c(1, 0), each = n_per_group)
set.seed(1)

gender_num <- rbinom(n = n_per_group * 2, size = 1, prob = 0.7)
gender <- NA
gender[gender_num == 0] = "female"
gender[gender_num == 1] = "male"

set.seed(2)
home_ownership <- sample(c("own", "rent", "friend"), n_per_group * 2, replace = T)

set.seed(3)
resilience <- rnorm(mean = resilience_mean, sd = resilience_sd, n = n_per_group*2)
```

```

set.seed(6)
anxiety_base <- rnorm(mean = base_anxiety_mean, sd = base_anxiety_sd, n = n_per_group*2)
anxiety_baseline <- anxiety_base + resilience * resilience_effect + gender_num * gender_effect + rnorm(1)
anxiety_post <- anxiety_base + treatment * treatment_effect + resilience * resilience_effect + gender_num * gender_effect + rnorm(1)
participant_ID <- paste0("ID_", 1:(n_per_group*2))

set.seed(5)
height_base <- rnorm(mean = base_height_mean, sd = base_height_sd, n = n_per_group*2)
height <- height_base + gender_num * gender_effect_on_height

group <- rep(NA, n_per_group*2)
group[treatment == 0] = "control"
group[treatment == 1] = "treatment"

health_status <- rep(NA, n_per_group*2)
health_status[anxiety_post < 11] = "cured"
health_status[anxiety_post >= 11] = "anxious"

data <- data.frame(participant_ID)
data = cbind(data, gender, group, resilience, anxiety_baseline, anxiety_post, health_status, home_ownership)
data = as_tibble(data)

data = data %>%
  mutate(gender = factor(gender))

data = data %>%
  mutate(group = factor(group))

data = data %>%
  mutate(health_status = factor(health_status))

data = data %>%
  mutate(home_ownership = factor(home_ownership),
         anxiety_baseline = round(anxiety_baseline, 2),
         anxiety_post = round(anxiety_post, 2),
         resilience = round(resilience, 2),
         height = round(height, 2))

```

Az adatok egy (kepzeletbeli) randomizalt kontrollalt klinikai kutatas eredményeibol szarmaznak, ahol a **pszichoterapia hatékonyságát** teszteltek. Olyan személyeket vontak be a kutatasba, akik egy **hurrikan áldozatai** voltak, és **szorongással** küszködtek. A személyeknél felmérték a reziliencia (pszichés ellenálló képesség) szintjét, majd véletlenszerűen osztották a személyeket egy kezelési vagy egy kontrol csoportba. Ezt követően a kezelési csoport **pszichoterápiát kapott 6 heten keresztül** heti egyszer, míg a kontrol csoport nem kapott kezelést. A vizsgálat végén megmérték a személyek **szorongásszintjét**, és a klinikai kritériumok alapján meghatározták, hogy a személy **gyógyultnak, vagy szorongónak** számít-e.

Láthatjuk, hogy 8 változó van az adattáblában.

- participant_ID - részvevő azonosítója
- gender - nem
- group - csoporttagság, ez egy faktor változó aminek két szintje van: “treatment” (kezelt csoport), és “control” (kontrol csoport). A “treatment” csoport kapott kezelést, míg a “control” csoport nem kapott kezelést.

- resilience - reziliencia: a nehézségekkel való megküzdés képessége, ez egy személyes képesség, olyasmint mint a személyiségvonások
- anxiety_baseline - szorongás szint a terápia előtt
- anxiety_post - szorongás szint a terápia után
- health_status - a klinikai kritériumok alapján szorongónak vagy gyógyultnak tekinthető a személy
- home_ownership - lakhatási helyzet: három szintje van az alapján hogy a személy hol lakik: “friend” - barát nál vagy családnál lakik, “own” - saját tulajdonú lakásban lakik, “rent” - bérlet lakásban lakik,
- height - magasság

Adatellenőrzés

Mint mindig, elemzés előtt **ellenőrizzuk**, hogy az adattal minden rendben van-e!

```
data
```

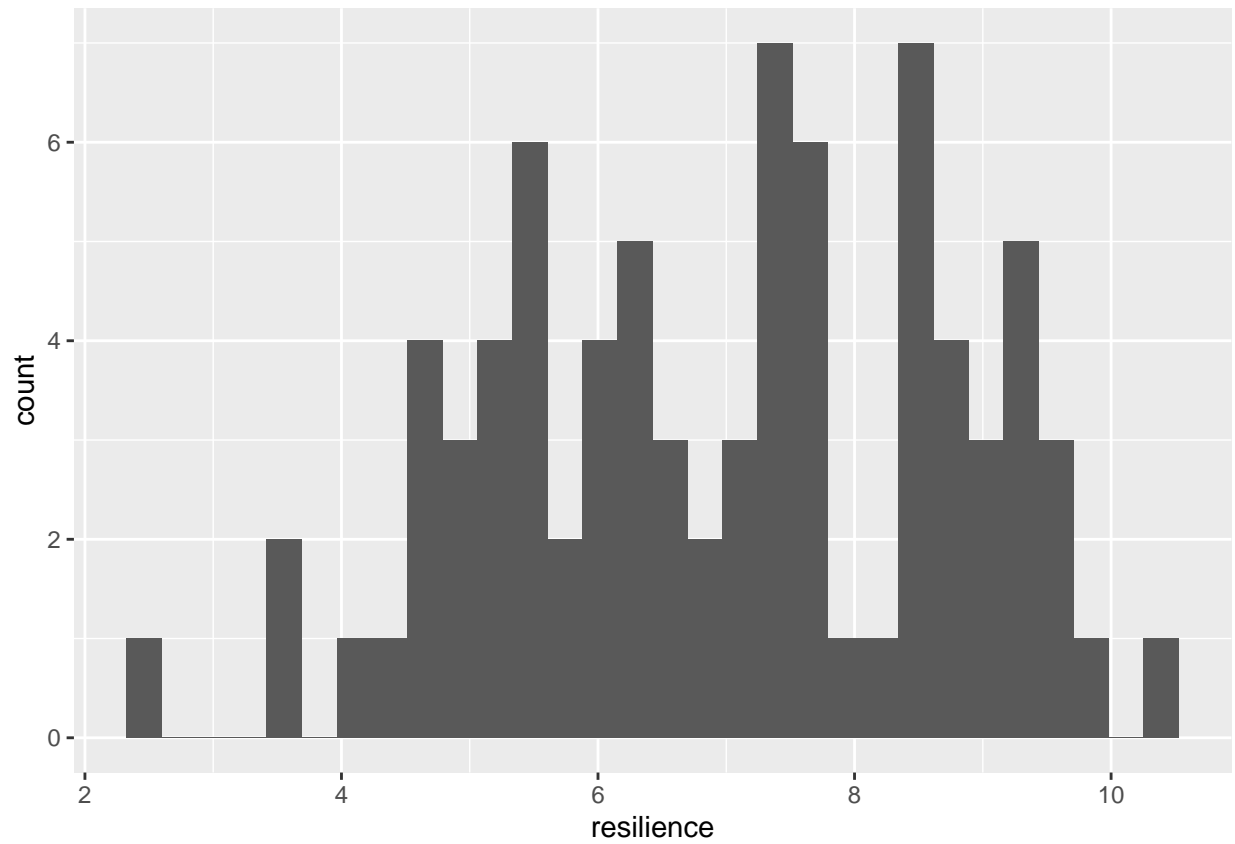
```
## # A tibble: 80 x 9
##   participant_ID gender group   resil~1 anxie~2 anxie~3 healt~4 home_~5 height
##   <chr>          <fct> <fct>   <dbl>   <dbl>   <dbl> <fct>   <fct>   <dbl>
## 1 ID_1          male  treatme~ 5.08    18.7    10.5  cured   own     168.
## 2 ID_2          male  treatme~ 6.41    10.5     7.61  cured   friend  190.
## 3 ID_3          male  treatme~ 7.52    17.2     9.72  cured   rent    163.
## 4 ID_4          female treatme~ 4.7     17.7    14.7  anxious rent    165.
## 5 ID_5          male  treatme~ 7.39     8.04    8.14  cured   own     193.
## 6 ID_6          female treatme~ 7.06    10.3    10.1  cured   own     158.
## 7 ID_7          female treatme~ 7.17    12.6     6.64  cured   own     159.
## 8 ID_8          male  treatme~ 9.23    10.6     8.09  cured   own     170.
## 9 ID_9          male  treatme~ 4.56    12.8    10.4  cured   own     173.
## 10 ID_10         male  treatme~ 9.53     5.82     4.28  cured   rent    177.
## # ... with 70 more rows, and abbreviated variable names 1: resilience,
## # 2: anxiety_baseline, 3: anxiety_post, 4: health_status, 5: home_ownership
```

```
data %>%
  summary()
```

```
## participant_ID      gender      group      resilience
## Length:80          female:25   control :40   Min.   : 2.470
## Class :character   male :55     treatment:40 1st Qu.: 5.518
## Mode  :character                               Median : 7.125
##                                                  Mean   : 6.981
##                                                  3rd Qu.: 8.477
##                                                  Max.   :10.400
## anxiety_baseline  anxiety_post  health_status home_ownership  height
## Min.   : 4.650    Min.   : 3.910  anxious:32    friend:22    Min.   :142.2
## 1st Qu.: 9.668    1st Qu.: 8.223  cured :48     own :31      1st Qu.:163.4
## Median :11.155    Median :10.110                rent :27      Median :173.0
## Mean   :11.393    Mean   :10.212                               Mean   :172.3
## 3rd Qu.:12.730    3rd Qu.:12.255                               3rd Qu.:179.7
## Max.   :19.320    Max.   :16.710                               Max.   :198.2
```

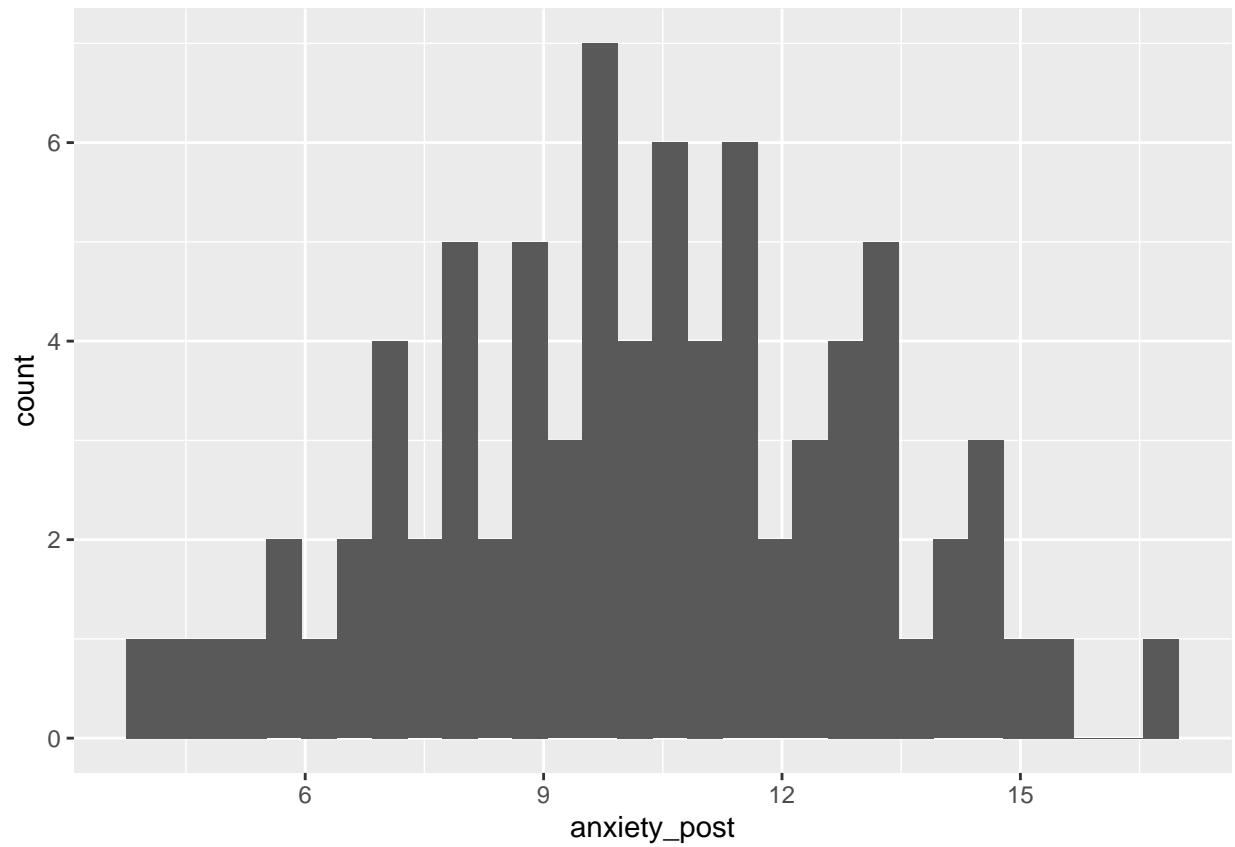
```
data %>%
  ggplot() +
    aes(x = resilience) +
    geom_histogram()
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```

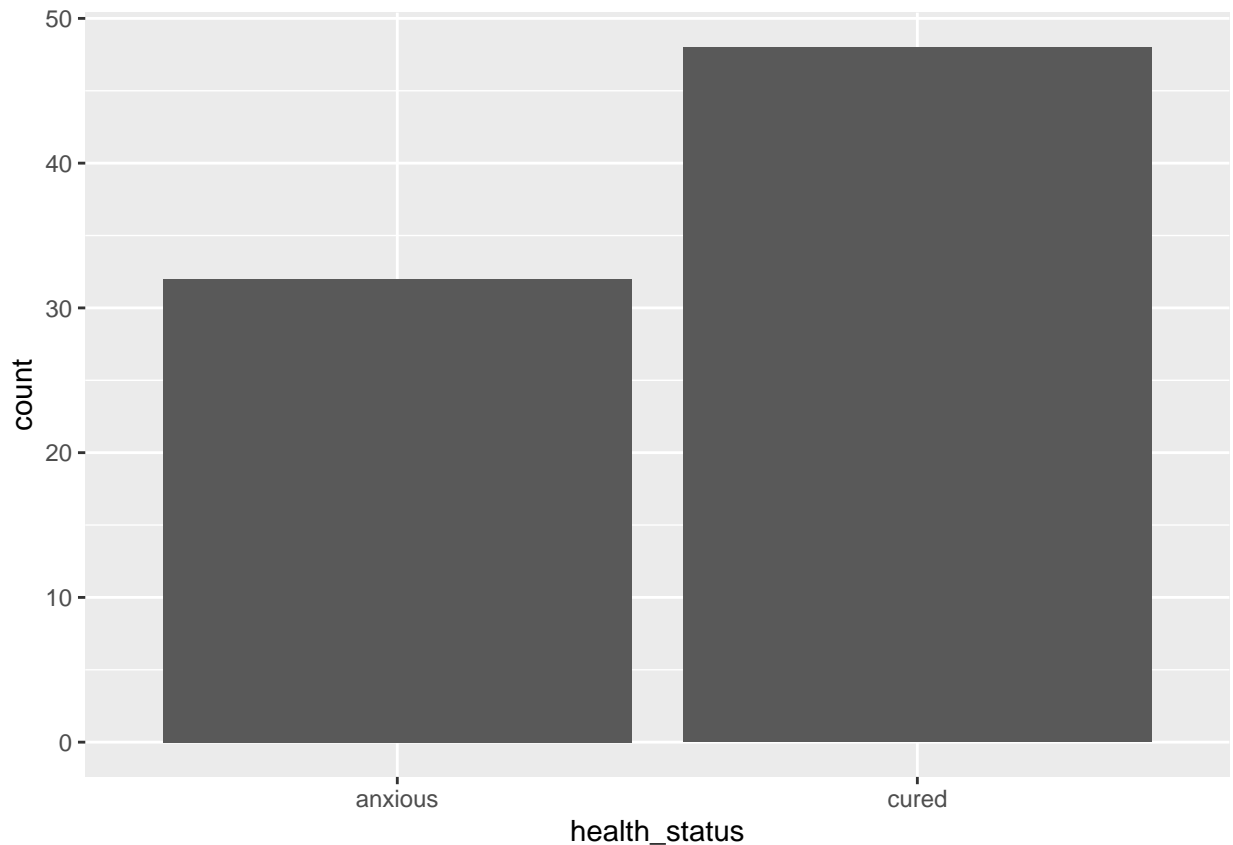


```
data %>%  
  ggplot() +  
    aes(x = anxiety_post) +  
    geom_histogram()
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```

```
data %>%  
  ggplot() +  
    aes(x = health_status) +  
    geom_bar()
```



```
set.seed(Sys.time())
```

Hipotezisek

Vizsgáljuk meg kutatásban szereplő változók összefüggését a hipotezisek mentén.

A kutatás hipotezise a következők voltak:

1. Több a férfi mint a nő ebben a klinikai mintában (**gender** vs. 50%).
2. A pszichoterápiát kapó csoportban a terápia után kevesebb lesz a klinikai kritériumok alapján szorongónak számító személy (**health_status** vs. **group**)
3. A terápiás csoportban alacsonyabb lesz a szorongás átlaga a kutatás végére mint a kontrol csoportban (**anxiety_post** vs. **group**)
4. A reziliencia és a kutatás végén mért szorongásszint negatív összefüggést fog mutatni (vagyis aki reziliensebb, annál alacsonyabb szorongásszintet fognak mérni a kutatás végén) (**anxiety_post** vs. **resilience**)

Gyakorlás

Teszteld a hipotézist, hogy “Több a férfi mint a nő ebben a klinikai mintában” (**gender** változó)

- Ezt ugyan úgy teheted meg, mint a fenti példában, hiszen a null-hipotézis az, hogy a férfiak (“male”) elvárt valószínűsége 50% vagy kevesebb ($p = 0.5$). Szóval a férfiak ekvivalensek a “fejekkel” a pénzfeldobásos példában.

- Meg kell határozni a férfiak számát a mintában, és a teljes mintaelemszámot, hogy ki tudja tölteni a `binom.test()` függvény paramtereit.
- Ez után végezd el a tesztet
- Es ird le a fentiek szerint az eredményeket.

Két kategorikus változó kapcsolata: Khi-negyzet próba (Chi-squared test)

Két kategorikus változó kapcsolatának vizsgálatára a **Khi-negyzet próba** javasolt.

Peldául megvizsgálhatjuk, hogy van-e kapcsolat abban, hogy a személyek lakhatási helyzete (**home_ownership**) és a között, hogy a kutatás végen az egyes személyek meggyógyultak-e (**health_status**).

A Khi-negyzet próba előfeltetelei:

- Minden megfigyeles független a többi megfigyelestől (pl. egy megfigyeles személyenként)
- A kategória-kombinációk ábrázolásával kapott táblázatban nem több mint a cellák 20%-ában kisebb a várható érték 5-nél, és minden cellában magasabb a várható érték mint 1.

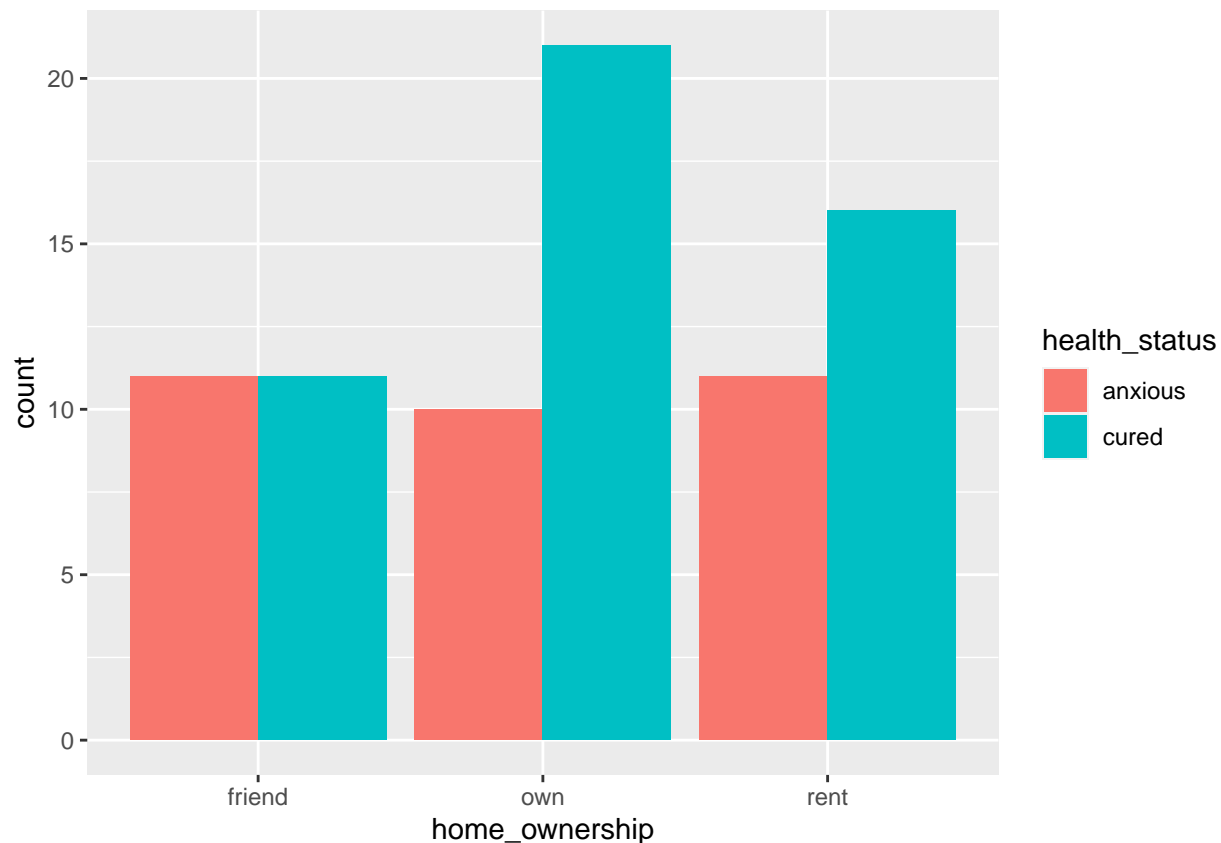
Eloszor feltáro elemzést végzünk:

- táblázatot rajzolunk a két változó kapcsolatáról
- ábrát készítünk (pl. `geom_bar`)

```
table(data$home_ownership, data$health_status)
```

```
##
##           anxious cured
## friend           11    11
## own              10    21
## rent             11    16
```

```
data %>%
  ggplot() +
    aes(x = home_ownership, fill = health_status) +
    geom_bar(position = "dodge")
```



Ez után elvegezzük a Khi-negyzet probát. Ehhez először készítenünk kell egy **tablazatot a két változó kapcsolatáról**, amit egy új objektumban elmentünk.

A Khi-negyzet próba azt a **null-hipotézist** teszteli, hogy **a csoportokban ugyan olyan a másik kategorikus változó eloszlása** (vagyis a mi esetünkben a null hipotézis hogy ugyan olyan arányban gyógyulnak meg akik barát nál laknak, akiknek saját lakasuk van, és akik berlik a lakást).

(Mivel itt egy 3x2-es tablázatunk van, a Khi-negyzet próba helyett a Fisher exact tesztet érdemes használni, de az alábbi kódban megtalálod a khi negyzet próba kódját is.)

```
ownership_health_status_table = table(data$home_ownership, data$health_status)
ownership_health_status_table
```

```
##
##           anxious cured
## friend         11    11
## own            10    21
## rent           11    16
```

```
chisq.test(ownership_health_status_table)
```

```
##
## Pearson's Chi-squared test
##
## data:  ownership_health_status_table
## X-squared = 1.697, df = 2, p-value = 0.428
```

A teszt eredményét így írhatjuk le: “Nem volt szignifikáns eltérés abban, hogy a különbozó lakhatási csoportokban (barátnál, saját lakásban, vagy berlemben lakók) milyen arányban voltak azok akik meggyógyultak a kutatás végére ($X^2 = 1.70$, $df = 2$, $p = 0.428$).”

Fentebb láthattuk hogy a Khi-negyzet próba alkalmazásának egyik feltétele, hogy ne legyen a cellák várható értéke 5-nél kisebb. A cellák várható értéket úgy lehet kiszámolni, hogy a cella sorának számainak az összeget megszorozzuk a cella oszlopának számainak az összegevel, majd ezt elosztjuk a táblázat összes számának összegevel. Ezt szerencsére nem kell kézzel kiszámolnunk minden cellára, mert a `chisq.test()` függvény kiszámolja nekünk. Ezt az információt a `chisq.test()` függvény eredményének az `$expected` elemében találjuk meg, az alábbi kód bemutat egy példát:

```
chi = chisq.test(ownership_health_status_table)
chi$expected
```

```
##
##           anxious cured
## friend      8.8  13.2
## own        12.4  18.6
## rent       10.8  16.2
```

Ha ebben a várható értékeket tartalmazó táblázatban a számok több mint 20%-a kisebb mint 5, vagy ha bármelyik kisebb mint 1, akkor a Khi-negyzet teszt helyett a Fisher tesztet kell használni. Mivel a fenti példában a táblázat csak 6 számot tartalmaz, annak a 20%-a 1.2, vagyis ha akár egy szám is 5 alatti lenne, a Fisher tesztet kellene használni. A példánkban a számok nem 5 alattiak, ezért a Khi-negyzet teszt eredménye a mervado, de az alábbi kód mutat egy példát arra, hogyan kellene a Fisher tesztet használni, ha a Khi-negyzet teszt előfeltétele nem teljesülne.

```
fisher.test(ownership_health_status_table)
```

```
##
## Fisher's Exact Test for Count Data
##
## data:  ownership_health_status_table
## p-value = 0.424
## alternative hypothesis: two.sided
```

A Fisher exact teszt eredményét így írhatjuk le:

“Nem volt szignifikáns eltérés abban, hogy a különbozó lakhatási csoportokban (barátnál, saját lakásban, vagy berlemben lakók) milyen arányban voltak azok akik meggyógyultak a kutatás végére (Fisher exact $p = 0.424$).”

Gyakorlás

Teszteld a 2. hipotézist, hogy “A pszichoterápiát kapó csoportban a terápia után kevesebb lesz a klinikai kritériumok alapján szorongónak számító személy” (**health_status** vs. **group**)

- Ezt ugyan úgy teheted meg, mint a fenti példában, hiszen a null-hipotézis az, hogy nincs különbség a csoporttagság szerint (treatment vs. control) abban hogy milyen arányban gyógyultak meg a kutatás végére.
- Eloszor vegezzünk egy feltárási elemzést egy táblázattal a két változó kapcsolatáról a `table()` funkcióval és egy ábrával (mondjuk `geom_bar()` használatával)

- A tablazatot mentsd el egy új objektumba
- Ez után vegezd el a tesztet, `chisq.test()`
- Es ird le a fentiek szerint az eredményeket.

Egy numerikus változó átlagának különbsége csoportok között: anova és t-teszt

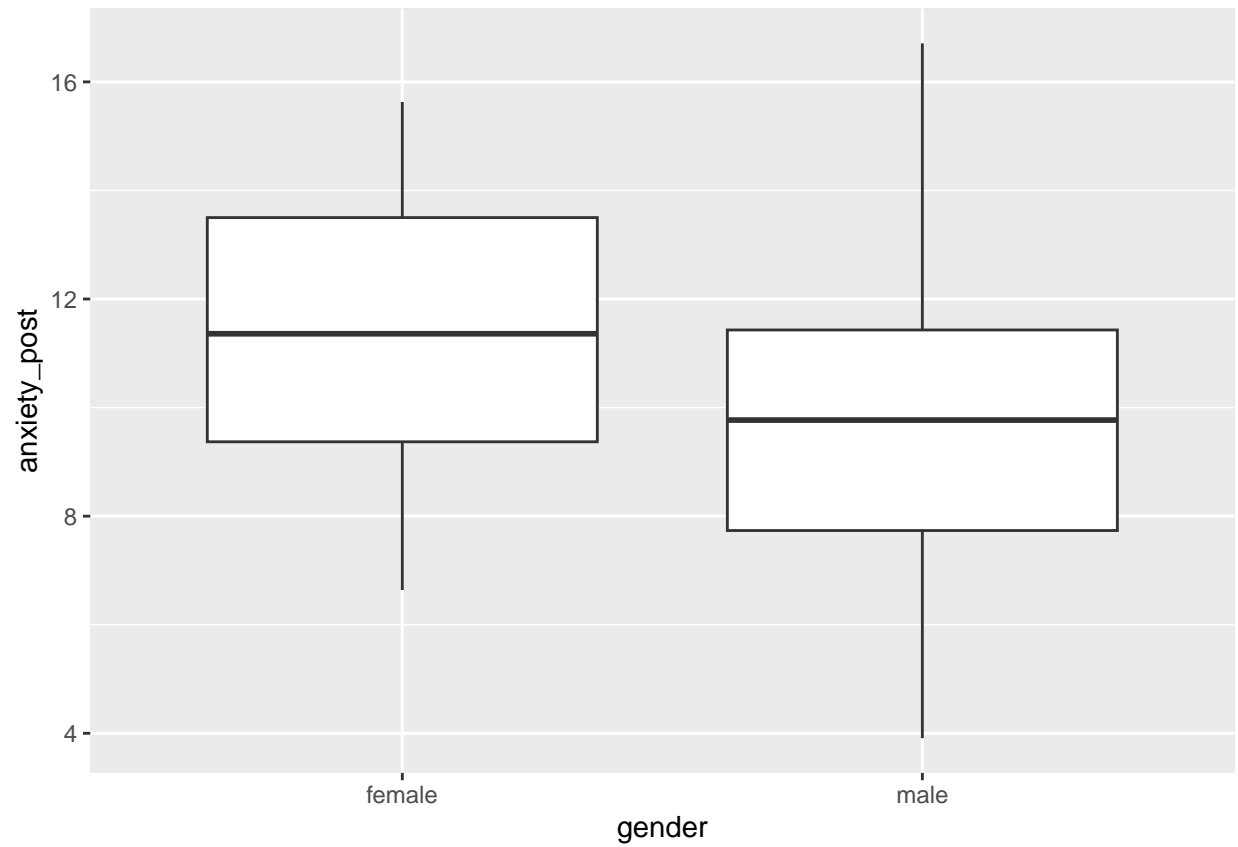
Tesztelhetjük például, hogy van-e különbség a nemek között (**gender**) a kutatás végén mért szorongás szintjében (**anxiety_post**).

Eloszor szokás szerint feltáro elemzést végzünk átlagok csoportonkénti összehasonlításával és ábrával. Erre pl. remek a `geom_boxplot()` és a `geom_density()`

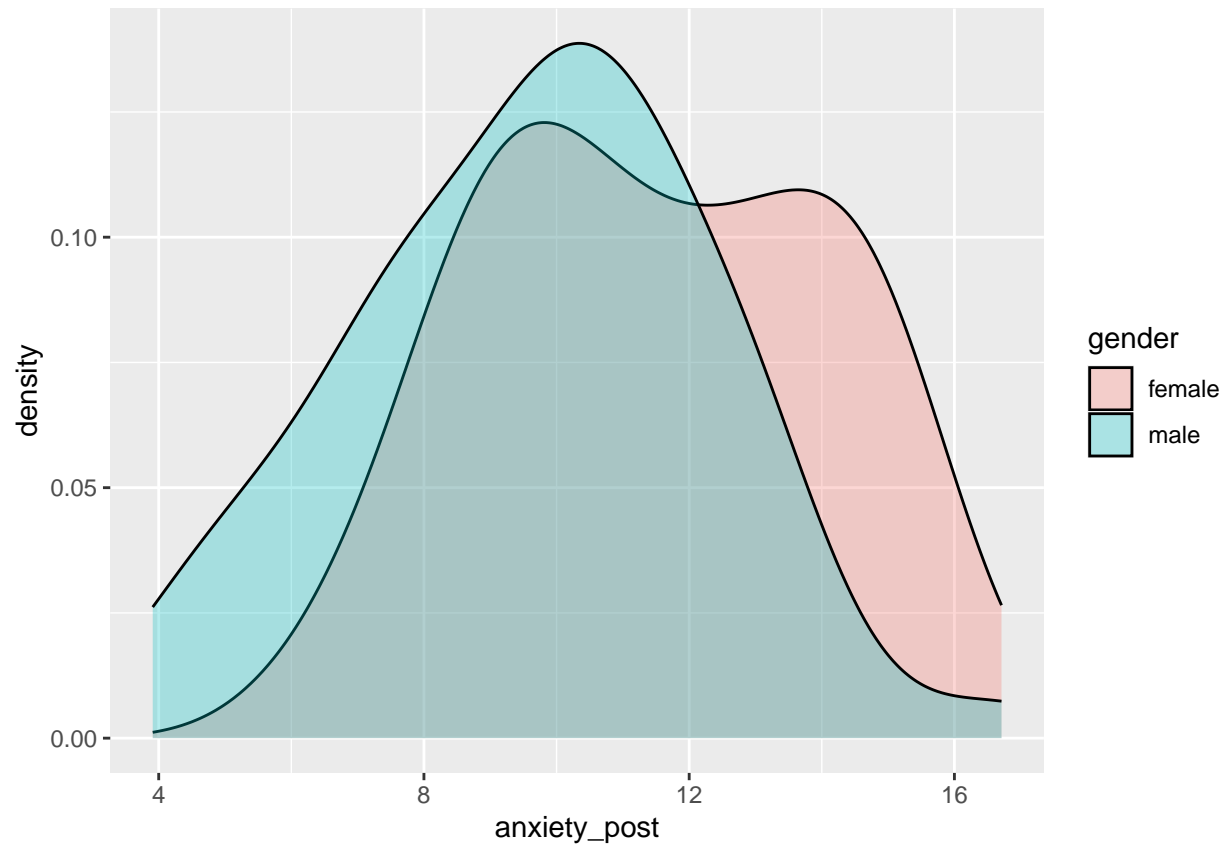
```
summary = data %>%
  group_by(gender) %>%
    summarize(mean = mean(anxiety_post), sd = sd(anxiety_post))
summary
```

```
## # A tibble: 2 x 3
##   gender mean    sd
##   <fct> <dbl> <dbl>
## 1 female 11.5    2.58
## 2 male   9.64    2.70
```

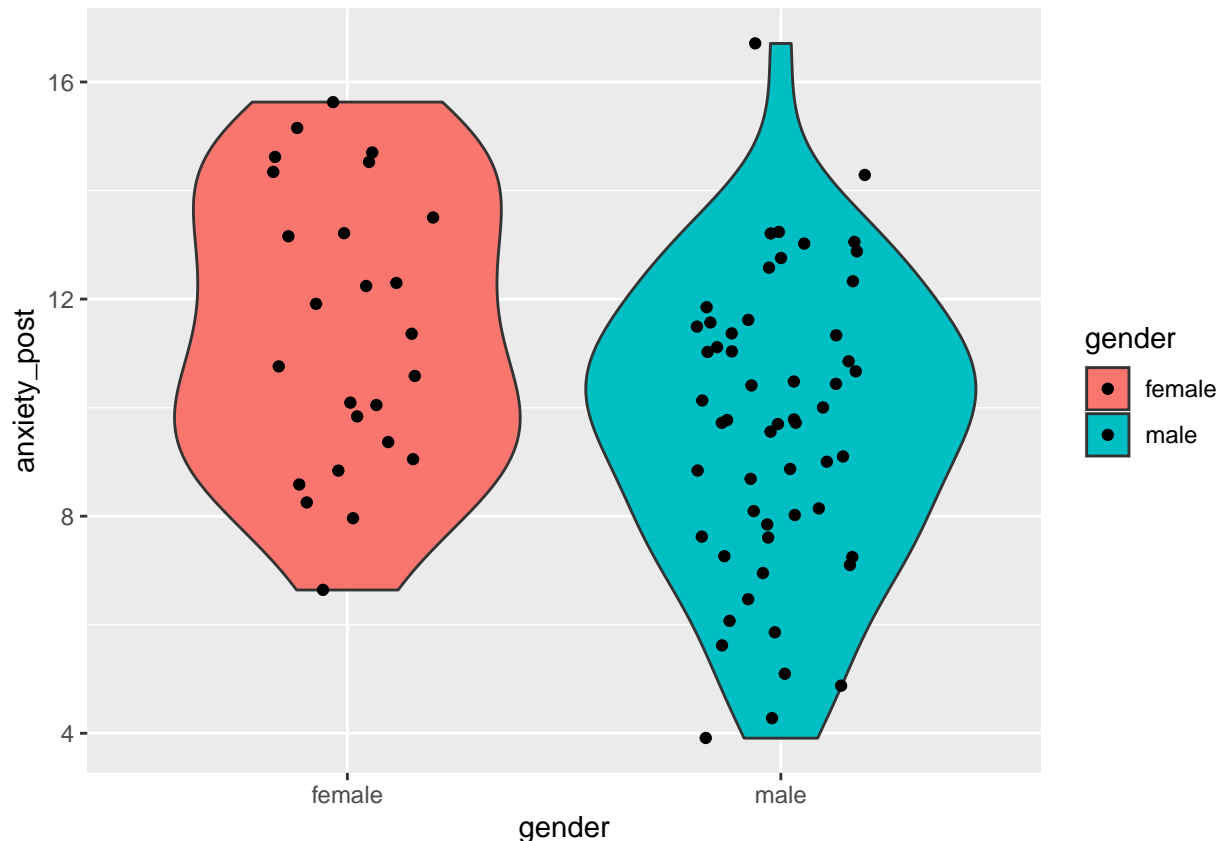
```
data %>%
  ggplot() +
    aes(x = gender, y = anxiety_post) +
    geom_boxplot()
```



```
data %>%  
  ggplot() +  
    aes(x = anxiety_post, fill = gender) +  
    geom_density(alpha = 0.3)
```



```
data %>%  
  ggplot() +  
    aes(x = gender, y = anxiety_post, fill = gender) +  
    geom_violin() +  
    geom_jitter(width = 0.2)
```

Lathatjuk a feltaro elemzes alapján, hogy a nok szorongasszintje nagyobb valamivel mint a ferfiaké atlagosan. Most nezzuk meg, ez a kulonbseg statisztikailag szignifikans-e.

Fuggetlen mintas t-teszt

Arra, hogy meghatarozzuk van-e kulonbseg ket csoport kozott valamilyen numerikus valtozo atlagaban, hasznalhatjuk a fuggetlen mintas **t-tesztet**, `t.test()`.

A t-teszt elofeltetelei:

- A fuggo valtozo intervallum vagy aranyiskalan mozog
- A fuggetlen valtozo ket egymastol fuggetlen kategorikus csoportot reprezental
- A megfigyelesek fuggetlenek egymastol. Minden megfigyeles csak az egyik csoportba sorolható, és a csoportok között nincs összefüggés az egyes megfigyelesek között.
- Nincsenek jelentős kiugró esetek.
- Csoportonként normalis eloszlást mutat a fuggo valtozo eloszlása.
- Variancia-homogenitas: a fuggo valtozo varianciaja azonos a ket csoportban. (A welch t-teszt-et lehet alkalmazni, ha ez a feltétel sérül).

```
t_test_results = t.test(anxiety_post ~ gender, data = data)
t_test_results
```

```
##
##  Welch Two Sample t-test
##
```

```
## data: anxiety_post by gender
## t = 2.8895, df = 48.518, p-value = 0.005754
## alternative hypothesis: true difference in means between group female and group male is not equal to
## 95 percent confidence interval:
## 0.5553479 3.0941793
## sample estimates:
## mean in group female    mean in group male
##          11.466400          9.641636
```

```
mean_dif = summary %>%
  summarize(mean_dif = mean[1] - mean[2])
mean_dif
```

```
## # A tibble: 1 x 1
##   mean_dif
##   <dbl>
## 1      1.82
```

Az eredményt így írhatjuk le:

“A férfiak és nők szignifikánsan különböztek a egymástól a szorongás szintjükben ($t = 2.89$, $df = 48.52$, $p = 0.006$). A csoportok szorongás szintjének átlaga és szórása a következő volt:”nők: 11.47 (2.58), férfiak: 9.64 (2.70). A nők átlagosan 1.82 ponttal voltak szorongóbbak (95% CI = 0.56, 3.09).”

Egyszempontos ANOVA

Ha egy kategorikus változon belül **három vagy több csoportunk** is van, a t-test nem használható. Helyette használhatjuk az **egyszempontos ANOVA**-t (one-way ANOVA) az aov() funkcióval. A formula ugyan úgy néz ki, mint a t-teszt esetén.

Az egyszempontos ANOVA előfeltételei majdnem ugyan azok, mint a független mintás t-tesztei:

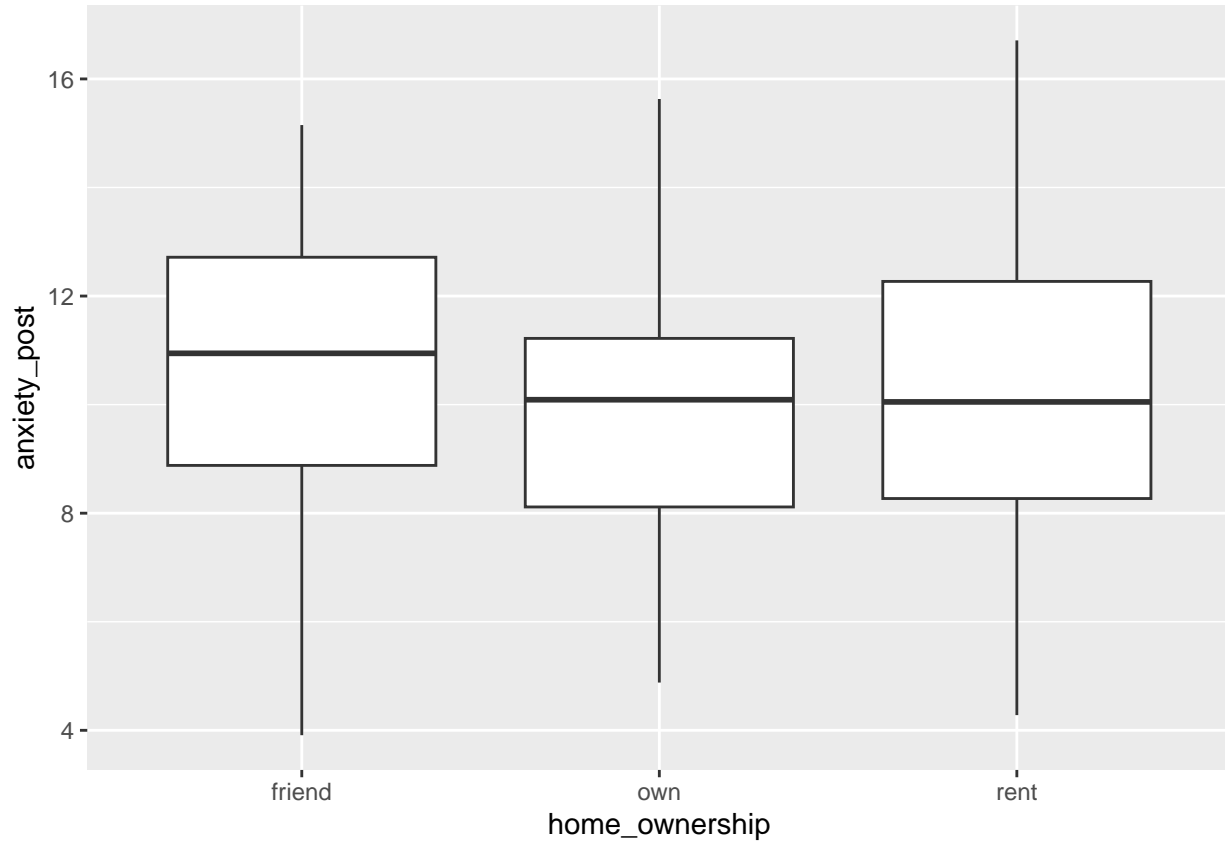
- A függő változó intervallum vagy arányskálán mozog
- A független változó két vagy több egymástól független kategorikus csoportot reprezentál
- A megfigyelések függetlenek egymástól. Minden megfigyelés csak az egyik csoportba sorolható, és a csoportok között nincs összefüggés az egyes megfigyelések között.
- Nincsenek jelentős kiugró esetek.
- Csoportonként normalis eloszlást mutat a függő változó eloszlása.
- Variancia-homogenitás: a függő változó varianciája azonos a csoportokban.

Igy teszteljük hogy van-e különbség a **lakhatási helyzet csoportjai** között a **szorongásszintben**.

```
summary_home_ownership_vs_anxiety_post = data %>%
  group_by(home_ownership) %>%
  summarize(mean = mean(anxiety_post), sd = sd(anxiety_post))
summary_home_ownership_vs_anxiety_post
```

```
## # A tibble: 3 x 3
##   home_ownership mean    sd
##   <fct>          <dbl> <dbl>
## 1 friend         10.6   2.93
## 2 own            9.86   2.57
## 3 rent          10.3   2.94
```

```
data %>%
  ggplot() +
    aes(x = home_ownership, y = anxiety_post) +
    geom_boxplot()
```



```
ANOVA_result = aov(anxiety_post ~ home_ownership, data = data)
summary(ANOVA_result)
```

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## home_ownership  2    7.5   3.760    0.48  0.621
## Residuals      77  603.3   7.835
```

Az eredményt így írhatjuk le:

“A lakhatási csoportonk szerint nem volt szignifikáns különbség a szorongás átlagos szintjében ($F(2, 77) = 0.48$, $p = 0.621$). A szorongás átlagát és szórását az egyes csoportok szerinti bontásban lásd az 1. táblázatban”

Alább látható, hogyan produkálnak a megfelelő táblázatot a szorongás átlagával home_ownership csoportok szerint.

Egyoldalu vs. ketoldalu tesztek

Fontos, hogy ha van előzetes elképzelésünk a hipotézisalkotáskor arról, hogy **milyen irányu** lesz a hatás, akkor **egy-oldalu (one-sided) tesztet** kell használnunk az alapértelmezett ket-oldalu teszt helyett.

Peldaul tegyük fel hogy amikor a hipotézisünket meghatároztuk (ideális esetben ez meg az adatgyűjtés előtt megtörténik), úgy gondoltuk, hogy a nőknek magasabb lesz a szorongásszintjük, mint a férfiaknak. Ezt az `alternative = "greater"` paraméterrel határozhatjuk meg.

Ha összehasonlítjuk ezt az eredményt a korábbi t-teszt eredményével, észrevehetjük hogy minden szám változatlan maradt, kivéve a **p-értéket**, ami pontosan felére csökkent, és a 95%-os **konfidencia intervallumot**, aminek a felső határa most egy végtelen nagy szám (`Inf`).

A p-érték azért feleződött meg, mert azzal, hogy meghatároztuk, melyik irányban fog a két csoport különbözni egymástól fele akkora lett az esélye hogy a most megfigyelt, vagy annál nagyobb különbséget kapunk a null-hipotézis helyesseget feltételezve. Vagyis amikor tudjuk, milyen irányu hatást várunk el, mindig érdemes egy-oldalu tesztet alkalmazni, mert ezzel nő a statisztikai erőnk a hatás kimutatására.

Az egyoldalu tesztek esetén amikor az a hipotézisünk, hogy a **referencia-csoport** átlaga magasabb lesz, (`alternative = "greater"`), akkor a konfidencia intervallumnak csak az alsó hatarat számoljuk ki. Ezért írja a teszt eredménye hogy a 95% CI 1.11, `Inf`, vagyis felfelé a végtelenségig tart a konfidencia intervallum.

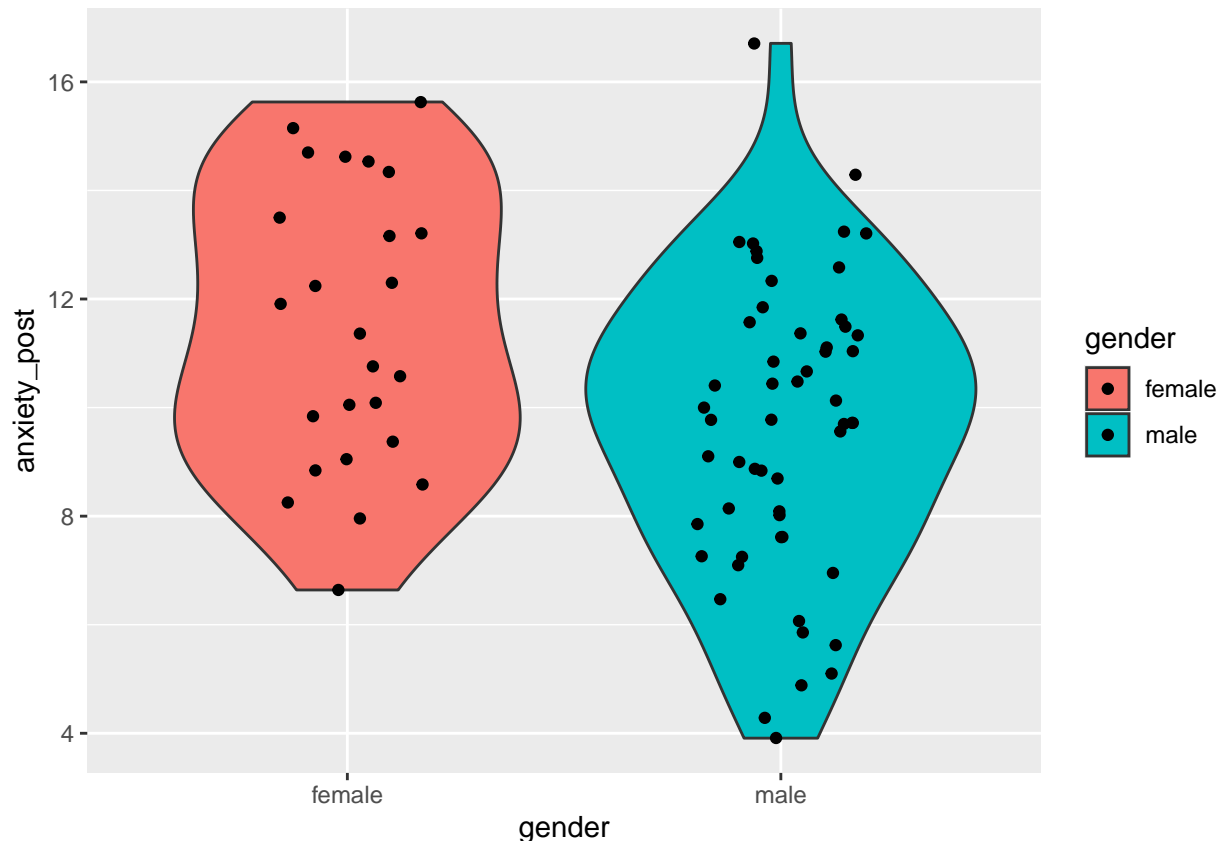
Fontos megjegyezni, hogy amikor azt írjuk a tesztben hogy **alternative = "greater"**, ez alatt azt értjük hogy az alternatív hipotézisünk az, hogy a **referencia-csoport** átlaga magasabb lesz. Ha az alternatív hipotézisünk az lenne hogy a referencia-csoport átlaga alacsonyabb lesz, akkor azt kellene írunk: **alternative = "less"**.

Ahogy korábban már volt róla szó, a referencia-csoport (vagy referencia-szint egy faktor változóban) alapértelmezett módon a faktorszintek névenek ABC sorrendje alapján dolgozik, az ABC sorrendben az első faktorszint lesz a referencia-szint. A példánkban a gender változóban a "female" a referencia-szint, mert az ABC sorrendben a "male" előtt van. Azt, hogy mi legyen a referencia-szint a korábban tanultak szerint a **factor()** **funkcióban a levels =** paraméter beállításával lehet befolyásolni. Nagyon fontos, hogy amikor kategorikus/csoportosított változókkal dolgozunk, mindig tudjuk, mi a referencia-szint.

```
summary = data %>%
  group_by(gender) %>%
  summarize(mean = mean(anxiety_post), sd = sd(anxiety_post))
summary
```

```
## # A tibble: 2 x 3
##   gender mean    sd
##   <fct> <dbl> <dbl>
## 1 female 11.5    2.58
## 2 male   9.64    2.70
```

```
data %>%
  ggplot() +
  aes(x = gender, y = anxiety_post, fill = gender) +
  geom_violin() +
  geom_jitter(width = 0.2)
```



```
t_test_results_one_sided = t.test(anxiety_post ~ gender, data = data, alternative = "greater")
t_test_results_one_sided
```

```
##
## Welch Two Sample t-test
##
## data: anxiety_post by gender
## t = 2.8895, df = 48.518, p-value = 0.002877
## alternative hypothesis: true difference in means between group female and group male is greater than
## 95 percent confidence interval:
##  0.7657774      Inf
## sample estimates:
## mean in group female    mean in group male
##           11.466400           9.641636
```

Az eredményt így írhatjuk le:

“A férfiak és nők szignifikánsan különböztek a egymástól a szorongás szintjükben ($t = 2.89$, $df = 48.52$, $p = 0.003$). A csoportok szorongás szintjének átlaga és szórása a következő volt:

A csoportok szorongás szintjének átlaga és szórása a következő volt: “nők: 11.47 (2.58), férfiak: 9.64 (2.70). A nők átlagosan 1.82 ponttal voltak szorongóbbak (95% CI = 0.77, inf).”

Nezzük meg, mi történne, ha azt tippeltük volna a hipotézisalkotáskor, hogy a nőknek alacsonyabb lesz a szorongásszintjük. Ezt úgy határozhatjuk meg, hogy a `t.test()` funkcióban `alternative = “less”` paramétert állítunk be.

A p-érték itt majdnem eléri az 1-et, vagyis nagyon nagy a valószínűsége, hogy a null-hipotézis helyesreget feltételezve ilyen, vagy ennél extremer különbséget figyelünk meg. Nem is csoda, hiszen a null hipotézisünk itt az hogy a nők szorongásának átlaga nem fog különbözni, vagy nagyobb lesz mint a férfiaké, és azt tapasztaltuk, hogy valóban nagyobb volt, vagyis a megfigyelés egyáltalán nem segít abban, hogy elutasítsuk a null-hipotézist.

```
t_test_results_one_sided = t.test(anxiety_post ~ gender, data = data, alternative = "less")
t_test_results_one_sided

##
## Welch Two Sample t-test
##
## data: anxiety_post by gender
## t = 2.8895, df = 48.518, p-value = 0.9971
## alternative hypothesis: true difference in means between group female and group male is less than 0
## 95 percent confidence interval:
##      -Inf 2.88375
## sample estimates:
## mean in group female    mean in group male
##           11.466400           9.641636
```

Azt is érdemes megjegyezni, hogy a “greater” és a “less” mind a kategóriás változó **referencia-szintjére** vonatkozik. Ha ezt nem állítottuk be maskepp, pl. a **factor (levels =)** funkcióval, akkor a referencia-szint az ABC sorrendben előbb levő szint lesz. A fenti esetben a két szint a “female” és a “male”, amik közül a “female” jön előbb ABC sorrendben. Ha azt tippeltük volna, hogy az lenne a hipotézisünk, hogy a férfiak (“male”) szorongásszintje lesz magasabb, akkor **alternative = “less”**-t kellene beállítanunk, mert ezzel egyben azt tippeljük, hogy a referenciaszint (“female”) átlaga lesz az alacsonyabb. Vagy át kellene állítani a referenciaszintet.

Gyakorlás

Teszteld a 3. hipotézist, hogy “A terápiás csoportban alacsonyabb lesz a szorongás átlaga a kutatás végére mint a kontrol csoportban” (**anxiety_post** vs. **group**)

- Eloszor vegezzünk egy feltároló elemzést egy táblázattal a két változó kapcsolatáról a `summarize(mean(), sd())` funkciókkal, és készítsünk ábrát, mondjuk `geom_boxplot()` segítségével.
- egy- vagy kétoldali tesztet kell alkalmaznunk? (gondolj arra, hogy a hipotézisünkben megjelöljük-e a hatás vagy különbség irányát vagy sem)
- Mi a null-hipotézis ebben az esetben?
- Melyik tesztet érdemes használni, az egyváltozós ANOVA-t, vagy a t-tesztet? (gondolj arra, hogy hány csoport (szint) van a kategóriás változón belül)
- Ez után vegezd el a tesztet
- Es írd le a fentiek szerint az eredményeket.

Két numerikus változó közötti kapcsolat, korreláció, `cor.test()`

Vizsgáljuk meg, van-e együttjárás a reziliencia (**resilience**) és a magasság (**height**) között.

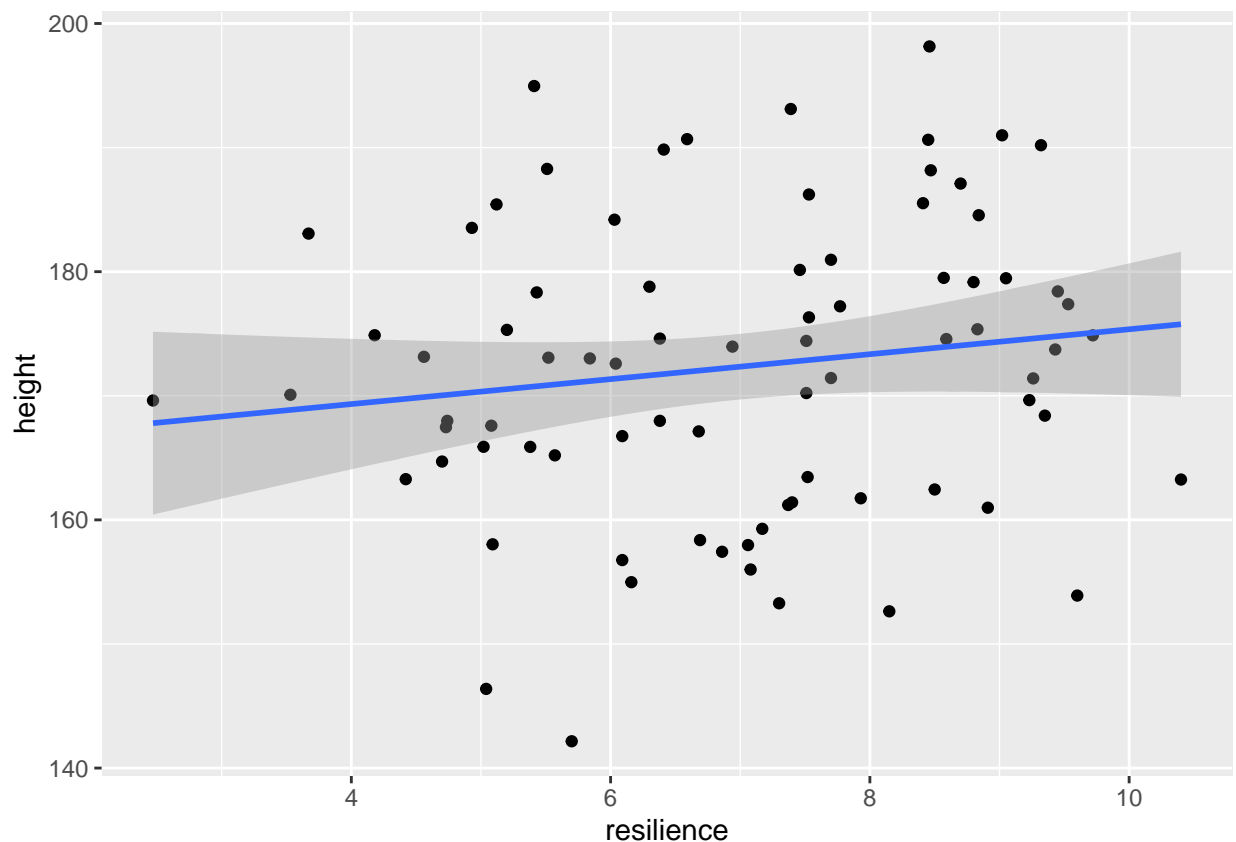
Eloszor vegezzünk feltároló elemzést a korrelációs együttható kiszámításával, és egy pontdiagrammal. Használjunk `geom_point()` és `geom_smooth()` geomokat egyszerre, és használjuk az “lm” módszert a trendvonal megjelölésére.

```
data %>%
  select(resilience, height) %>%
  cor()
```

```
##           resilience  height
## resilience  1.000000 0.146929
## height      0.146929 1.000000
```

```
data %>%
  ggplot() +
    aes(x = resilience, y = height) +
    geom_point() +
    geom_smooth(method = "lm")
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```



A két változó függetlennek tűnik egymástól a feltárolt elemzés alapján, de elképzelhető, hogy a hatás, bármilyen kicsi is, mégis statisztikailag szignifikáns, szóval végezzük el a statisztikai tesztet is.

Ezt a **pearson korrelációs teszt** segítségével tehetjük meg.

A Pearson korrelációs teszt előfeltételei:

- Két folytonos skálájú változó. Ha bármelyik változó ordinalis skálájú, akkor a Spearman korrelációt lehet használni.

- Minden megfigyelesi egységhez két érték tartalmazzon.
- Nincsenek jelentős kiugró értékek
- Linearitás. A két változó kapcsolata egy egyenes vonallal jellemezhető.
- Normalitás: mindkét változó normális eloszlást mutat. Nem normális eloszlás esetén a Spearman korreláció használható.

A tesztet a `cor.test()` funkcióval végezhetjük el a következő képpen:

```
correlation_result = cor.test(data$resilience, data$height)
correlation_result
```

```
##
## Pearson's product-moment correlation
##
## data: data$resilience and data$height
## t = 1.3119, df = 78, p-value = 0.1934
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.07521612 0.35517967
## sample estimates:
## cor
## 0.146929
```

Az eredményt így írhatjuk le:

“A reziliencia és a magasság között nem találtunk szignifikáns együttjárást ($r = 0.15$, 95% CI = -0.08, 0.36, $df = 78$, $p = 0.193$)”

Hasonlóan a t-teszthez, a korrelációs teszt esetében is érdemes egyoldalu tesztet használni amikor a hipotézisünk megmondja a kapcsolat irányát is, nem csak azt, hogy van kapcsolat a két változó között.

Például feltételezzük, hogy a két változó közötti **kapcsolat pozitív irányú** lesz. Vagyis egy ember minél magasabb, annál magasabb a rezilienciája. Ezt úgy adhatjuk meg a statisztikai teszt specifikációjakor, hogy a formulához hozzátesszük az **alternative = “greater”** paramétert. Ha az eredményt összehasonlítjuk az elozo korrelációs teszt eredményével, láthatjuk, hogy a p-érték is megváltozott. A konfidencia intervallumnak itt is csak az alsó határa érdekes, a felső határa a lehető legmagasabb értéket veszi fel ilyenkor, ami a korrelácional 1.

```
correlation_result_greater = cor.test(data$resilience, data$height, alternative = "greater")
correlation_result_greater
```

```
##
## Pearson's product-moment correlation
##
## data: data$resilience and data$height
## t = 1.3119, df = 78, p-value = 0.0967
## alternative hypothesis: true correlation is greater than 0
## 95 percent confidence interval:
## -0.03942784 1.00000000
## sample estimates:
## cor
## 0.146929
```

Gyakorlas

Teszteld a 4. hipotézist, hogy “A reziliencia és a kutatás végen mért szorongásszint negatív összefüggést fog mutatni (vagyis aki reziliensebb, annál alacsonyabb szorongásszintet fognak mérni a kutatás végen)” (**anxiety_post** vs. **resilience**)

- Először vegezzünk egy feltárási elemzést a korrelációs együttható meghatározásával és egy pontdiagrammal a két változó kapcsolatáról.
- egy- vagy kétoldali tesztet kell alkalmaznunk? (gondolj arra, hogy a hipotézisünkben megjelöljük-e a hatás vagy különbség irányát vagy sem)
- Mi a null-hipotézis ebben az esetben?
- Ez után végezd el a tesztet
- És írd le a fentiek szerint az eredményeket.

A statisztikai tesztek eredményének közléséről általában

A statisztikai tesztek eredményének közlése során a következő információkat szoktuk megadni általánosságban. Ez tesztől tesztre változhat, de az alábbiak közül minél több információt megadnunk, annál jobb.

- az eredmény szöveges leírása
- teszt-statisztika
- szabadságfok (ez egyszerű teszteknel általában az elemszámmal is megadható)
- p-érték
- hatás mértéke (parameterbecslés)
- hatásmérték 95%-os konfidencia intervalluma