

Seminar_5

Zoltan Kekecs

March 9, 2021

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
```

Hipotezistesztes

A statisztikai inferencia, és hipotézis tesztelés során az a célunk, hogy megállapítsuk, letezik-e egy bizonyos hatás vagy kapcsolat. De ezt a **null-hipotézis szignifikancia tesztelés (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 sejtésem, hogy **egy penzérme 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 **érme nem cinkelt**. Vagyis a null-hipotézis szerint ugyanakkora a valószínűsége fejnek és irást 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 kódrészlet 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
```

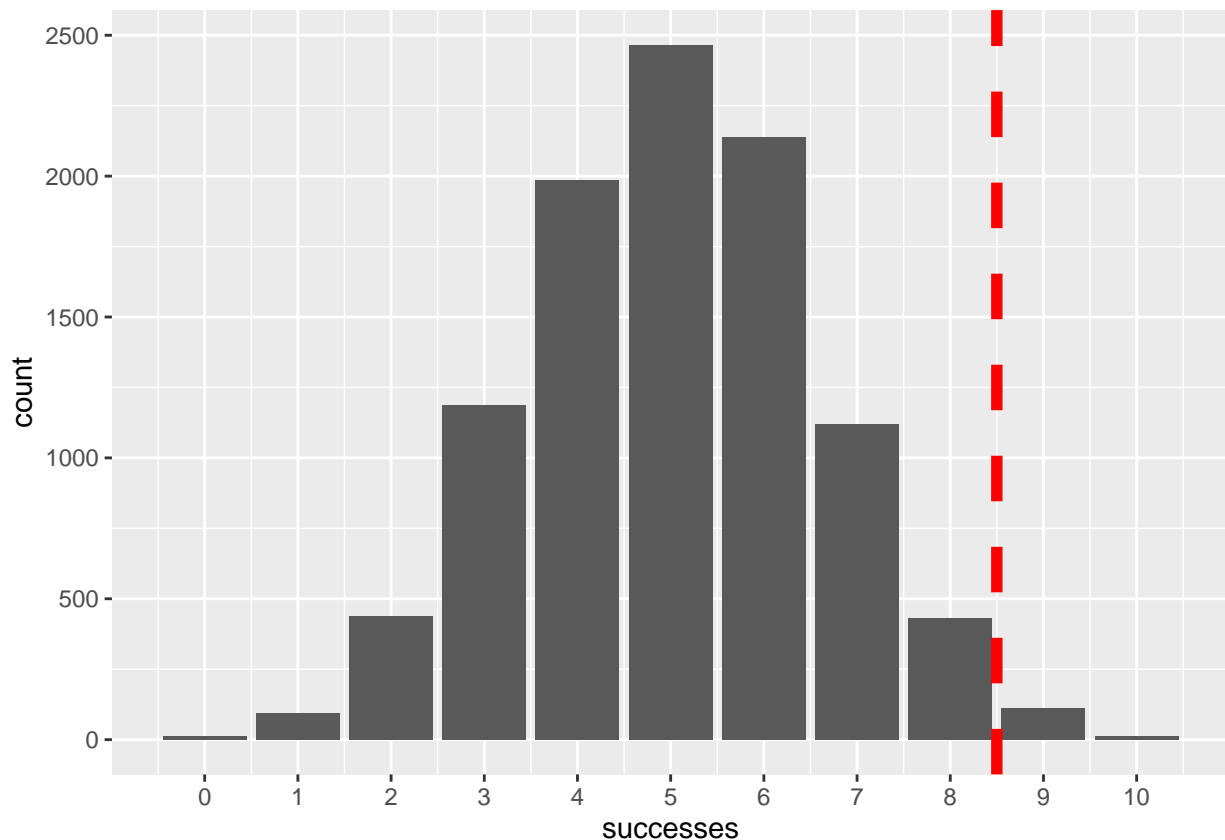
```
## [1] 0.01074219
```

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 varnának, hogy 9 vagy több fejet kapjunk.

Ezt le is ellenőrizhetjük, ha **randomizálunk 10.000 hasonló kísérletet** az `rbinom()` funkcióval. Az ábrán látható hogy csak a kísérletek igen kis százalékában kaptunk 9 vagy több “sikert”. (Ezt a kódrészt sem fontos megérteni, a lényege hogy az `rbinom()` funkcióval 10000-szer azt szimuláltuk, hogy egymás után 10-szer feldobtuk egy érmet (vagyis hogy véletlenszerűen választottunk egy számot 0 és 1 közül), ez után ennek a 10000 kísérletnek az eredményét ábrázoltuk a `ggplot`-tal)

```
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)
```



Vagyis a 9 fej 10 feldobásból 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, hogy az NHST-t teljesen meg tudjuk érteni.

Tegyük fel hogy egy **“igen-vagy-nem” dontest** kell hoznunk arról, hogy cinkelt-e az érme vagy sem.

Mondjuk minket biztak meg hogy ellenorizzuk az ermet egy fontos penzfeldobas előtt, es el kell dntenunk, hogy megbizunk-e ennek az ermenek a hitelessegeben, vagy kerjunk egy uj ermet a penzfeldobashoz, mert ezt cinkeltnek iteljuk. Itt jon az NHST **teszt** resze. Ezt a dontest az NHST-ben egy elore meghatarozott valoszinusegi kuszobertek, **dontesi kuszobertek**, figyelembevetelevel hozzuk meg. Ha az altalunk megfigyelt eredmény **kelloen meglepo, kelloen ritka** a null hipotezis helyesseget feltetelezve, akkor elvetjuk azt a feltetelezest, hogy a null-hipotezis helyes. Ilyenkor kizarasos alapon az alternativ hipotezis helyesseget fogadjuk el.

A pszichologia tudomanyaban a dontesi kuszobertek tradicionalisan 5%, vagyis ha annak a valoszinusege hogy az altalunk megfigyelt eredményt (vagy annal extremer eredményt) kapjunk a null hipotezis helyessege eseten **kisebb mint 5%** ($p < 0.05$), akkor **elvetjuk a null-hipotezist**.

Fontos azonban hangsulyozni, hogy egy-egy NHTS teszt során nem tudjuk meg a null hipotezis helyessegenek, vagy az alternativ hipotezis helyessegenek a valodi valoszinuseget. Csak azt tudjuk, hogy mennyire valoszinu vagy valoszinutlen hogy az altalunk megfigyelt eredményt latjuk “egy olyan vilagban” ahol a null hipotezis helyes. Se tobbet, se kevesebbet. Es ez alapon hozzuk meg a dontesunket a null-hipotezis elveteserol, vagy megtartasarol.

Az NHST modszer fo elonye, hogy ha **konzisztensen használjuk a fent említett dontesi kuszobot** a kutatásainkban, akkor **elegge biztosak** lehetünk abban, hogy a statisztikai donteseinknek **csak 5%-aban vetjuk el hibasan a null hipotezist**. Vagyis a statisztikai donteseknek csak 5%-a lesz hibás, ha a null hipotezis valojaban igaz, így tehet az elsofaju hiba (alpha-error) valoszinusege 5%. (Masszoval csak a teszteknek csak 5%-aban allitjuk hibasan, hogy van hatas, amikor valojaban nincs hatas.)

Két fontos kitetelt erdemes megfigyelni a fenti allitasban. Egyreszt hogy azt irtam hogy “elegge biztosak” lehetünk. Azert csak “elegge biztosak” lehetünk ebben, es nem teljesen biztosak, mert ahhoz hogy ez az allitas helyes legyen, az altalunk hasznalt statisztikai tesztek **elofelteveseinek teljesulnie kell**, es ebben nem lehetünk teljesen biztosak a populacio szintjen. A masik, hogy **“ha a null hipotezis valojaban igaz”**. Arrol az NHST-ben nem kapunk garanciat, hogy a statisztikai donteseinknek hany szazaleka hibás ha az alternativ hipotezis az igaz. Azt is fontos megerteni, hogy ez nem jelenti azt, hogy az osszes publikalt null hipotezis-tesztesben csak 5%-nyi lenne az elsofaju hiba, mert nem minden statisztikai dontest publikálnak.

Statisztikai tesztek

Nem kell jonak lennunk valoszinusegszamitasbol hogy jo statisztikai donteseket tudjunk hozni. A megfigyeles valoszinuseget a null-hipotezis helyesseget feltetelezve altalaban egy **statisztikai teszt** mondja meg nekunk. Ezen az oran 5 statisztikai tesztet fogunk megismerni.

- binomialis teszt
- khi-negyzet teszt
- t-teszt
- egyszempontos ANOVA
- korrelacios teszt

binomialis teszt

A hipotezist, hogy az erme cinkelt, tesztelhetjuk a **binomialis teszt**tel, aminek R-ben `binom.test()` a funkcioja. Az x helyere a megfigyelt “celmegfigyelesek” vagy “sikerek” szamat (a mi esetunkben a fejek szamat, $x = 9$), az n helyere az osszes megfigyeles szamat ($n = 10$), a p helyere pedig a **null-hipotezis** helyesseget feltetelezve a “celmegfigyelesek” eleresenek valoszinuseget kell beírni (mivel a hipotezisunk az hogy az erme cinkelt, az null hipotezisunk az, hogy az erme “nem cinkelt”). Ezt valoszinusegkent kell megadni, amit egy 0 és 1 kozotti szammal 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 jelentese hogy a megfigyelesek 60%-a lesz “siker”). A mi esetunkben a null hipotezis helyessege eseten a fej valoszinusege 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 H1, 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 való hatásmértéket (ami a mi esetünkben a “siker”/fej valószínűsége). Fontos, hogy nem tudjuk, hogy a mi konkrét kísérletünkben a konfidencia intervallum tartalmazza-e a való hatásmértéket.
- sample estimates: A “siker” (“celmegfigyeles”, a mi esetünkben a fej) valószínűségenek becsült mértéke 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 pénzfeldobásból (90%). Ez alapján úgy ítéltük, hogy annak a valószínűsége, hogy fejet dobunk az érmevel szignifikánsan 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(
anxiety_post <- anxiety_base + treatment * treatment_effect + resilience * resilience_effect + gender_n
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 hatékonysagát** teszteltek. Olyan személyeket vontak be a kutatasba, akik egy **hurrikan**

aldozatai voltak, es **szorongással** kuszzkodtek. A személyeknel felmerte a reziliencia (psziches ellenal-lokepesseg) szintjet, majd veletlenszeruen osztottak a személyeket egy kezelesi vagy egy kontrol csoportba. Ezt kovetoen a kezelesi csoport **pszichoterapiat kapott 6 heten keresztul** heti egyszer, mig a kontrol csoport nem kapott kezelest. A vizsgalat vege megmerte a személyek **szorongasszintjet**, es a klinikai kriteriumok alapjan meghataroztak, hogy a személy **gyogyultnak, vagy szorongonak** szamit-e.

Lathatjuk, hogy 8 valtozo van az adattablaba.

- participant_ID - reszvevo azonositoja
- gender - nem
- group - csoporttagsag, ez egy faktor valtozo aminek ket szintje van: “treatment” (kezelt csoport), es “control” (kontrol csoport). A “treatment” csoport kapott kezelest, mig a “control” csoport nem kapott kezelest.
- resilience - reziliencia: a nehezsegekkel valo megkuzdes kepessege, ez egy személyes kepesseg, olyasmi mint a személyisegvonasok
- anxiety_baseline - szorongas szint a terapia elott
- anxiety_post - szorongas szint a terapia utan
- health_status - a klinikai kriteriumok alapjan szorongonak vagy gyogyultnak tekintheto a személy
- home_ownership - lakatasi helyzet: harom szintje van az alapjan hogy a személy hol lakik: “friend” - baratnal vagy csaladnal lakik, “own” - saját tulajdonu lakasban lakik, “rent” - berelt lakasban lakik,
- height - magassag

Adatellenorzes

Mint mindig, elemzes elott **ellenorizzuk**, hogy az adattal minden rendben van-e!

data

```
## # A tibble: 80 x 9
##   participant_ID gender group resilience anxiety_baseline anxiety_post
##   <chr>          <fct> <fct>      <dbl>          <dbl>          <dbl>
## 1 ID_1          male  trea~      5.08           18.7           10.5
## 2 ID_2          male  trea~      6.41           10.5           7.61
## 3 ID_3          male  trea~      7.52           17.2           9.72
## 4 ID_4          female trea~      4.7            17.7          14.7
## 5 ID_5          male  trea~      7.39            8.04           8.14
## 6 ID_6          female trea~      7.06           10.3           10.1
## 7 ID_7          female trea~      7.17           12.6           6.64
## 8 ID_8          male  trea~      9.23           10.6           8.09
## 9 ID_9          male  trea~      4.56           12.8           10.4
## 10 ID_10         male  trea~      9.53            5.82           4.28
## # ... with 70 more rows, and 3 more variables: health_status <fct>,
## #   home_ownership <fct>, height <dbl>
```

data %>%

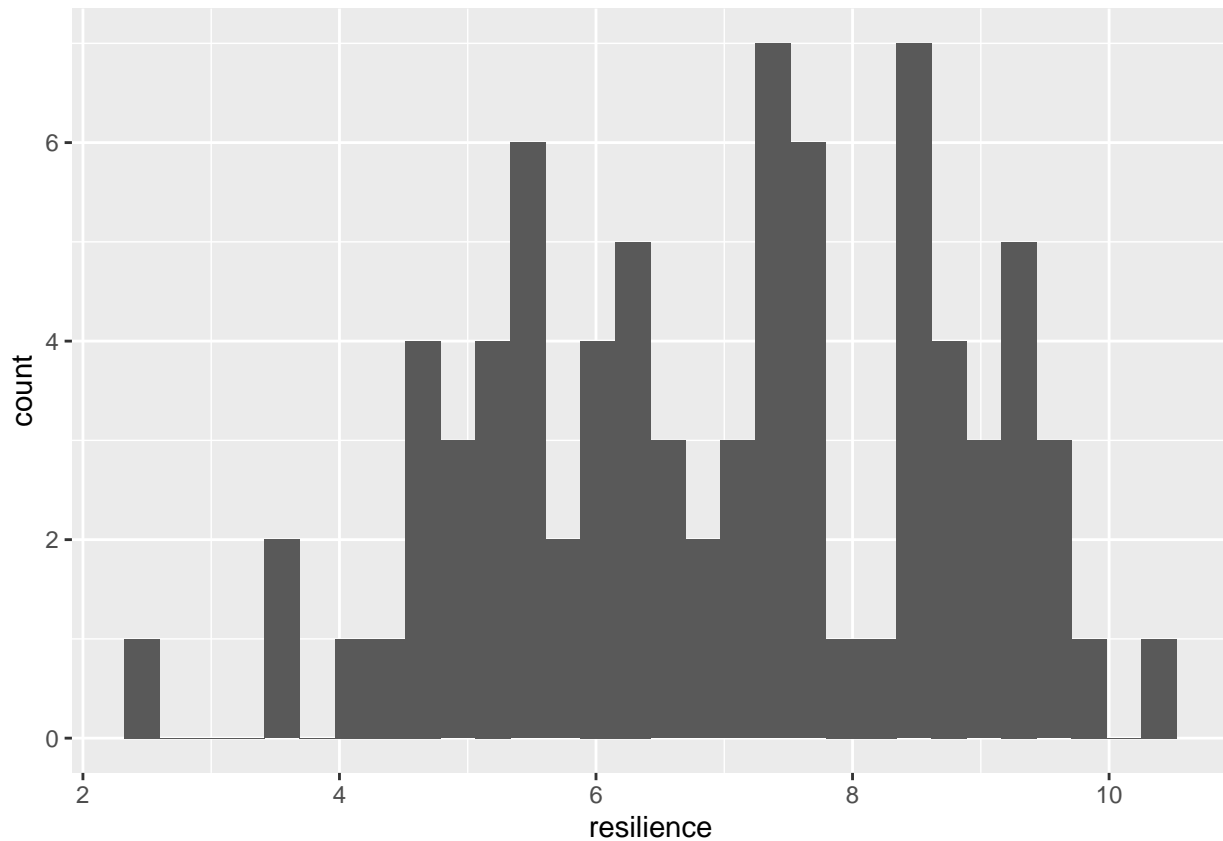
summary()

```
## 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
```

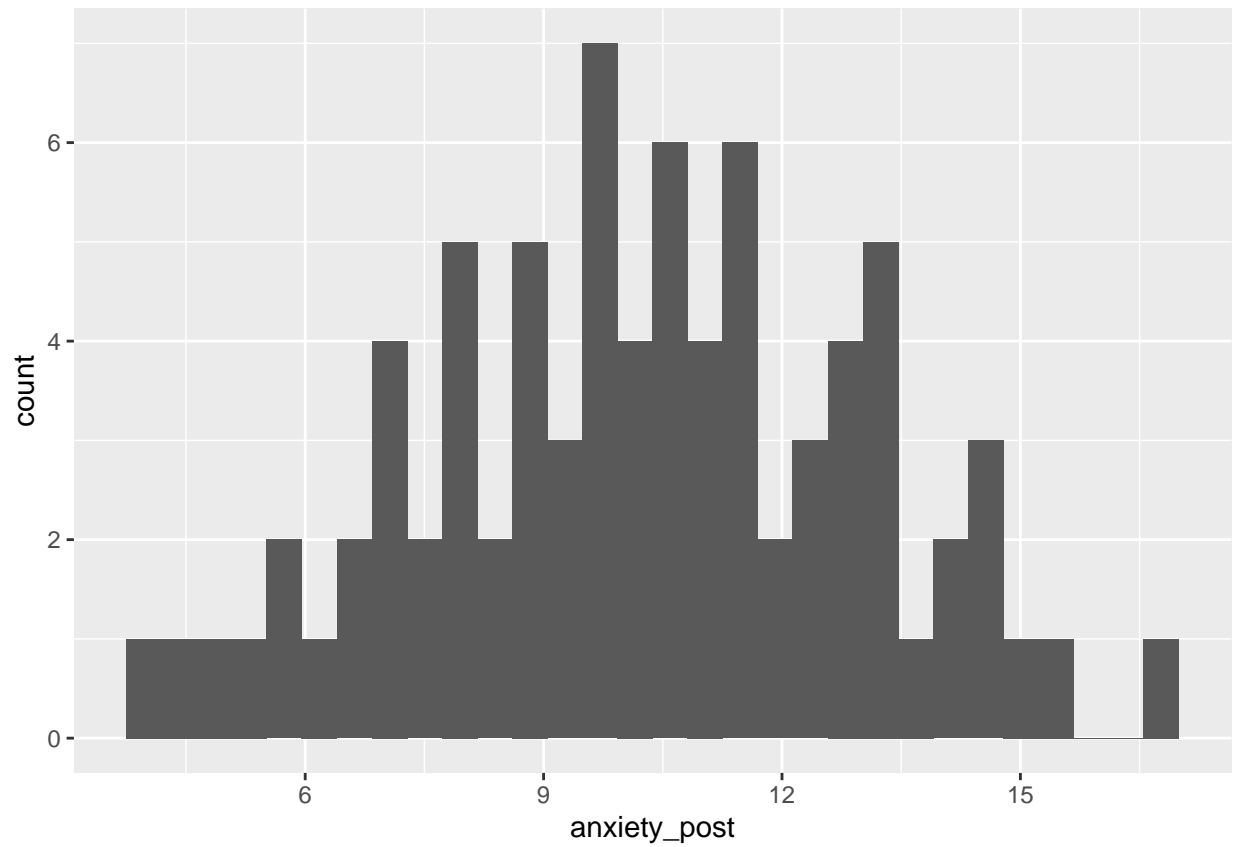
```
data %>%
  ggplot() +
    aes(x = resilience) +
    geom_histogram()
```

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

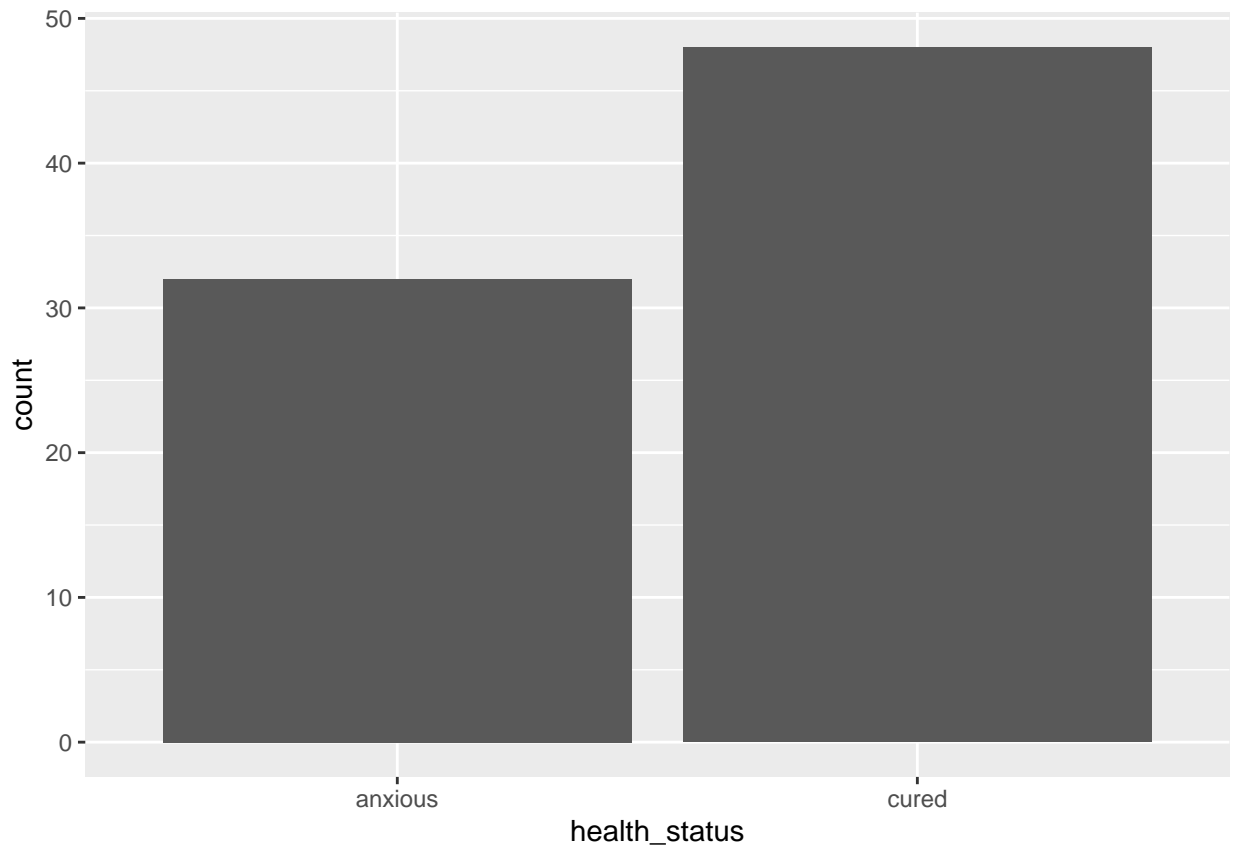


```
data %>%
  ggplot() +
    aes(x = anxiety_post) +
    geom_histogram()
```

```
## `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 hipotézise 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 végen mért szorongásszint negatív összefüggést fog mutatni (vagyis aki reziliensebb, annál alacsonyabb szorongásszintet fognak mérni a kutatás végen) (**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 tudj tölteni a `binom.test()` függvény paramétereit.
- Ez után vedd el a tesztet

- Es ird le a fentiek szerint az eredményeket.

Ket kategorikus valtozo kapcsolata: Khi-negyzet proba (Chi-squared test)

Ket kategorikus valtozo kapcsolatának vizsgalatara tobbfajta teszt is alkalmazhato. Olyan esetben, ahol mind a ket kategorikus valtozonak ket-ket szintje van csak (vagyis a kulonbozo kategoria-szint kombinaciok egy 2x2-es tablazatban abrazolhatoak) a Fisher tesztet erdemes alkalmazni vagy a likelihood ratio tesztet, de olyan tablazatokra ahol a kategoriai kombinacioja tobb mint 2x2-es tablazatot alkot, vagyis ahol az egyik csoportosito valtozonak tobb mint 2 szintje van, a **Khi-negyzet proba** javasolt.

Peldaul megvizsgalhatjuk, hogy van-e kapcsolat abban, hogy a személyek lakhatasi helyzete (**home_ownership**) es a kozott, hogy a kutatas vegen az egyes személyek meggyogyultak-e (**health_status**).

A Khi-negyzet proba elofeltetelei:

- Minden megfigyeles fuggetlen a tobbi megfigyelestol (pl. egy megfigyeles személyenkent)
- A kategoria-kombinaciok abrazolasaval kapott tablazatban nem tobb mint a cellak 20%-aban kisebb a varhato ertek 5-nel, es minden cellaban magasabb a varhato ertek mint 1.

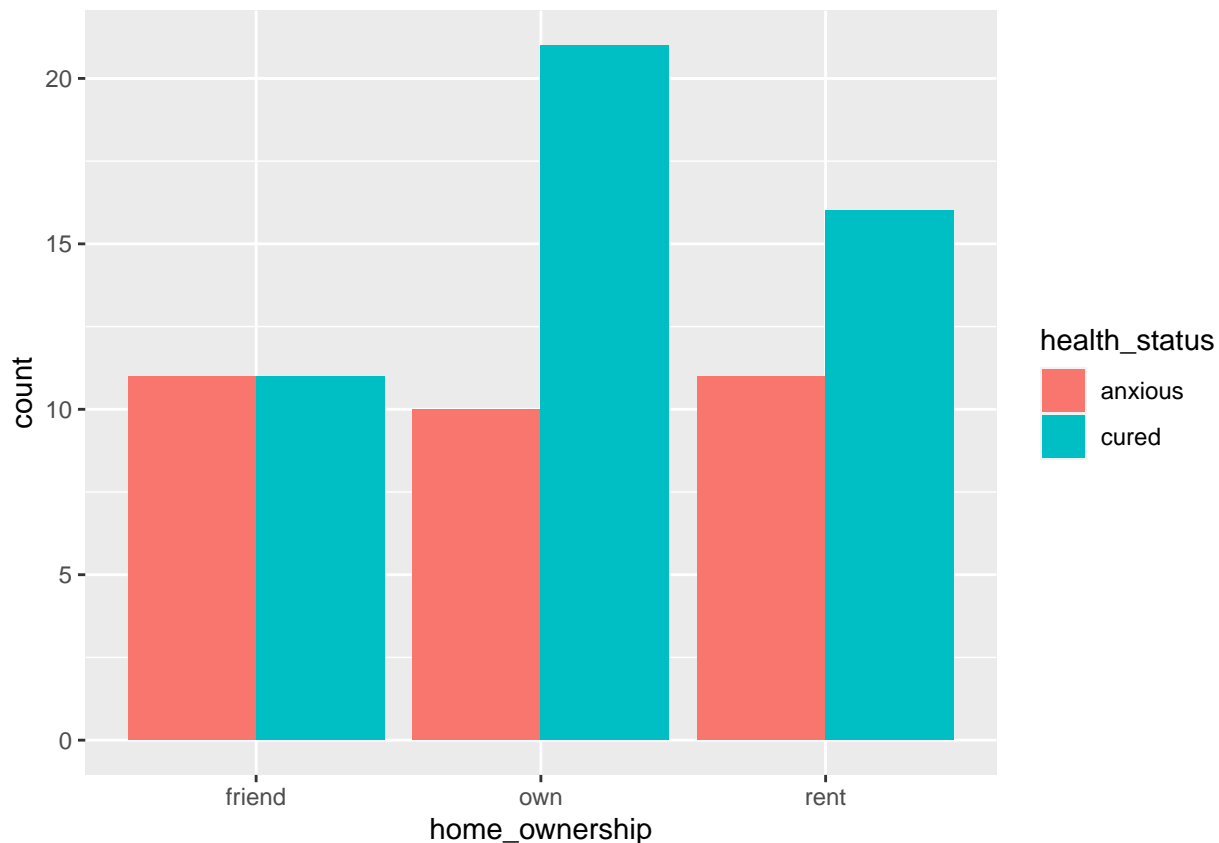
Eloszor feltaro elemzest vegzunk:

- tablazatot rajzolunk a ket valtozo kapcsolatarol
- abrat keszitunk (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).

```
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
```

Az eredményt így írhatjuk le:

“Nem volt szignifikáns eltérés abban, hogy a különböző lakhatási csoportokban (barátnál, saját lakásban, vagy bérlésben lakók) milyen arányban voltak azok akik meggyógyultak a kutatás végére ($X^2 = 1.7$, $df = 2$, $p = 0.428$).”

Gyakorlas

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 táblázatot 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**).

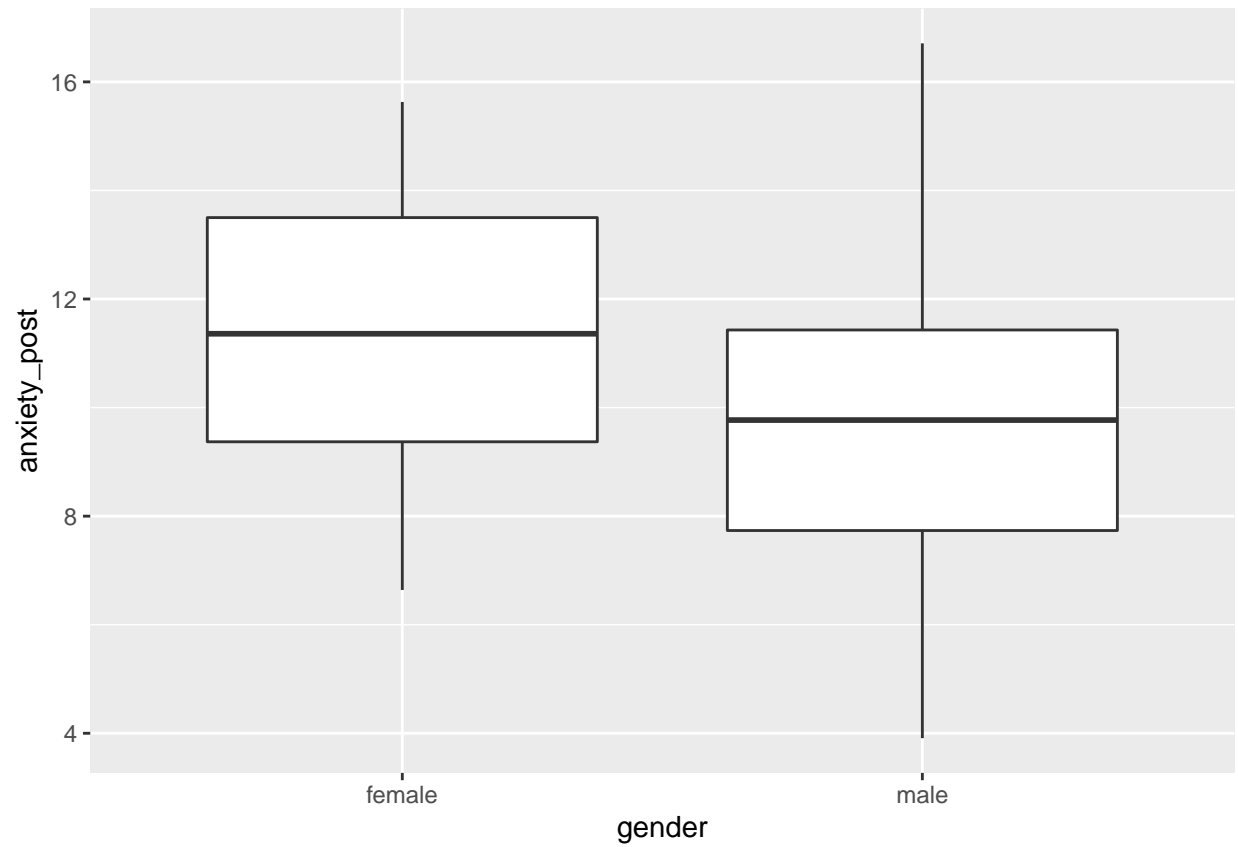
Eloszor szokás szerint feltárási elemzést végzünk átlagok csoportonkénti összehasonlításával és ábrával. Erre pl. remek a `geom_boxplot()` és a `geom_density()`

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

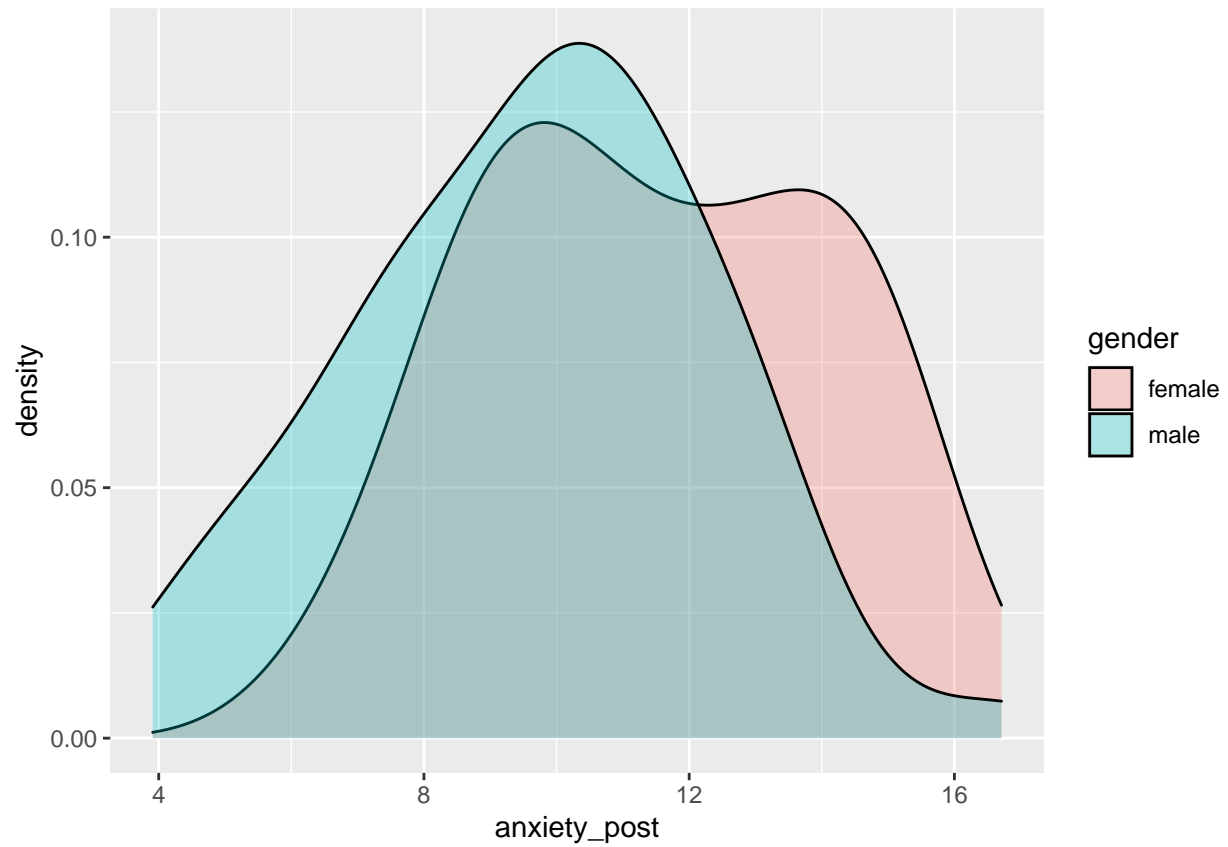
## `summarise()` ungrouping output (override with `.groups` argument)
summary

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

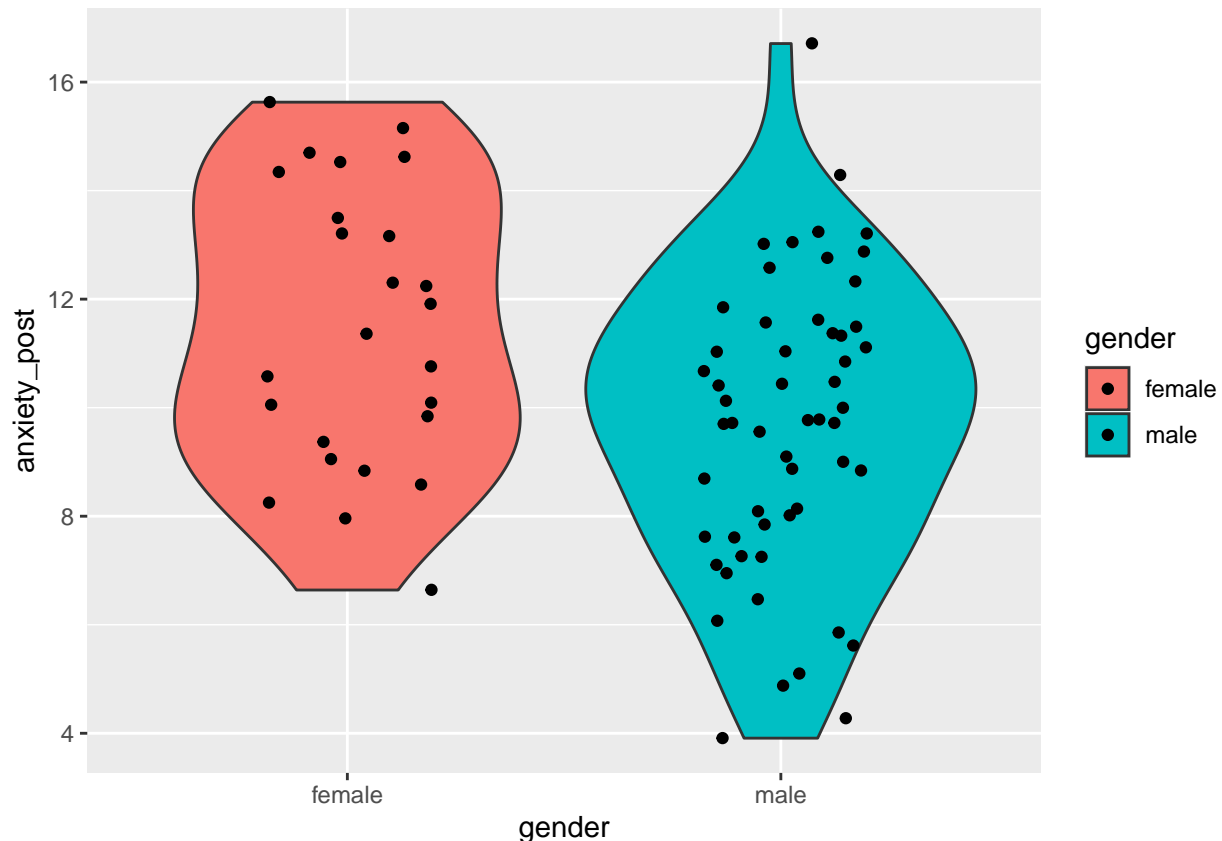
data %>%
  ggplot() +
  aes(x = gender, y = anxiety_post) +
  geom_boxplot()
```



```
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 alapjan, hogy a nok szorongasszintje nagyobb valamivel mint a ferfiake 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 aranysskalan mozog
- A fuggetlen valtozo ket egymastol fuggetlen kategorikus csoportot reprezental
- A megfigyelesek fuggetlenek egymastol. Minden megfigyeles csak az egyik csoportba sorolható, es a csoportok kozott nincs osszefugges az egyes megfigyelesek között.
- Nincsenek jelentos kiugro esetek.
- Csoportonkent normalis eloszlast 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 feltetel serul).

```
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 is not equal to 0
```

```
## 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 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(2.7). 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()` függővel. A formula ugyan úgy néz ki, mint a t-teszt esetén.

Az egyszempontos ANOVA előfeltételei majdnem ugyan azok, mint a független mintas 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**.

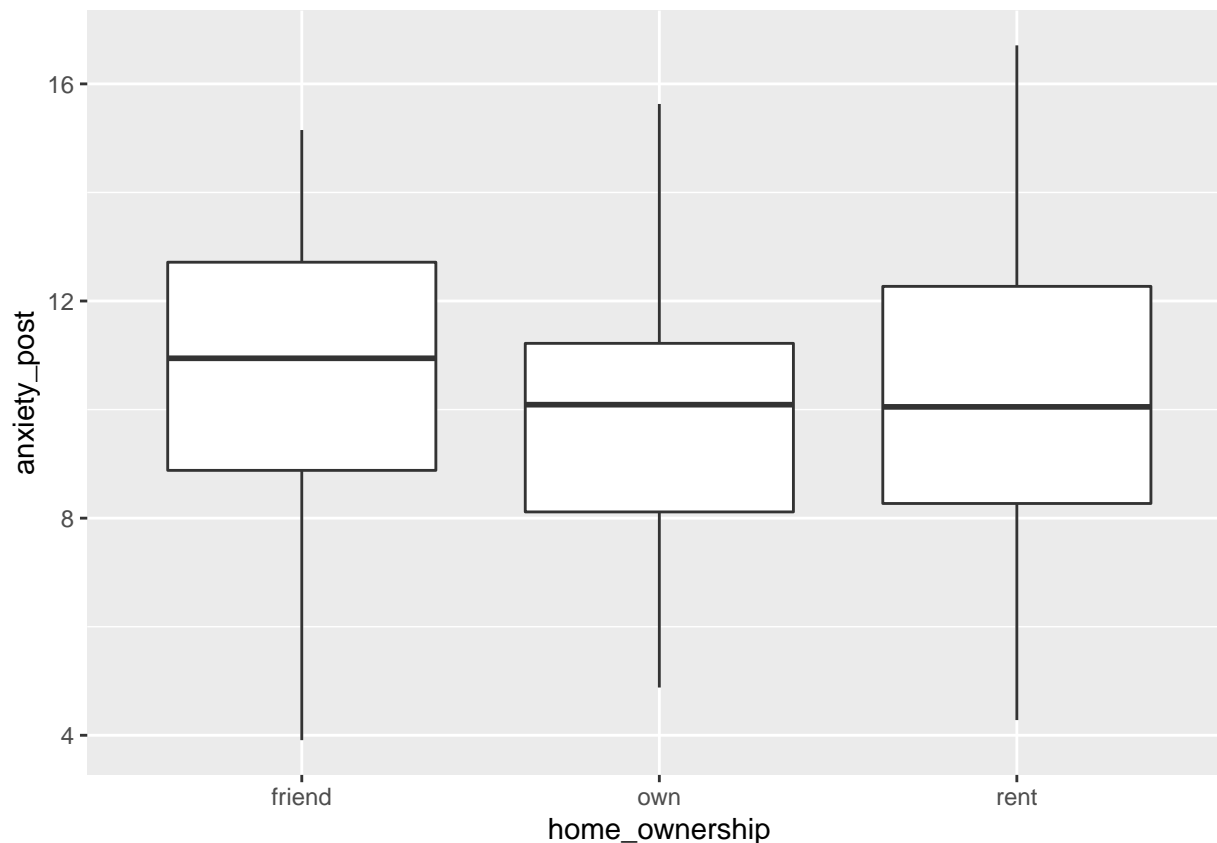
```
summary_home_ownership_vs_anxiety_post = data %>%
  group_by(home_ownership) %>%
  summarize(mean = mean(anxiety_post), sd = sd(anxiety_post))
```

```
## `summarise()` ungrouping output (override with `.groups` argument)
```

```
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 produkálnánk a megfelelő táblázatot a szorongás átlagával home_ownership csoportok szerint.

Egyoldalu vs. kétoldalu tesztek

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

Például tegyük fel, hogy amikor a hipotézisünket meghatároztuk (ideális esetben ez még az adatgyűjtés előtt megtörténik), úgy gondoltuk, hogy a nőknek magasabb lesz a szorongásszintjük, mint a férfiaknak. Ezt az $\text{alternative} = \text{“greater”}$ 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 (∞).

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-hipotezis helyesseget feltételezve. Vagyis amikor tudjuk, milyen irányu hatást várunk el, mindig érdemes egy-oldali tesztet alkalmazni, mert ezzel nő a statisztikai erőnk a hatás kimutatására.

Az egyoldali tesztek esetén amikor az a hipotézisünk, hogy **a referencia-csoport** átlaga magasabb lesz, (alternative = “greater”), akkor a konfidencia intervallumnak csak az alsó határt 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”**.

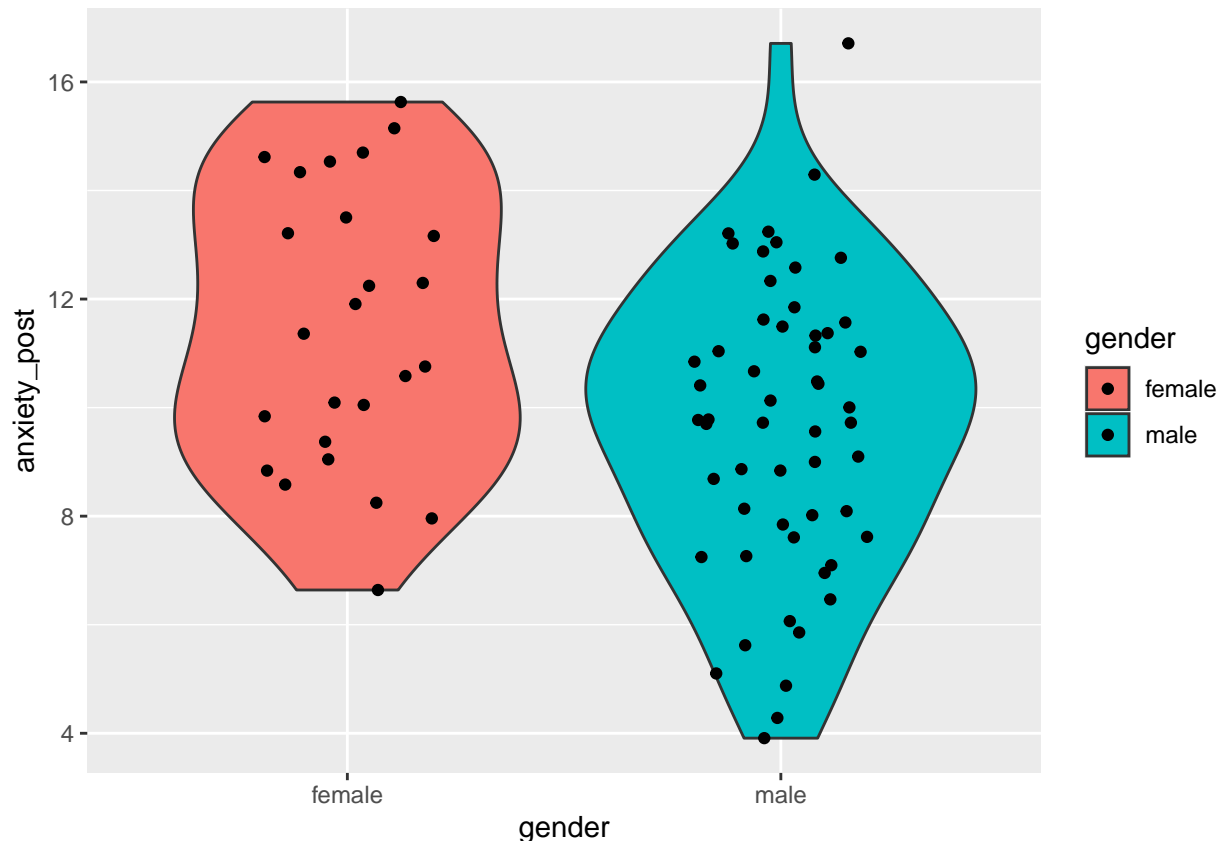
Ahogy korábban már volt róla szó, a referencia-csoport (vagy referencia-szint egy faktor változóban) alapértelmezett módon a faktorszintek nevének ABC sorrendje alapján dolgozik, az ABC sorrendben az első faktorszint lesz a referencia-szint. A példánkban a gender változóban a “female” a referencia-szint, mert az ABC sorrendben a “male” előtt van. Azt, hogy mi legyen a referencia-szint a korábban tanultak szerint a **factor()** **funkcióban a levels =** parameter 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))

## `summarise()` ungrouping output (override with `.groups` argument)
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 is greater than 0
## 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 szorongás szintjükben ($t = 2.89$, $df = 48.52$, $p = 0.003$). A csoportok szorongás szintjének átlaga és szórása a következő volt: nők: $11.47(2.58)$, férfiak: $9.64(2.7)$. 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 helyeslegelt feltételezve ilyen, vagy ennél extremerbb 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 megfigyeles egyaltalan nem segit abban, hogy elutasitsuk a null-hipotezist.

```
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 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 erdemes megjegyezni, hogy a “greater” es a “less” minding a kategorikus valtozo **referencia-szintjere** vonatkozik. Ha ezt nem allitottuk be maskepp, pl. a **factor (levels =)** funkcioval, akkor a referencia-szint az ABC sorrendben elorebb levo szint lesz. A fenti esetben a ket szint a “female” es a “male”, amik kozul a “female” jon elobb ABC sorrendben. Ha azt tippeltuk volna, hogy az lenne a hipotezisunk, hogy a ferfiak (“male”) szorongasszintje lesz magasabb, akkor alternative = “less”-t kellene beallitanunk, mert ezzel egyben azt tippeljuk, hogy a referenciaszint (“female”) atalaga lesz az alacsonyabb. Vagy at kellene allitani a referenciaszintet.

Gyakorlas

Teszteld a 3. hipotezist, hogy “A terapias csoportban alacsonyabb lesz a szorongas atlaga a kutatás vegere mint a kontrol csoportban” (**anxiety_post** vs. **group**)

- Eloszor vegezzunk egy feltaro elemzest egy tablazattal a ket valtozo kapcsolatarol a summarize(mean(), sd()) funkeiokkal, es keszitsunk abrat, mondjuk geom_boxplot() segitsegevel.
- egy- vagy ketoldalu tesztet kell alkalmaznunk? (gondolj arra, hogy a hipotezisunkben megjosoljuk-e a hatas vagy kulonbseg iranyat vagy sem)
- Mi a null-hipotezis ebben az esetben?
- Melyik tesztet erdemes hasznalni, az egyvaltozos ANOVA-t, vagy a t-tesztet? (gondolj arra, hogy hany csoport (szint) van a kategorikus valtozon belül)
- Ez utan vegezd el a tesztet
- Es ird le a fentiek szerint az eredmenyeket.

Ket numerikus valtozo kozotti kapcsolat, korrelacio, cor.test()

Vizsgaljuk meg, van-e egyuttjaras a reziliencia (**resilience**) es a magassag (**height**) kozott.

Eloszor vegezzunk feltaro elemzest a korrelacios egyutthato kiszamitasaval, es egy pontdiagrammal. Használjunk geom_point() es geom_smooth() geomokat egyszerre, es használjuk az “lm” modszert a trendvonal megrajzolasara.

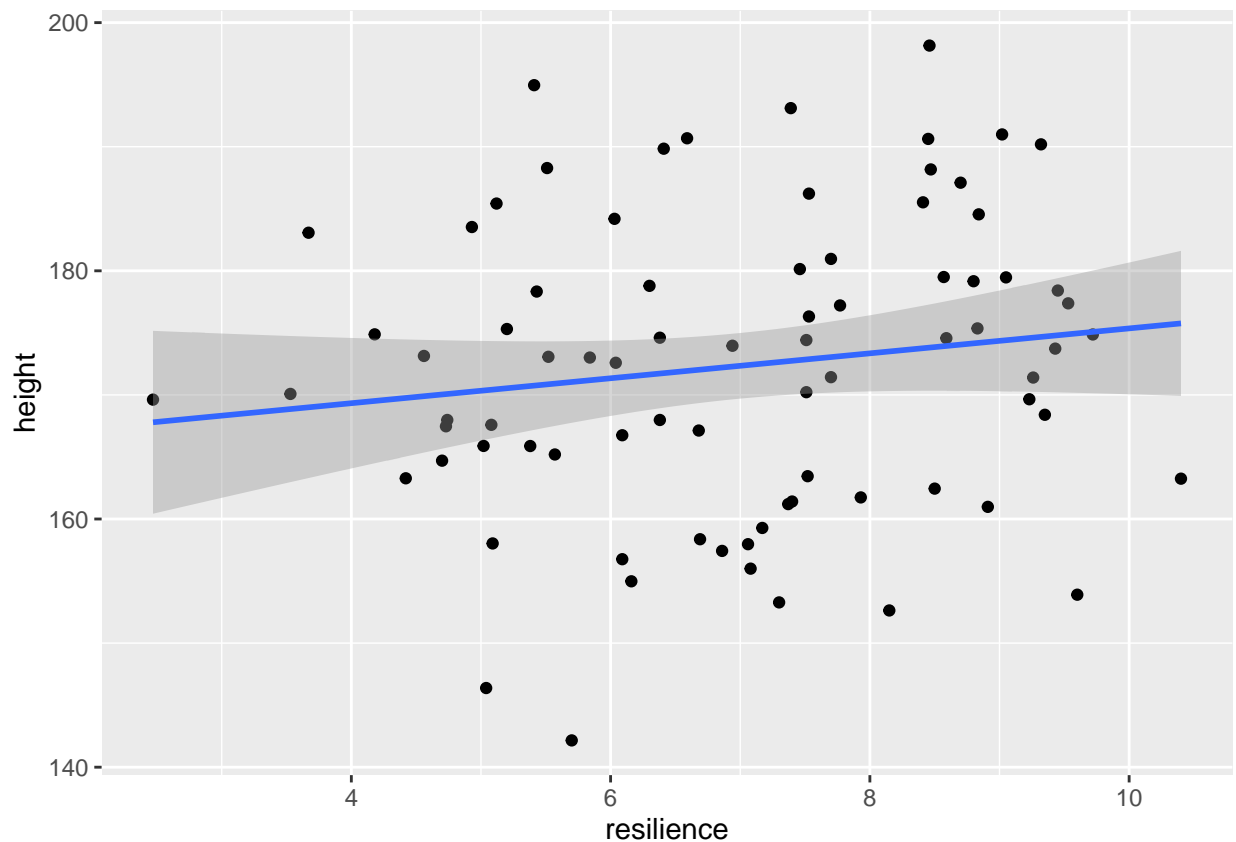
```
data %>%
  select(resilience, height) %>%
  cor()
```

```
##           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 ket változó függetlennek tűnik egymástól a feltáró elemzés alapján, de elképzelhető, hogy a hatas, bármilyen kicsit is, mégis statisztikailag szignifikáns, szóval vegezzük el a statisztikai tesztet is.

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

The assumptions of pearson's correlation are:

- Ket folytonos skalaju változó. Ha bármelyik változó ordinalis skalaju, akkor a spearman korrelációt lehet használni.
- Minden megfigyelesi egységhez ket érték tartalmazzon.
- Nincsenek jelentős kiugró értékek
- Linearitás. A ket változó kapcsolata egy egyenes vonallal jellemezhető.
- Normalitás: mindket változó normalis eloszlást mutat. Nem normalis eloszlás esetén a Spearman korreláció használható.

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

```
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.

Peldául 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 specifikációjában, 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
```

Gyakorlás

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**)

- Eloszor vegezzünk egy feltárási elemzést a korrelációs együttható meghatározásával és egy pontdiagrammal a két változó kapcsolatáról.
 - egy- vagy kétoldalu 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 vedd el a tesztet
 - Es írd le a fentiek szerint az eredményeket.
-

A statisztikai tesztek eredményenek közleserol altalaban

A statisztikai tesztek eredményenek közlese során a következő információkat szoktuk megadni általánosságban. Ez tesztrol 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
- hata mértéke (parameterbecslés)
- hatásmérték 95%-os konfidencia intervalluma