

Kevert modellek - alapok

Zoltan Kekecs

17 November 2020

Contents

1	Absztrakt	2
2	Adatmenedzsment es leiro statisztikak	2
2.1	Package-ek betoltese	2
2.2	Sajat funkcio	2
2.3	A Bully-zas adatbazis betoltese	2
2.4	Adatellenorzes es adattisztitas	4
3	A kevert modellek alapfogalmai	4
3.1	Clustering (csoportosulas) feltarasa	4
3.2	Kevert modellek	6
3.3	A hatasok (prediktorok) ket tipusa	6
3.4	Kevert modellek felepitesi az R-ben	8
3.5	Melyik modell reprezentalja legjobban a valosagot?	9
3.6	Mit kell kozolni az elemzesrol	13

1 Absztrakt

Az eddig tanult lineáris regressziós modellek a csoportokba rendeződött adatokat úgy kezelik, hogy prediktorként bevonják azokat a modellbe. Ez remekül működik ha keves csoport van (a csoportosító változónak keves szintje van) és minden csoportot van módunk megfigyelni. Pl. kísérleti vs. kontroll csoport. De ezek a modellek nem jól működnek olyan esetekben ha az adataink csoportokba/klaszterekbe rendeződnek egy olyan változó mentén aminek a kutatásunk célpopulációjában **sok csoportszintjét különíthetjük el, de a mi kutatásunkban ennél kevesebb figyelhető meg**. Ilyen eset például ha a vizsgálati személyeink különböző iskolákból érkeznek, és elképzelhető hogy az iskolának hatása van a kimeneti változóra, de néhány iskolából vannak adataink és nem tudunk az ország összes iskolájából mintát venni, így sok lehetséges iskola hiányzik az adatok közül. Ilyen esetekben **kevert modelleket** célszerű használni.

Ebben a gyakorlatban megismerheted a kevert modellekkel kapcsolatos alapfogalmakat, valamint hogy hogyan lehet őket felépíteni.

2 Adatmenedzsment és leíró statisztikák

2.1 Package-ek betöltése

Ebben a gyakorlatban a következő package-ekre lesz szükség:

```
library(psych) # for describe\td
library(tidyverse) # for tidy code and ggplot\td
library(cAIC4) # for cAIC\td
library(r2glmm) # for r2beta\td
library(lme4) # for lmer
library(lmerTest) # to get singificance test in lmer
library(MuMIn) # for r.squaredGLMM
```

2.2 Saját funkció

Ezzel a funkcióval kinyerhetjük a standardizált Beta együtthatót a kevert modellekből. Ez a funkció innen lett áttemelve: <https://stackoverflow.com/questions/25142901/standardized-coefficients-for-lmer-model>

```
stdCoef.lmerMod <- function(object) {
  sdy <- sd(getME(object, "y"))
  sdx <- apply(getME(object, "X"), 2, sd)
  sc <- fixef(object) * sdx/sdy
  se.fixef <- coef(summary(object))[, "Std. Error"]
  se <- se.fixef * sdx/sdy
  return(data.frame(stdcoef = sc, stdse = se))
}
```

2.3 A Bully-zas adatbázis betöltése

Ebben a gyakorlatban az általános iskolai bully-zasról (magyarul talán “zaklatás”?) teszünk fel kutatási kérdéseket. Ez egy szimulált adatbázis, vagyis nem valódi adatokat tartalmaz, de képzeljük el, hogy az adatok a következő kutatásból származnak: Ebben a kutatásban az érdekel minket, hogy a **testsúly** hogyan befolyásolja a gyerekek **serulekenységet a bully-zással szemben**. A kutatók azt feltételezik hogy a testsúly összefügg az elvett szendvicsek számával.

Változók:

- **sandwich_taken** - A bullyzással kapcsolatos serulekenység méroszáma. A kutatásban megkerdezték a vizsgáltai személyeket (általános iskolai gyerekek) hogy az elmúlt hónapban hányszor kenyészerítették ki toluk a bully-k az ébredre hozott szendvicseit

- **weight** - testsúly
- **class** - faktor változó ami azt mutatja melyik iskolai osztályba jár a vizsgalati személy. Faktorszintek: class_1, class_2, class_3, class_4.

Két adatfajlt is betöltünk. Mindket adatfajl úgy lett legenerálva, hogy a diákok különböznek abban, hogy mennyi szendvicset vesznek el tőlük attól függetlenül, hogy milyen a testsúlyuk és attól függetlenül is, hogy melyik iskolai osztályba járnak. Vagyis mind a testsúlynak, mind az osztálynak van hatása az elvett szendvicsek számára.

Vizsgálat a két adatbázis különbözik abban, hogy az, hogy a diák melyik osztályba jár, befolyasolja-e, hogy a testsúlynak mekkora hatása van az elvett szendvicsek számára. A **data_bully_int.csv** adatfajlban *a testsúly hatása ugyanakkor minden osztályban* (függetlenül az osztálytól), míg a **data_bully_slope.csv** adatfajlban *a testsúly hatása különbözik osztályonként* (néhány osztályban a testsúly hatása nagyobb mint másokban).

```
# load data
data_bully_int = read_csv("https://raw.githubusercontent.com/kekecsz/PSYP13_Data_analysis_class-2018/main/data_bully_int.csv")

## Parsed with column specification:
## cols(
##   sandwich_taken = col_double(),
##   weight = col_double(),
##   class = col_character()
## )

# assign class as a grouping factor
data_bully_int %>% mutate(class = factor(class))

## # A tibble: 80 x 3
##   sandwich_taken weight class
##           <dbl>   <dbl> <fct>
## 1             8     30 class_1
## 2            10     28 class_1
## 3            11     27 class_1
## 4            12     33 class_1
## 5             7     36 class_1
## 6            10     30 class_1
## 7             9     30 class_1
## 8            13     23 class_1
## 9             8     33 class_1
## 10            7     36 class_1
## # ... with 70 more rows

data_bully_slope = read_csv("https://raw.githubusercontent.com/kekecsz/PSYP13_Data_analysis_class-2018/main/data_bully_slope.csv")

## Parsed with column specification:
## cols(
##   sandwich_taken = col_double(),
##   weight = col_double(),
##   class = col_character()
## )

data_bully_slope %>% mutate(class = factor(class))

## # A tibble: 80 x 3
##   sandwich_taken weight class
##           <dbl>   <dbl> <fct>
## 1             8     44 class_1
```

```
## 2      8      38 class_1
## 3      7      40 class_1
## 4      8      42 class_1
## 5      8      36 class_1
## 6     10      29 class_1
## 7      6      43 class_1
## 8      8      35 class_1
## 9      8      31 class_1
## 10     7      36 class_1
## # ... with 70 more rows
```

2.4 Adatellenorzes es adattisztitas

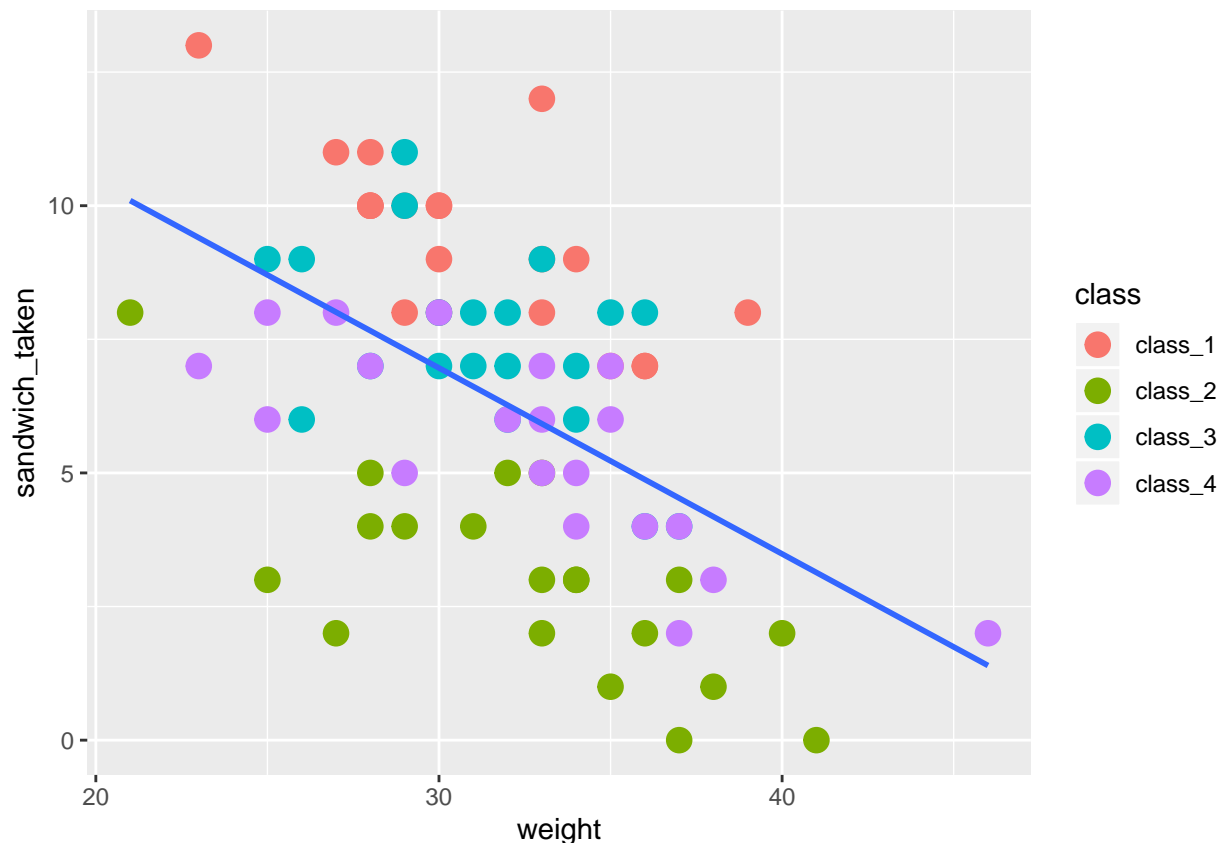
Ahogy mindig, eloszor kezd az adatok ellenorzesével es az esetleges adattisztitással. Ehhez használhatod a `View()`, `summary()`, es `describe()` funkciókat, es a `ggplot()` funkciót vizualizalashoz.

3 A kevert modellek alapfogalmai

3.1 Clustering (csoportosulas) feltarasa

Vizualizaljuk a `sandwich_taken` es `weight` változók összefüggését egy pontdiagram (scatterplot) segítségével. Az adatok egy egyértelmű negatív összefüggést mutatnak a `sandwich_taken` es `weight` változók között, de az adatok variabilitása nagyon nagy.

```
data_bully_int %>% ggplot() + aes(y = sandwich_taken, x = weight) +
  geom_point(aes(color = class), size = 4) + geom_smooth(method = "lm",
  se = F, formula = "y ~ x")
```



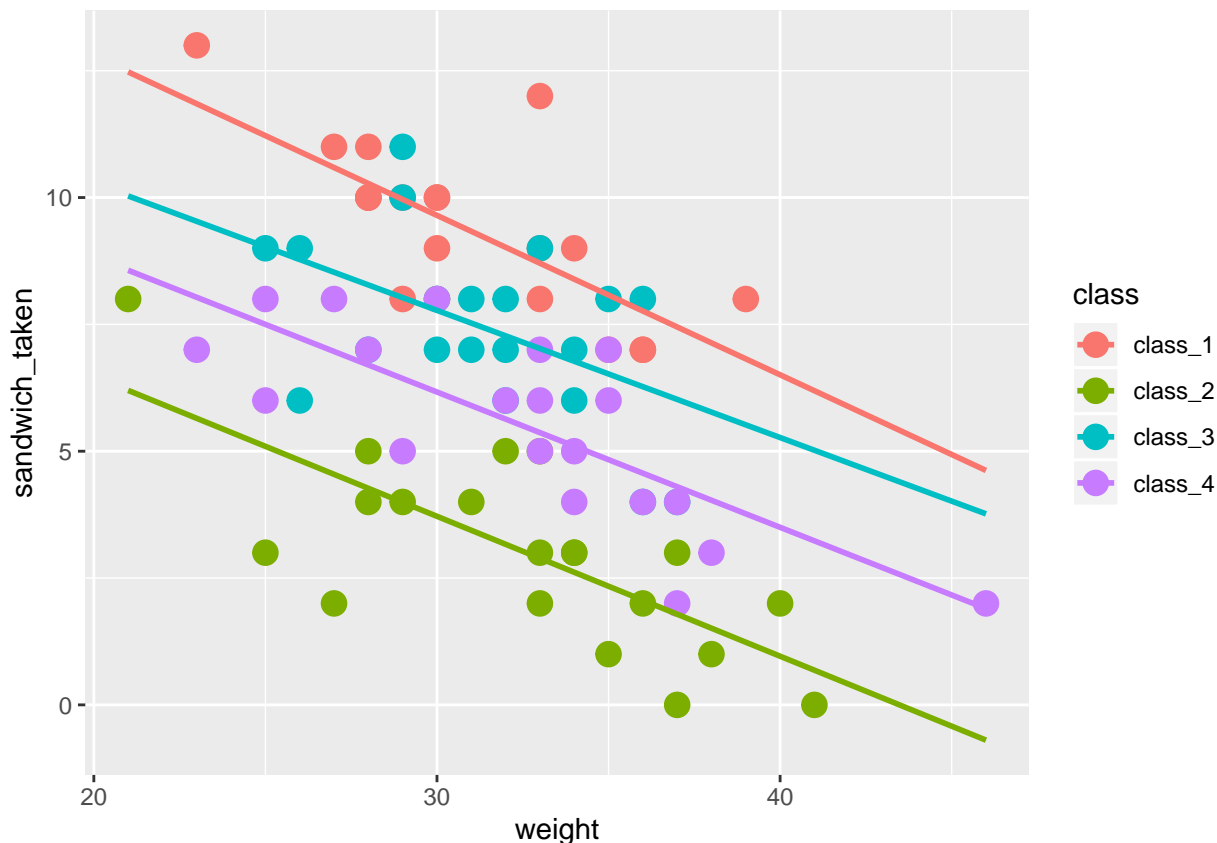
A pontok színe az ábrán azt mutatja, a diák **melyik iskolai osztályba** jár (class_1, class_2, class_3, vagy class_4). Ha jobban megnézzük, úgy tűnik, hogy az azonos színű pontok egymáshoz közel helyezkednek el az ábrán, nem pedig random módon elszórva, ami arra utal, hogy az adatpontok nem teljesen függetlenek egymástól, hanem csoportosulnak (klasztereket alkotnak).

Nezzük meg, hogy az iskolai osztály meg tudja-e magyarázni a variabilitás egy részét. Például felrajzolhatjuk a **regressziós egyeneseket csoportonként**. Ez úgy tűnik, hogy megmagyarázza a variabilitás egy részét, hiszen a regressziós vonalak közelebb kerülnek a valós megfigyelésekhez. Szóval úgy tűnik, hogy érdemes lenne a class változót is figyelembe venni a modellünk megépítésénél.

(Alább az ábrát elmentjük egy int_plot nevű objektumba, hogy később ugyan ezt az ábrát könnyen előhívhassuk.)

```
int_plot = data_bully_int %>% ggplot() + aes(y = sandwich_taken,
  x = weight, color = class) + geom_point(size = 4) + geom_smooth(method = "lm",
  se = F, fullrange = TRUE, formula = "y ~ x")

int_plot
```



3.2 Kevert modellek

Akkor használjuk a kevert modelleket amikor olyan prediktor változónk van, aminek sok szintje/lehetséges értéke van a valóságban, de nekünk ebből a sok lehetséges értékből csak kevésről kapunk információt a mintánkban.

Jelen kutatásban csak az érdekel minket hogy a testsúly befolyasolja-e az elvett szendvicsek számát, és ha igen, mennyire. Az iskolai osztályok hatása nem része a fő kutatási kérdésnek, és még ha az is lenne, az információt amit ezekről az iskolai osztályokról szerzünk **nem tudnánk általánosítani más iskolákban**. Nem lenne sok értelme megtudni, hogy mi a hatása annak ha valaki a “class 1”-be jár, amikor a regressziós modellünket új mintán szeretnénk bejósolni használni egy új iskolában, hiszen a többi iskolában más osztályok vannak, amiknek vélhetően mások a karakterisztikái. Szóval az iskolai osztály ebben az esetben egy “zavaró tényező” (**nuisance variable**). Vagyis ezt a class változót nem szeretnénk figyelembe venni a regressziós egyenletben, hiszen akkor más iskolákban nem tudnánk felhasználni az egyenletet.

A **kevert modellek** segítségével **figyelembe vehetjük az adatok ilyen fajta csoportosulását anélkül hogy a regressziós egyenletünkbe be kellene tennünk ezeket a zavaró tényezőket**.

3.3 A hatások (prediktorok) két típusa

Itt fontos megkülönböztetnünk a **hatások két típusát**.

Fix hatások (Fixed effects) - Azokat a hatásokat amikkel eddig a lineáris modellekben foglalkoztunk, “fix hatásoknak” (fixed effects) nevezzük. Ezek azok a hatások/prediktorok, amikre regressziós együtthatókat számítunk ki. Ezek a predikciós regressziós egyenletünk részei, amiket a későbbiekben is felhasznalunk majd a bejósoláshoz.

Random hatások (Random effects) - Azokat a hatásokat, amiket a “zavaró tényezőknek” tulajdonítunk,

modellezhetjük random hatásként. A random hatásokat úgy modellezzük, hogy bár a megfigyelések különböznek a klaszterek (csoportok) mentén, de az egyes csoportok (mint például itt az osztályok) **hatása nem szisztematikus**, hanem egyfajta véletlenszerű különbségből fakad a csoportok között. A csoportok közötti ilyen véletlenszerű különbség felismerése segít abban, hogy pontosabban kiszámítsuk a fix hatások regressziós együtthatóival kapcsolatos bizonytalanságot (konfidencia intervallumot). Ezek a random hatások **nem kerülnek bele a regressziós egyenletünkbe**, és nem kapunk velük kapcsolatban regressziós együtthatókat.

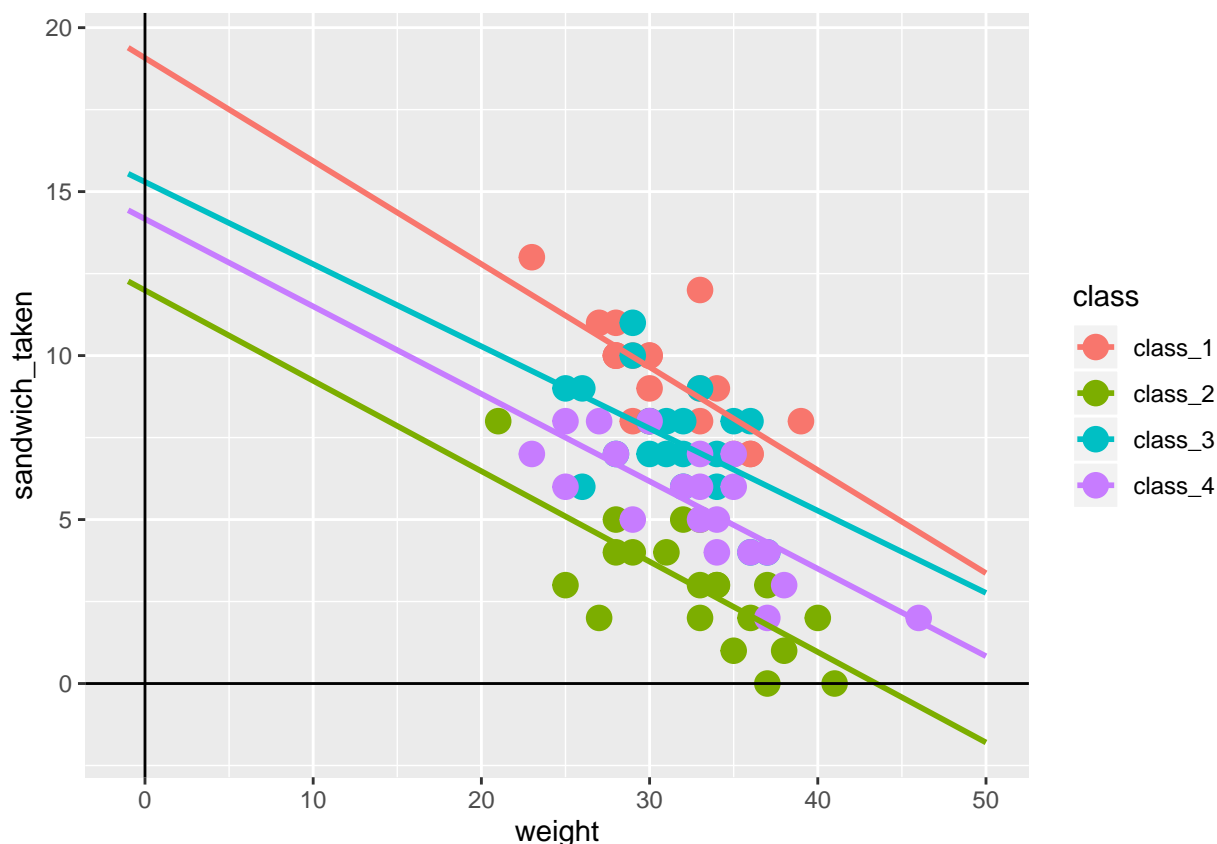
Ezért azokat a modelleket, amik mind fix, mind random hatásokat tartalmaznak, **kevert modelleknek (mixed models)** nevezzük.

3.3.1 A random hatások előfordulási fajtái

Általánosságban a random hatások két módon lehetnek hatással a **kimeneti változóra**. Az egyik hogy **direkt hatást fejtenek ki rá (random intercept)**, a másik hogy a **fix hatások mértékét és irányát befolyásolják (random slope)**.

random intercept, random slope nélkül: Lehetőséges hogy a csoportok (klaszterek) csak abban különböznek egymástól, hogy a **kimeneti változon átlagosan milyen értéket vesznek fel**, de a **fix hatások azonosak** a klaszterek között. Ez igaz a `data_bully_int` adatbázisra. Megfigyelhetjük az ábrán, hogy a regressziós egyenesek meredeksége (slope) nem különbözik az osztályok között, ami arra utal, hogy a testsúly hatása ugyanakkor az egyes osztályokban. Az osztályok csak abban különböznek, hogy milyen „magasan” vannak a regressziós egyenesek, vagyis abban, hogy a regressziós egyenesek milyen értékkel metszik az Y tengelyt. Ez látható az alábbi ábrán is.

```
int_plot + xlim(-1, 50) + geom_hline(yintercept = 0) + geom_vline(xintercept = 0)
```



random intercept, es random slope: A fentiekben csak a `data_bully_int` adatbázist használtuk. Most

vizsgáljuk meg a másik adatbázist (`data_bully_slope`). Ahogy fent említettük, ebben az adatbázisban azt szimuláltuk, hogy az osztálynak nem csak az elvett szendvicsek számára van hatása, hanem **a testsúly hatása is különbözik az osztályok között**.

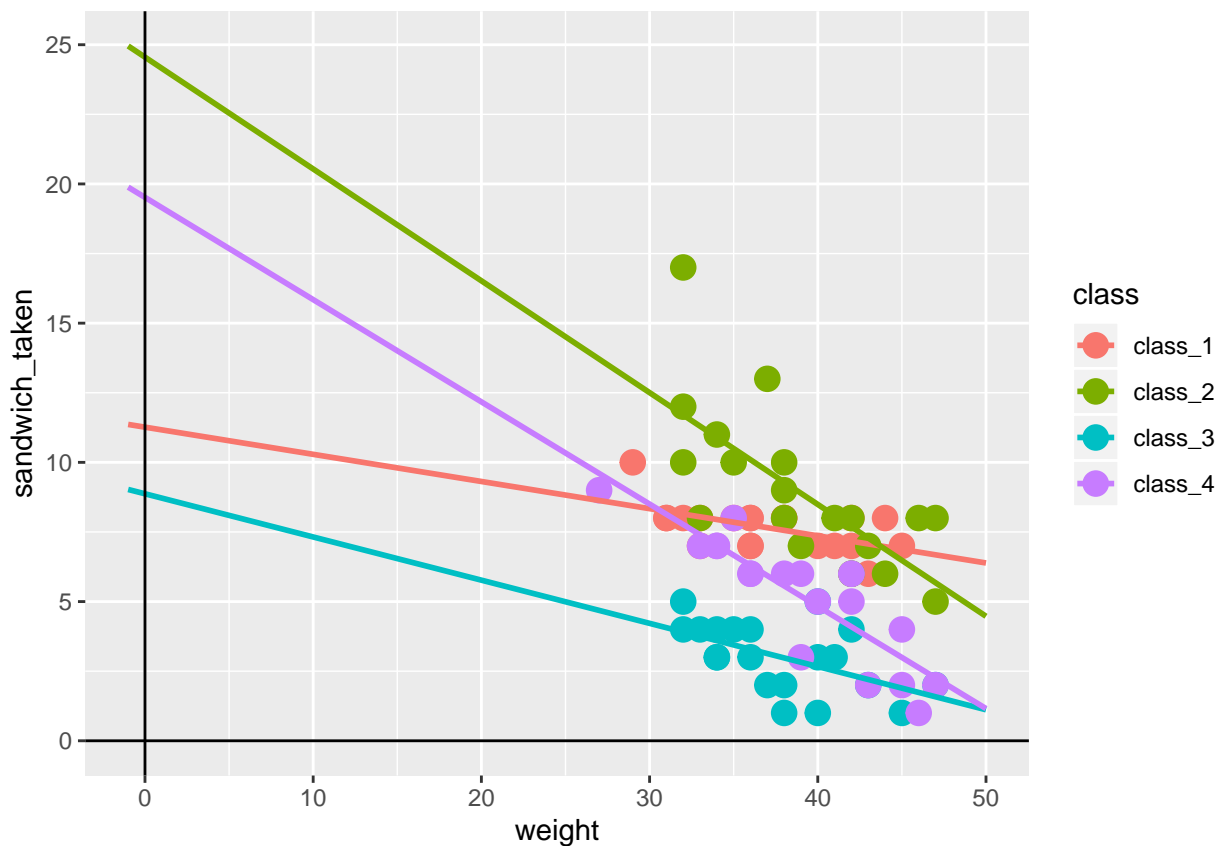
Az ábrán jól látszik, az osztályok nem csak abban különböznek, hogy a hozzájuk tartozó regressziós egyenes hol metszi az Y tengelyt, de **a regressziós egyenesek meredeksége is különbözik**.

Peldaul a `class_1`-ben a testsúly hatása elhanyagolhatónak tűnik abból a szempontból, hogy kitől mennyi szendvicset vesznek el, míg a `class_2`-ben és `class_4`-ben a testsúly hatása számottevő.

```
slope_plot = data_bully_slope %>% ggplot() + aes(y = sandwich_taken,
  x = weight, color = class) + geom_point(size = 4) + geom_smooth(method = "lm",
  se = F, fullrange = TRUE) + xlim(-1, 50) + geom_hline(yintercept = 0) +
  geom_vline(xintercept = 0)
slope_plot
```

```
## Warning: Removed 1 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```



3.4 Kevert modellek felepítése az R-ben

Az alábbi példa bemutatja, hogy hogyan lehet a random hatásokat beépíteni a modellekbe.

Három modellt fogunk építeni. Először egy szimpla fix hatásokat tartalmazó modellt, majd egy **random intercept modellt**, és egy **random slope modellt**.

A fenti ábra alapján arra lehet következtetni, hogy a `data_bully_slope` adatbázisban az iskolai osztály egy olyan random hatás, ami mind a regressziós egyenes intercept-jét, mind a meredekséget (slope) befolyásolja.

Ezert normalis esetben csak a random slope modellt illesztünk. A többi modell csak demonstrációs célból építjük, hogy összehasonlítsuk azok formuláit és bejósoló erejét a random slope modellel.

Eloszor építünk egy egyszerű regressziós modellt, melyben egyetlen fix hatás prediktor van: `weight`. Ezt a modellt a `mod_fixed` objektumba mnetjük.

egyszerű regressziós modell (csak fix hatás)

```
mod_fixed = lm(sandwich_taken ~ weight, data = data_bully_slope)
```

random intercept modell (a random intercept megengedett, de a random slope nem)

A kevert modellek formulája nagyon hasonló a csak fix hatást tartalmazó modellekehez, de az `lm()` függő helyett az `lmer()` függőt használjuk.

A random intercept random hatást a **“+ (1|class)” hozzáadásával** tehetjük a modellbe.

Ez gyakorlatilag azt jelenti, hogy megengedjük a modellnek hogy **külön regressziós egyenest illesszen minden klaszterre** (a mi esetünkben minden iskolai osztályra), de azt meghatározzuk, hogy **minden regressziós egyenesnek ugyan olyan legyen a meredeksége**.

Ezt normalis esetben akkor tennénk, ha azt gyanítanánk hogy az osztályok között nincs lenyegi eltérés a fix hatásokban, csak a kimeneti változó átlagos szintjében. Ez a modell jól illeszkedne a `data_bully_int` adatbázisra, de a fenti ábra alapján azt várjuk hogy a `data_bully_slope` adatbázisra kevesbé jól illeszkedik majd.

```
mod_rnd_int = lmer(sandwich_taken ~ weight + (1 | class), data = data_bully_slope)
```

random slope modell (mind a random intercept, mind a random slope megengedett):

Ennek a modellnek a formulája szinte teljesen megegyezik a random intercept modellel, egyedül abban különbözik, hogy a random hatásról szóló részben **“+ (1|class)”** helyett **“+ (weight|class)”** szerepel. Ez arra utal, hogy a class random hatás nem csak az interceptre, hanem a `weight` prediktor hatására is kiterjed.

Ezzel megengedjük a modellünknek, hogy **külön regressziós egyenest illesszen minden klaszterre**, és hogy azoknak **mind** az Y tengellyel való metszéspontja (**intercept**), **mind a meredeksége (slope) különbözhet**.

```
mod_rnd_slope = lmer(sandwich_taken ~ weight + (weight | class),
  data = data_bully_slope)
```

```
## Warning in checkConv(attr(opt, "derivs"), opt$par, ctrl = control$checkConv, :
## Model failed to converge with max|grad| = 0.00317298 (tol = 0.002, component 1)
```

3.5 Melyik modell reprezentálja legjobban a valóságot?

Hogyan döntjük el hogy **melyik modellt használjuk**? Ahogy korábban is láthattuk, a modellválasztásnál mindig **az elméletileg leginkább megalapozott** modellt érdemes választani. Ha van okunk feltételezni hogy egy hatás különbözik lesz a különböző klaszterekben, akkor használjuk a random slope modellt. Ha elméleti alapon inkább úgy ítéljük, hogy a fix hatások valószínűleg allandoak a csoportok között, illesszünk random intercept modellt.

Ennek ellenére van olyan eset, **amikor elméletileg mindket eshetőség elképzelhető**. Ilyen esetben hagyatkozhatunk a **vizualizációra** és a **modellilleszkedési mutatókra**, hogy eldöntsük, melyik modellt érdemesebb használni.

3.5.1 Random hatások vizualizációja

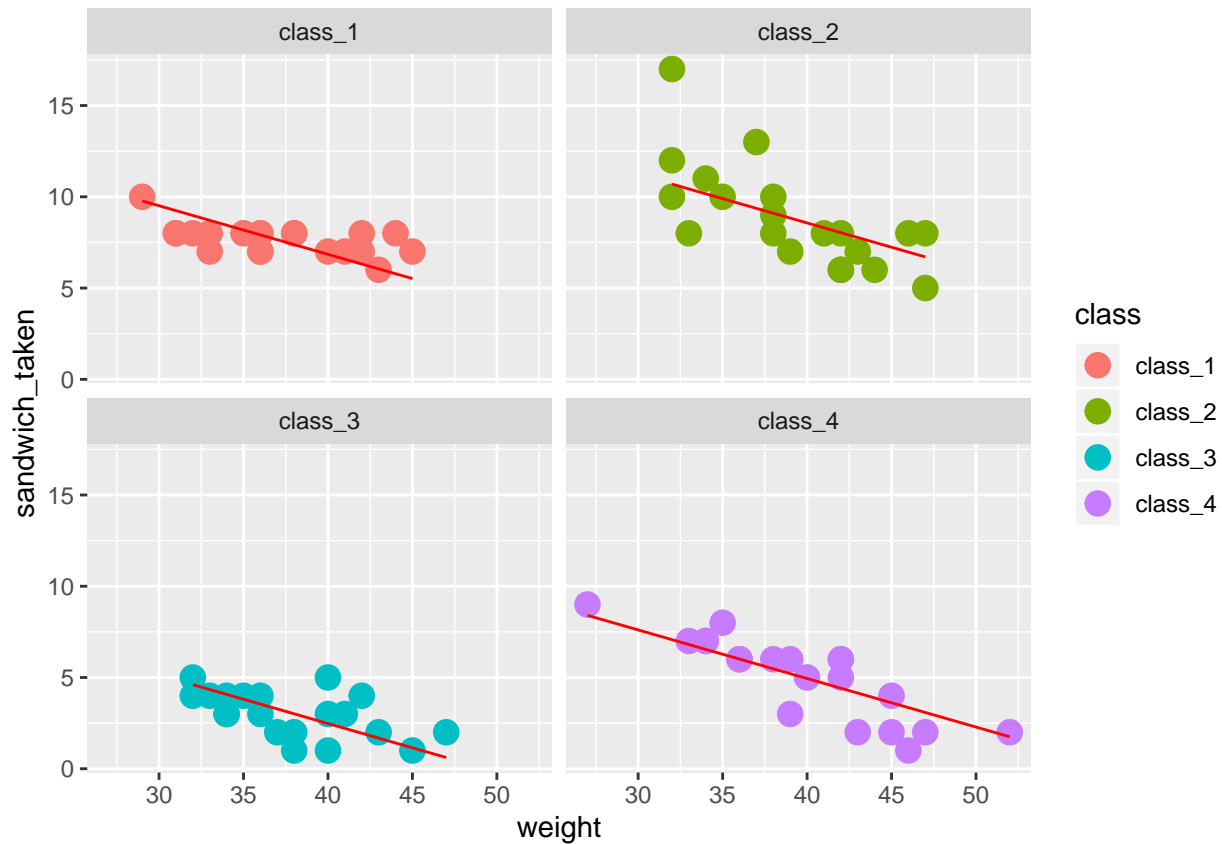
A random hatások explorációja esetén a vizualizáció kulcsszerepet tölt be.

Eloszor érdemes elmentenünk az intercept és a slope modellek által **bejósolt értékeket új változókba** (alább a `pred_int` és `pred_slope` változókba mentjük ezeket). Az eredeti adatbázisból származó predikciókat a `predict()` funkcióval nyerhetjük ki.

```
data_bully_slope = data_bully_slope %>% mutate(pred_int = predict(mod_rnd_int),
  pred_slope = predict(mod_rnd_slope))
```

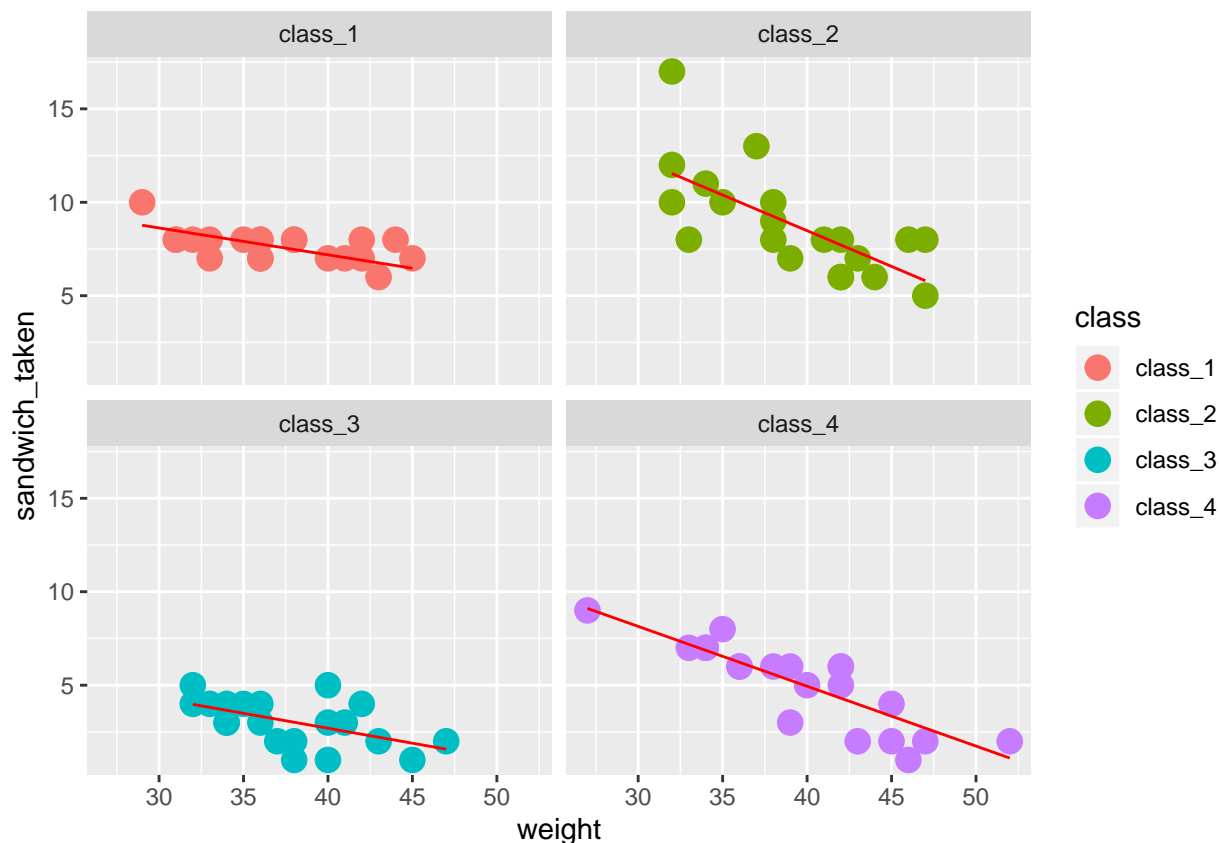
Igy vizualizáljuk a random **intercept modell** predikcióit:

```
data_bully_slope %>% ggplot() + aes(y = sandwich_taken, x = weight,
  group = class) + geom_point(aes(color = class), size = 4) +
  geom_line(color = "red", aes(y = pred_int, x = weight)) +
  facet_wrap(~class, ncol = 2)
```



Igy pedig a random **slope modell** predikcióit:

```
data_bully_slope %>% ggplot() + aes(y = sandwich_taken, x = weight,
  group = class) + geom_point(aes(color = class), size = 4) +
  geom_line(color = "red", aes(y = pred_slope, x = weight)) +
  facet_wrap(~class, ncol = 2)
```



Az ábrákon azt kell megvizsgálnunk, hogy a pontok mintázata kelle-e különbözni a csoportok között, hogy arra engedjen következtetni, hogy csoportonként különbözik a fix hatás.

Mivel a modellek által generált regressziós egyeneseket is tartalmazzák az ábrák, megtehetjük, hogy megvizsgáljuk, mi a hatása annak, hogy megengedjük a regressziós egyenes meredekségének változását a random intercept modellhez képest, ahol ez nincs megengedve. Az illesztés (szinte) mindig jobb lesz a random slope modellben. Ez szükséges, hiszen a modell flexibilisebb, több szabadsága van, ezért közelebb tud helyezkedni a pontokhoz minden csoportban. De **ha az illesztéskésébeli különbség a két modell között nem számottevő, biztonságosabb a random intercept modellnel maradni, hogy elkerüljük a túlillesztést**, ha csak az elmélet nem támogatja egyértelműen a random slope modellt.

3.5.2 Reziduais hiba összehasonlítása (ezt nem használjuk a gyakorlatban)

Talán első ránézésre csabítonak tűnhet, hogy egyszerűen arra hagyatkozunk a modellválasztás során, hogy melyik modell produkálta a legkevesebb reziduais hibát.

Ha összehasonlítjuk a három modell **reziduais hibáját** (residual sum of squares - RSS), láthatjuk, hogy a csak fix hatást tartalmazó modell használatakor marad a legtöbb hiba, a random intercept modell a második, és a **legkevesebb hibát a random slope modell esetén találjuk**.

```
sum(residuals(mod_fixed)^2)
```

```
## [1] 581.6364
```

```
sum(residuals(mod_rnd_int)^2)
```

```
## [1] 159.5818
```

```
sum(residuals(mod_rnd_slope)^2)
```

```
## [1] 132.2322
```

De ez nem igazan meglepo, hiszen a modell komplexitasa es ezzel **flexibilitasa egyre nott**, es errol tudjuk, hogy csokkenti a hibát azon az adatbazison amin a modellt epítettük, de a flexibilitas novelese miatt ez **tulilleszteshez** vezethet, ami új adatokon rosszabb bejoslasi hatekonysaghoz vezet. Ezert a nyers rezidualis hiba osszehasonlitas helyett olyan modell-illeszkedesi mutatohoz kell fordulnunk, amik korrigalva vannak a flexibilitasara (ezt ugy is mondhatjuk hogy a modell parameterek szamara.)

3.5.3 conditional AIC

Az egyik olyan modell-illeszkedesi mutato, amely korrigal a modell parameterek szamara az AIC. A kevert modellekhez egy specialis AIC mutato-t szamitunk ki, a **cAIC** mutatot (ami a conditional AIC roviditese). A cAIC-t megkaphatjuk peldaul a cAIC4 package cAIC() funkcioja segitsegevel.

```
AIC(mod_fixed)
```

```
## [1] 391.7357
```

```
cAIC(mod_rnd_int)$caic
```

```
## [1] 294.4023
```

```
cAIC(mod_rnd_slope)$caic
```

```
## [1] 285.6445
```

Ahogy korábban is lattuk az AIC eseten, ha az egyik modell cAIC mutatoja legalabb 2-vel alacsonyabb mint a masik modellhez tartozo cAIC, akkor azt mondhatjuk az alacsonyabb cAIC mutatoval biro modell szignifikansan jobban illeszkedik az adatokhoz.

3.5.4 likelihood ratio test

A masik bevett mod a modellek osszehasonlitasara a likelihood ratio test. Ez a modszer kevesbe elfogadott manapsag, de meg mindig sokan hasznaljak a szakirodalomban.

Ezt csakugy mint a nem kevert modelleknel, a kevert modelleknel is az anova() funkcioval vegezgetjuk el. Fontos, hogy ezt a likelihood ratio test-et csak a beagyazott modellek (nested models) eseten hasznalhatjuk (lasd a modell-osszehasonlitas gyakorlatot).

Az anova() funkcio hasznalatkor egy figyelmeztetést kapunk: ‘refitting model(s) with ML (instead of REML)’. ez aztert van mert a likelihood ratio test csak a Maximum likelihood (ML) becslessel dolgozo modellek osszehasonlitasara alkalmas. Viszont a kevert modelleket alapertelmezett modon a Restricted maximum likelihood (REML) becslessel dolgozunk, mert ez kevert modelleknel jobb becsleshez vezet. Ennek ellenere a REML es az ML becsleseket hasznalo modellek altalaban nagyon hasonloak egymashoz, ezert ezt a figyelmeztetést legtobbszor figyelmen kívül hagyható.

```
anova(mod_rnd_int, mod_rnd_slope)
```

```
## refitting model(s) with ML (instead of REML)
```

```
## Data: data_bully_slope
```

```
## Models:
```

```
## mod_rnd_int: sandwich_taken ~ weight + (1 | class)
```

```
## mod_rnd_slope: sandwich_taken ~ weight + (weight | class)
```

```
##           Df    AIC    BIC logLik deviance  Chisq Chi Df Pr(>Chisq)
## mod_rnd_int    4 310.00 319.53 -151.00   302.00
## mod_rnd_slope  6 307.94 322.23 -147.97   295.94 6.0595      2    0.04833 *
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

3.6 Mit kell kozolni az elemzesrol

A kozlendo informaciok nagyon hasonloak ahhoz, amit a fix hatas modellek eseten kozoltunk.

3.6.1 A statisztikai modszer leirasa:

“Ahhoz hogy a bullyzassal szembeni serulekenyseget meghatározzuk, egy **kevert linearis modellt illesztet-tunk**. A kevert modellben az elvett szendvicsek szamat mint **kimeneti valtozot** a testsullyal mint **fix hatasu prediktorral** jósoltuk be. A modellben ezen felul az iskolai osztaly **random hatasat** modelleztuk. Epitettunk mind egy **random slope** es egy **random intercept modellt**. Ahogy ezt a kutatasi tervunkben meghatároztuk, a ket modellt **összehasonlitottuk a cAIC** modellilleszkedeis mutatojuk alapjan, es ez alapjan határoztuk meg, melyik lesz a vegso bejoslo modellunk.”

A kovetkezo funkciokkal kapnank meg a kutatasi jelenteshez szukseges eredmenyeket:

cAIC:

```
cAIC(mod_rnd_int)$caic
```

```
## [1] 294.4023
```

```
cAIC(mod_rnd_slope)$caic
```

```
## [1] 285.6445
```

anova:

```
anova(mod_rnd_int, mod_rnd_slope)
```

```
## refitting model(s) with ML (instead of REML)
```

```
## Data: data_bully_slope
```

```
## Models:
```

```
## mod_rnd_int: sandwich_taken ~ weight + (1 | class)
```

```
## mod_rnd_slope: sandwich_taken ~ weight + (weight | class)
```

```
##           Df      AIC      BIC logLik deviance  Chisq Chi Df Pr(>Chisq)
```

```
## mod_rnd_int    4 310.00 319.53 -151.00   302.00
```

```
## mod_rnd_slope  6 307.94 322.23 -147.97   295.94 6.0595      2    0.04833 *
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

3.6.2 A teljes modell illeszkedesenek jellemzese

Az `r2beta()` funkcio kiszamitja a “marginalis R^2 ” mutatot Nakagawa, Johnson & Schielzeth (2017) cikkenek ajanlasi alapjan. Ez az R^2 mutato specialis fajtaja, ami azt mutatja meg, hogy mekkora a modell fix hatasu prediktorai által megmagyarazott varianciaarany. Ezt az R^2 mutatot erdemes hasznalni a modell bejoslo hatekonysaganak megadasara, hiszen a random hatasu prediktorokat uj adatokon nem tudjuk majd hasznalni bejoslasra.

Nincs egy klasszikus F-test aminek az eredmenyet fel lehetne hasznalni annak ertekelesere, hogy a teljes modell szignifikansen jobb bejoslast eredmenyez-e a null-modellnel, de az `r2beta` megadja a 95%-s konfidencia intervallumot, amit felhasznalhatunk szignifikanciatesztelesre. Ahogy korabban is, ha a konfidencia intervallim tartalmazza a 0-t, akkor a modell nem szignifikansen kulonbozik a null modelltol bejoslo hatekonysag tekinteteben.

Ezen felul mind a marginalis mind a kondicionalis R^2 erteket megkaphatjuk az `r.squaredGLMM()` funkcio hasznalataval a MuMIn package-bol. Ez a funkcio szinten a Nakagawa, Johnson & Schielzeth (2017) által publikalt formulat hasznalja ezen ertekek kiszamitasahoz.

Hivatkozás: Nakagawa, S., Johnson, P.C.D., Schielzeth, H. (2017) The coefficient of determination R^2 and intraclass correlation coefficient from generalized linear mixed-effects models revisited and expanded. J. R. Soc. Interface 14: 20170213.

```
# marginal R squared with confidence intervals
r2beta(mod_rnd_slope, method = "nsj", data = data_bully_slope)
```

```
## Effect Rsq upper.CL lower.CL
## 1 Model 0.147 0.304 0.036
## 2 weight 0.147 0.304 0.036
```

```
# marginal and conditional R squared values
r.squaredGLMM(mod_rnd_slope)
```

```
## Warning: 'r.squaredGLMM' now calculates a revised statistic. See the help page.
```

```
## R2m R2c
## [1,] 0.147134 0.8331627
```

Az eredmények reszben így írhatjuk le az eredményeket:

“A random slope modell jobb modell-illeszkedéshez vezetett mint a random intercept modell mind a likelihood ratio test ($X^2 = 6.06$, $df =$, $p = .048$) mind a cAIC alapján (cAIC intercept = 294.4, cAIC slope = 285.64). Ezért az alábbiakban a random slope modell eredményeit közöljük.

A kevert lineáris modell szignifikánsan jobb volt mint a null modell. A modellben a fix hatasu prediktorok az elvett szendvicsek varianciajának 14.7%-at magyaráztak meg ($R^2 = 0.15$ [95% CI = 0.04, 0.3]).”

3.6.3 Regressziós együtthatók közlése

Ezen felül a prediktorokhoz tartozó regressziós együtthatókról is közölnünk kell az eredményeket. Ezt a korábbiakhoz hasonlóan egy táblázatban szoktuk megtenni, ami minden prediktorra külön sorban közli az adatokat. (Itt csak egy fix hatasu prediktor van, szóval csak két sor lesz a táblázatban, egy az intercept-nek és egy a testsúly prediktornak.)

A végső táblázat valahogy így néz majd ki:

```
## b 95%CI lb 95%CI ub Std.Beta p-value
## (Intercept) 15.89 8.02 23.72 0 .021
## weight -0.25 -0.41 -0.09 -0.42 .04
```

A táblázat egyes elemeit itt találhatod meg:

Regressziós együtthatók és a hozzájuk tartozó p-értékek a `summary()` függvénnyel kaphatók meg (csak akkor fog p-értéket kiadni a `summary` függvény a kevert modellekre, ha az `lmerTest` package be van töltve)

```
summary(mod_rnd_slope)
```

```
## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula: sandwich_taken ~ weight + (weight | class)
## Data: data_bully_slope
##
## REML criterion at convergence: 297.4
##
## Scaled residuals:
## Min 1Q Median 3Q Max
## -2.3403 -0.5630 -0.0076 0.3896 4.0511
##
## Random effects:
## Groups Name Variance Std.Dev. Corr
```

```
## class      (Intercept) 44.61736 6.6796
##           weight      0.01693 0.1301   -0.94
## Residual              1.81799 1.3483
## Number of obs: 80, groups: class, 4
##
## Fixed effects:
##           Estimate Std. Error      df t value Pr(>|t|)
## (Intercept) 15.88887    3.55290  2.95835   4.472  0.0215 *
## weight      -0.25155    0.07224  2.98048  -3.482  0.0404 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##      (Intr)
## weight -0.940
## convergence code: 0
## Model failed to converge with max|grad| = 0.00317298 (tol = 0.002, component 1)
```

A regressziós egyutthatokhoz tartozó konfidencia intervallumok: (ez a funkció sokaig fut, mert sok iterációt vegez a kiszámításhoz)

```
confint(mod_rnd_slope)
```

A standardizált beta értékeket pedig a `stdCoef.merMod()` saját funkcióval lehet kinyerni:

```
stdCoef.merMod(mod_rnd_slope)
```

```
##           stdcoef      stdse
## (Intercept)  0.000000 0.0000000
## weight      -0.422144 0.1212255
```