

PSZB17-210 - Seminar_4

Zoltan Kekecs

Oct 02, 2024

4. Ora - Adatexploracio

Az ora celja az adatexploracios modszerek elsajatitasa.

Package-ek betoltese

A kovetkezo package-ekre lesz szuksegunk

```
if (!require("gridExtra")) install.packages("gridExtra")
library(gridExtra) # for grid.arrange
if (!require("psych")) install.packages("psych")
library(psych) # for describe
if (!require("tidyverse")) install.packages("tidyverse")
library(tidyverse) # for dplyr and ggplot2
```

Adatok betoltese

Beolvassuk a WHO által legutobb feltoltott COVID-19 adatokat a `read_csv()` funkcioval, es elmentjuk egy `my_data` nevu objektumba. A `read_csv()` funkcio a tidyverse resze, es egybol tibble formatumban menti el az adatainkat.

```
my_data <- read_csv("https://raw.githubusercontent.com/kekecsz/PSZB17-210-Data-analysis-seminar/refs/heads/master/seminar_04/StudentPerformanceFactors.csv")
```

Adatok attekintese

Mindig erdemes azzal kezdeni, hogy **megismerkedunk az adat** szerkezetével es tartalmával.

A **tibble objektum** meghivasaval kaphatunk nemi informaciot az adattabla szerkezetéről. Lathatjuk hany sor es hany oszlop van az adattablában, es lathatjuk milyen class-ba tartoznak (chr, dbl ...)

```
my_data
```

```
## # A tibble: 6,607 × 20
##   Hours_Studied Attendance Parental_Involvement Access_to_Resources
##   <dbl>         <dbl> <chr>                <chr>
## 1           23           84 Low                    High
## 2           19           64 Low                    Medium
## 3           24           98 Medium                 Medium
## 4           29           89 Low                    Medium
## 5           19           92 Medium                 Medium
## 6           19           88 Medium                 Medium
## 7           29           84 Medium                 Low
## 8           25           78 Low                    High
## 9           17           94 Medium                 High
## 10          23           98 Medium                 Medium
## # i 6,597 more rows
## # i 16 more variables: Extracurricular_Activities <chr>, Sleep_Hours <dbl>,
## #   Previous_Scores <dbl>, Motivation_Level <chr>, Internet_Access <chr>,
## #   Tutoring_Sessions <dbl>, Family_Income <chr>, Teacher_Quality <chr>,
## #   School_Type <chr>, Peer_Influence <chr>, Physical_Activity <dbl>,
## #   Learning_Disabilities <chr>, Parental_Education_Level <chr>,
## #   Distance_from_Home <chr>, Gender <chr>, Exam_Score <dbl>
```

Leiro statisztikak

Ha az egyes valtozok **leiro statisztikaira** (descriptive statistics) vagyunk kíváncsiak, kerhetjuk ezt a mar tanult modon.

Peldaul lekerhetjuk a valtozo alapveto legalacsonyabb es legmagasabb erteket, atlagat, medianjat, a kvartiliseket, es hogy hany hianyzo adat van (ha van) a **summary()** funkcioval (miutan a select funkcioval kivlasztottuk, melyik valtozora vagyunk kivancsiak)

```
my_data %>%
  select(Attendance) %>%
  summary()
```

```
##      Attendance
##  Min.   : 60.00
##  1st Qu.: 70.00
##  Median : 80.00
##  Mean   : 79.98
##  3rd Qu.: 90.00
##  Max.   :100.00
```

Vagy megkapthatjuk ugyanezt az osszes valtozora, ha ugyanezt az egesz adattablara futtatjuk le. Persze a karakter osztalyba tartozo valtozoknal mindezeknek a leiro statisztikaknak nincs ertelme, ott csak a class informaciot kaptjuk az output-ban.

```
my_data %>%
  summary()
```

Gyakorlas

- Mi volt a legalacsonyabb részvételi arány (*Attendance*)?
- Mi az átlagos alvásmennyiség (*Sleep_Hours*) azok között a hallgatók között, akiknek alacsony a hozzáférése az erőforrásokhoz (*Access_to_Resources* == "Low")?

Megtobb leiro statisztika

A **Psych** package segitsegevel a **describe()** funkcio megtobb hasznos informaciot adhat. Ez a funkcio elsosorban szam-valtozok leirasara szolgal, es karakter tipusu kategorikus valtozok eseten sok warning message-et ad, ezert erdemes a funkciot csak a szam-valtozokra lefuttatni. Itt harom ilyen valtozot valsztok ki a select() funkcioval.

```
my_data %>%
  select(Hours_Studied, Attendance, Exam_Score) %>%
  describe()
```

```
##           vars      n mean    sd median trimmed   mad min max range skew
## Hours_Studied    1 6607 19.98  5.99     20    19.97  5.93    1  44    43 0.01
## Attendance       2 6607 79.98 11.55     80    79.97 14.83   60 100    40 0.01
## Exam_Score       3 6607 67.24  3.89     67    67.11  2.97   55 101    46 1.64
##
##           kurtosis    se
## Hours_Studied     0.02 0.07
## Attendance       -1.19 0.14
## Exam_Score       10.56 0.05
```

Faktorok

Nehany karaktervaltozonak csak **korlatozott mennyisegu eleme** lehet, mint peldaul a Parental_Education_Level (ebben az adatbazisban csak College, High School, Postgraduate szinteket vesz fel). Ezeket megjelolhetjuk faktor (factor) osztalyu valtozokent, es akkor az R tobb informaciot fog adni rola.

```
table(my_data$Parental_Education_Level)
```

```
##
##      College  High School Postgraduate
##      1989      3223      1305
```

```
names(my_data)
```

```
## [1] "Hours_Studied"      "Attendance"
## [3] "Parental_Involvement" "Access_to_Resources"
## [5] "Extracurricular_Activities" "Sleep_Hours"
## [7] "Previous_Scores"      "Motivation_Level"
## [9] "Internet_Access"      "Tutoring_Sessions"
## [11] "Family_Income"        "Teacher_Quality"
## [13] "School_Type"          "Peer_Influence"
## [15] "Physical_Activity"     "Learning_Disabilities"
## [17] "Parental_Education_Level" "Distance_from_Home"
## [19] "Gender"               "Exam_Score"
```

```
my_data <- my_data %>%
  mutate(Parental_Education_Level = factor(Parental_Education_Level),
         Parental_Involvement = factor(Parental_Involvement),
         Access_to_Resources = factor(Access_to_Resources),
         Extracurricular_Activities = factor(Extracurricular_Activities),
         Motivation_Level = factor(Motivation_Level),
         Internet_Access = factor(Internet_Access),
         Family_Income = factor(Family_Income),
         Teacher_Quality = factor(Teacher_Quality),
         School_Type = factor(School_Type),
         Peer_Influence = factor(Peer_Influence),
         Learning_Disabilities = factor(Learning_Disabilities),
         Distance_from_Home = factor(Distance_from_Home),
         Gender = factor(Gender))
```

Miutan egy változót faktorként azonosítottunk, bizonyos funkciók képesek felhasználni ezt az információt.

Peldaul így már a fenti **summary()** funkció is kiadja az **egyes faktorszintekről** hogy hany megfigyeles tartozik az egyes kategoriakba (faktorszintekbe).

A **levels()** funkció megmutatja mik a faktorunk szintjei, de lathato ez akkor is ha csak meghivjuk a változót magát.

A **table()** funkció pedig tablazatot készít arról, hogy az egyes csoportokban hany megfigyeles talalhato. (A **table()** sima karakter változokkal is mukodik, nem csak faktorokkal)

Amikor kilistazzuk a faktor változót, akkor is kiírja az R a lista aljara, hogy milyen faktorszintek vannak.

```
levels(my_data$Parental_Education_Level)
```

```
## [1] "College"      "High School"  "Postgraduate"
```

```
table(my_data$School_Type)
```

```
##
##   Home Private  Public
##     6    2006   4595
```

```
my_data %>%
  select(Gender) %>%
  summary()
```

```
##      Gender
## Female:2793
## Male   :3814
```

```
my_data$School_Type
```

Van, hogy szeretnénk **kizárni** bizonyos **faktorszinteket** az elemzésből. Pl. ha valamelyik faktor szintből nagyon keves megfigyeles van, vagy csak a kutatási kérdésünk nem vonatkozik az adott részére a populációnak.

Az alábbi példában a **School_Type** változóból kizárjuk a "Home" szintet, vagyis azokat a válaszadókat akik otthoni iskolába járnak.

Ezt a már korábban tanult **filter()** funkció segítségével könnyedén megtehetjük, azonban arra figyelniünk kell, hogy az R megjegyzi a faktorszinteket, és azt azt követően is a **változohoz rendelve tartja**. A **faktorszintek meg akkor is megmaradnak ha nem marad egy megfigyeles sem** az adott faktorszinten az adattáblában.

```
levels(my_data$School_Type)
```

```
## [1] "Home" "Private" "Public"
```

```
my_data %>%
  filter(School_Type != "Home") %>%
  select(Exam_Score, School_Type) %>%
  summary()
```

```
##      Exam_Score      School_Type
## Min.   : 55.00   Home   :    0
## 1st Qu.: 65.00   Private:2006
## Median : 67.00   Public :4595
## Mean   : 67.23
## 3rd Qu.: 69.00
## Max.   :101.00
```

Igy ezeket a szinteket ejtethetjük a **droplevels()** funkcióval.

```
my_data_noHomeSchooled = my_data %>%
  filter(School_Type != "Home") %>%
  mutate(School_Type = droplevels(School_Type))
```

```
my_data_noHomeSchooled %>%
  select(Exam_Score, School_Type) %>%
  summary()
```

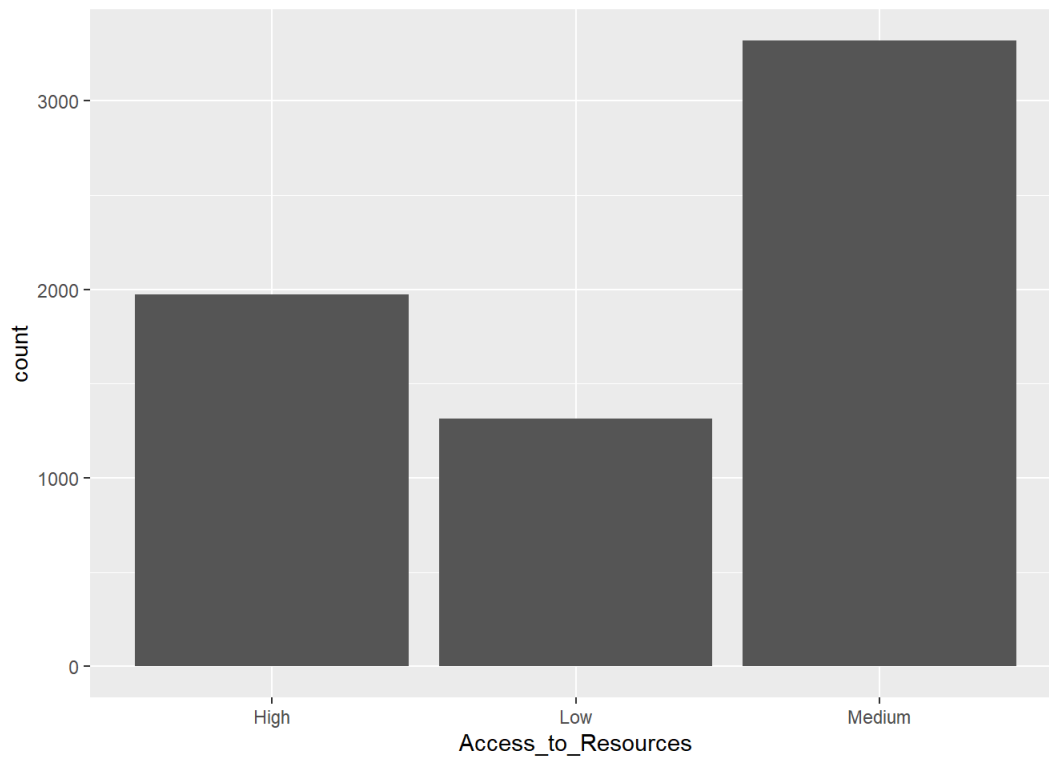
```
##      Exam_Score      School_Type
## Min.   : 55.00   Private:2006
## 1st Qu.: 65.00   Public :4595
## Median : 67.00
## Mean   : 67.23
## 3rd Qu.: 69.00
## Max.   :101.00
```

Faktorszintek egymashoz viszonyított értéke

Legtöbbször a faktorszintek között nincs “értékbeli” különbség, egyszerűen csoportnevekről van szó, de néha egy meghatározott reláció van közöttük, pl. a legmagasabb iskolai végzettség lehet végzettsége nélküli < általános iskolai < középiskolai < felsőfokú ... Itt faktorszinteknek van egy meghatározott hierarchiaja, vagy sorrendje. Ebben az adatbázisban sok ilyen változó van, pl. `Access_to_Resources`.

Amikor ábrát rajzolunk erről a változóról, láthatjuk hogy a faktorszintek sorrendje “High”, “Low”, “Medium”.

```
my_data %>%
  ggplot() +
  aes(x = Access_to_Resources) +
  geom_bar()
```

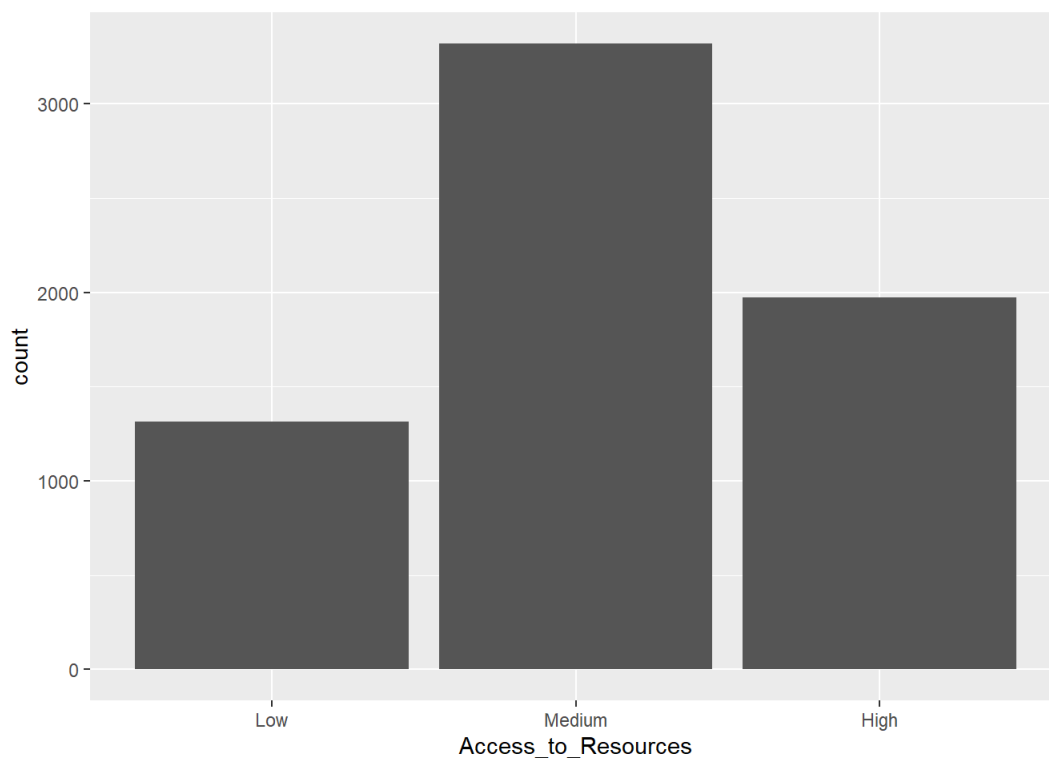


Ez nem feltétlenül intuitív ábrázolás, hiszen általában a kisebbtől a nagyobbig szoktunk haladni balról jobbra. De az R nem tudja mit jelentenek a faktorszintek nevei. A faktorszintek sorrendjének meghatározásánál ezért alapértelmezett módon **abc sorrendet használ**.

Specifikálhatjuk maskepp is a faktorszintek sorrendjét a factor funkcióban a **levels = c()** parameteren keresztül egy vektorban megadva.

```
my_data = my_data %>%
  mutate(Access_to_Resources = factor(Access_to_Resources, levels = c(
    "Low",
    "Medium",
    "High")))

my_data %>%
  ggplot() +
  aes(x = Access_to_Resources) +
  geom_bar()
```



Attól meg hogy megadjuk a levels-el a faktorszintek listázási sorrendjét, az R meg mindig egyenrangúként kezeli a faktorszinteket, csak most már jó sorrendben írja ki őket.

Ha azt szeretnénk ha az R úgy értékelne hogy a faktorszintek valamilyen hierarchikus sorrendben van, vagyis **ordinalis változokent**, akkor ezt a `factor()` funkcion belül az **ordered = T** paraméter beállításával tehetjük meg.

Ha ezt teszük, a faktor változó kilistázásakor reláció-jelek kerülnek a faktorszintek közé, és más funkciók is fel tudják majd használni ezt az információt.

```
my_data = my_data %>%
mutate(Access_to_Resources = factor(Access_to_Resources, ordered = T, levels = c(
  "Low",
  "Medium",
  "High")))

head(my_data$Access_to_Resources)
```

```
## [1] High    Medium Medium Medium Medium Medium
## Levels: Low < Medium < High
```

Kategorikus változó létrehozása és újrakodolása

Ha egy folytonos változó alapján szeretnénk egy kategorikus változót létrehozni, használhatjuk a **mutate()** és **case_when()** funkciók kombinációját hogy csináljunk egy új változót.

Mondjuk a vizsgán elért százalékok alapján hozzunk létre értékelesi csoportokat.

Ebbe a kódba beleépítettem a **factor()** funkciót is, hogy azonnal meghatározzuk, hogy ez az új változó egy faktor, és hogy ordinalis változó, hiszen a különböző szinteknek van érték-relációja.

```
my_data = my_data %>%
  mutate(Grade = factor(
    case_when(Exam_Score < 60 ~ "Poor",
              Exam_Score >= 60 & Exam_Score < 80 ~ "Good",
              Exam_Score > 80 ~ "Excellent"), levels = c("Poor", "Good", "Excellent"), ordered = T))

levels(my_data$Grade)
```

```
## [1] "Poor"      "Good"      "Excellent"
```

Egy másik funkció amivel manipulálhatjuk a faktorszinteket, a **recode()**. Ha kategorikus változókat szeretnénk átkodolni, mondjuk ha szeretnénk az ímént létrehozott Grade változó alapján egy újrakódolt új változót létrehozni, azt a következő képpen tehetjük:

```
my_data = my_data %>%
  mutate(Grade_passfail = factor(recode(Grade,
    "Poor" = "Failed",
    "Good" = "Passed",
    "Excellent" = "Passed"))))

levels(my_data$Grade_passfail)
```

```
## [1] "Failed" "Passed"
```

Gyakorlas

- szurd az adatokat úgy hogy ne legyenek benne az otthon tanuló (School_Type == "Home") hallgatók.
- csinálj egy új kategorikus változót (nevezzük ezt *Sleep_Categorical*-nak) a `mutate()` funkció használatával amiben azok az országok ahol a *Sleep_Hours* változó 6 alatt van "inadequate", ahol 6 vagy a felett van "adequate" kategóriába kerüljenek.
- figyeld oda hogy faktorként jelöld meg ezt az új változót (Ezt lehet az előző lépésben a `mutate()` funkcion belül, vagy egy külön lépésben, de mindenképpen a `factor()` vagy az `as.factor()` funkciókat érdemes hozzá használni)
- mentsd el ezt a változót az eredeti adatobjektumban úgy hogy később is lehessen vele dolgozni
- készíts egy táblázatot arról, hogy hányszor esnek a *Sleep_Categorical* egyes kategóriába.
- Add meg a faktorszintek helyes sorrendjét: az "inadequate" szint legyen előbb sorolva, mint az "adequate" szint (Írd fel a *Sleep_Categorical* korábbi változatát ezzel a változattal ahol a szintek már helyes sorrendben vannak, vagy ezt a sorrendezést is bele vonhatod az eredeti funkcióba, amivel a változót generáltad)
- Ellenőrizd, hogy valóban helyes sorrendben szerepelnek-e a faktor szintjei.

Exploracio vizualizacion keresztül

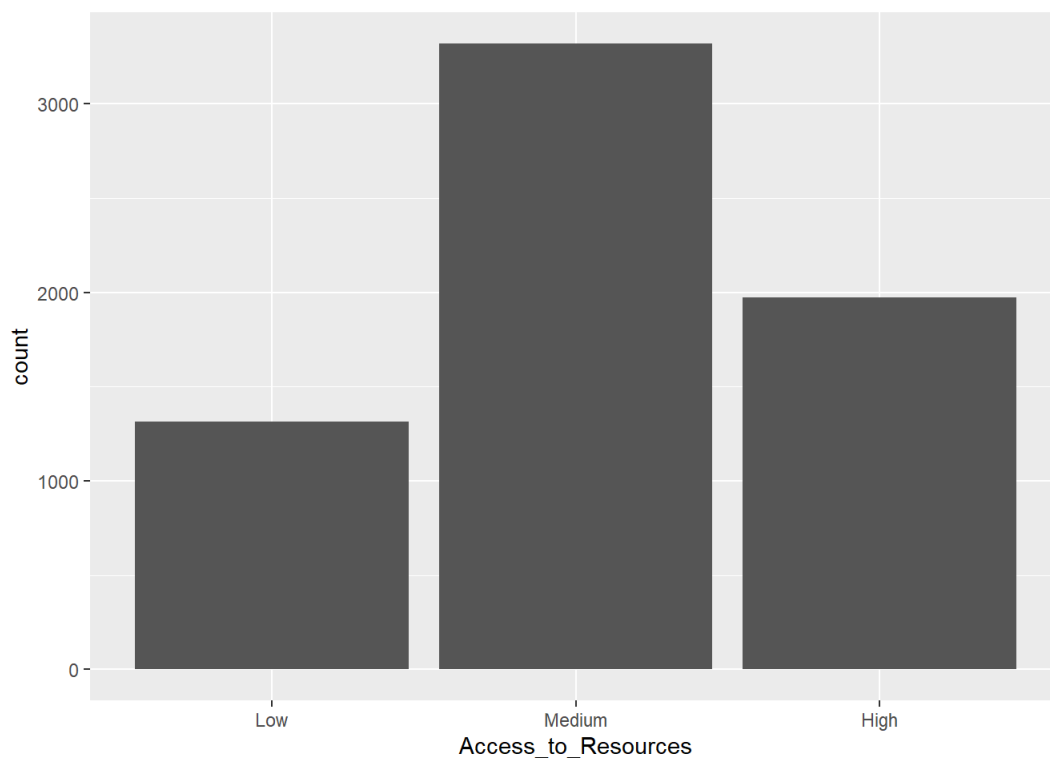
Az egyes változók vizualizációja és a leíró statisztikák átvizsgálása elengedhetetlen, hogy azonosítsuk az esetleges adatbeviteli **hibákat** és **egyéb nemvárt furcsaságokat** az adataink között.

Egyes változók vizualizációja

Az egyes változók például **abrak** segítségével megvizsgálhatók.

A **kategorikus** változókat gyakran oszlopdiagrammal (**geom_bar**) ábrázoljuk,

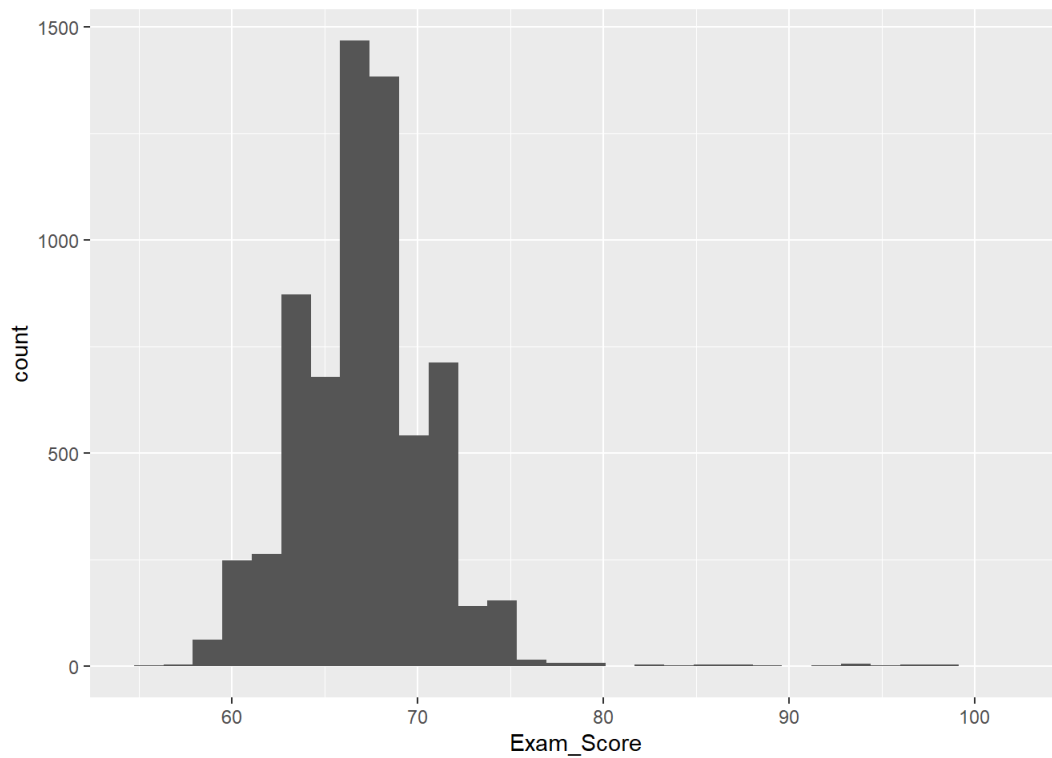
```
my_data %>%  
ggplot() +  
  aes(x = Access_to_Resources) +  
  geom_bar()
```



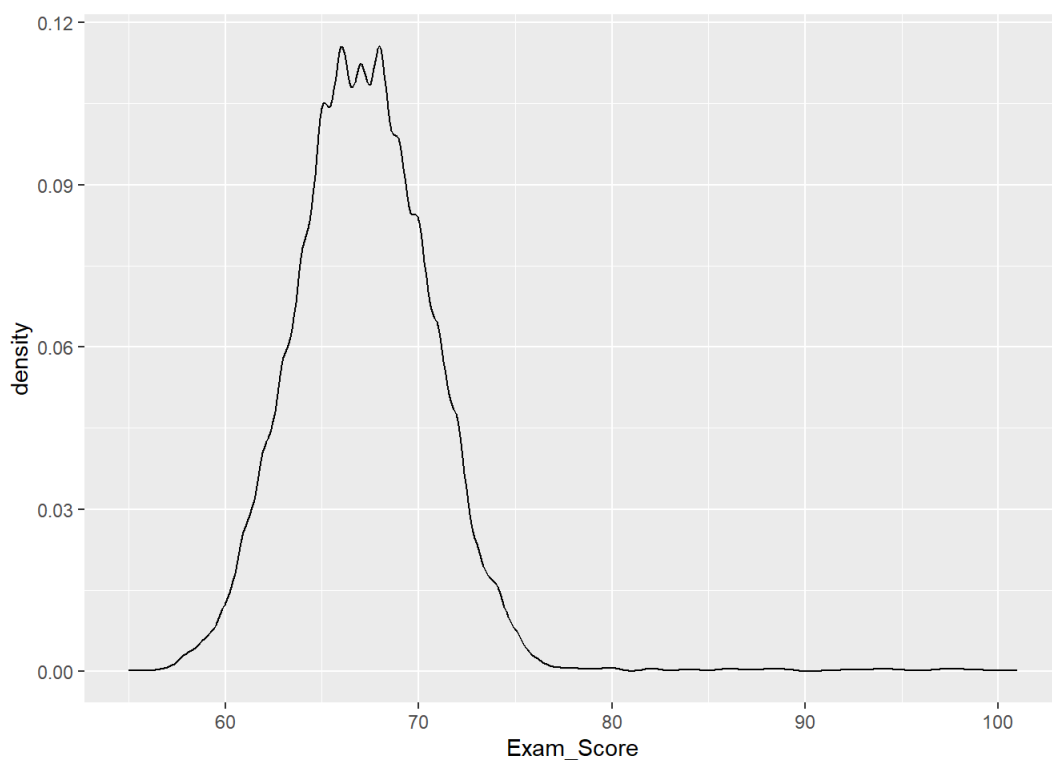
A **folytonos** változókat gyakran histogrammal (**geom_histogram**) vagy sűrűségábrával (**geom_density**) ábrázoljuk,

```
my_data %>%  
ggplot() +  
  aes(x = Exam_Score) +  
  geom_histogram()
```

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



```
my_data %>%  
ggplot() +  
  aes(x = Exam_Score) +  
  geom_density()
```



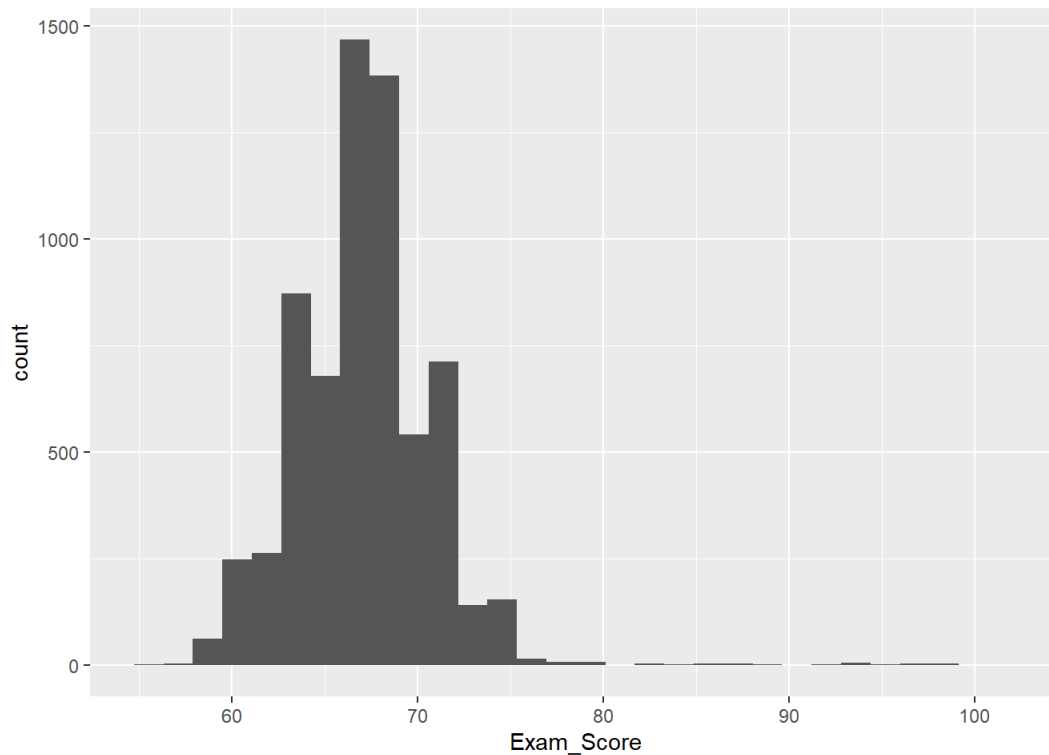
Hibaellenorzes

MINDING ellenorizd az adataidat mielőtt komolyabb adatelemzésbe kezdesz, hogy meggyozodj róla, hogy az adatok tisztak és megfelelnek az elvárásaidnak.

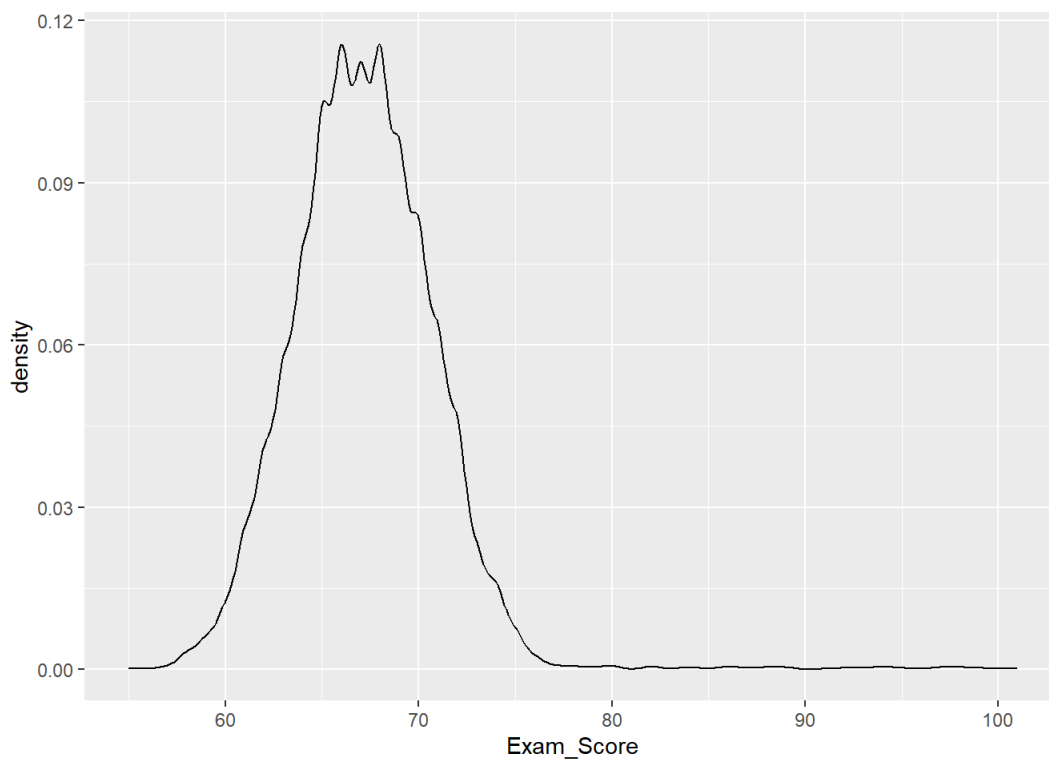
Ehhez használhatsz mind adatvizualizációt, mind a fentebb tanult leíró statisztikát (`summarize()`, `summary()`, `describe()` funkciókkal).


```
my_data %>%  
ggplot() +  
  aes(x = Exam_Score) +  
  geom_histogram()
```

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



```
my_data %>%  
ggplot() +  
  aes(x = Exam_Score) +  
  geom_density()
```



```
my_data %>%
  select(Exam_Score) %>%
  summary()
```

```
##      Exam_Score
##  Min.   : 55.00
## 1st Qu.: 65.00
##  Median : 67.00
##   Mean  : 67.24
## 3rd Qu.: 69.00
##   Max.  :101.00
```

Gyakorlas

Szurd az adatokat úgy hogy csak a Previous_Scores, Peer_Influence, Parental_Involvement, es a Sleep_Hours változokkal dolgozz.

Hasznald a fent tanult módszereket, hogy **azonosítsd a my_data adattablában levo hibakat** vagy nem vart furcsasagokat.

- A vizualizacian tul a View(), describe(), table(), es summary() funkciokat erdemes hasznalni az adatok elso attekintesere
- A numerikus (vagy eppen folytonos) változoknál vizsgál meg a minimum es maximum értéket es a hiányzó adatok mennyiséget, valamint az eloszlást, esetleg a felvett értékek mennyiséget, ha nincs túl sok felvehető érték.
- A kategorikus változoknál vizsgál meg az összes faktorszintet es az egyes szintekhez tartozó megfigyelesek mennyiséget.

A hibákat a következokeppen javíthatjuk.

A **mutate()** es a **replace()** funckciok hasznalataval **cserélhetünk ki** értékeket más értékekre. Azt, hogy ilyenkor hiányzó adatra (NA), vagy egy másik, valószínű értékre kell megváltoztatni az értéket, a szituációtól függ. Általában a biztosabb megoldás ha hiányzó adatnak jelöljük a kérdéses értéket (NA), de ez sok adatvesztéshez vezethet. Ha elég valószínű hogy mi a helyes érték, beírhatjuk, DE **minden javítást fel kell tüntetni** a kutatási jelentésben (es a ZH során is), hogy az olvasó számára tiszta legyen, hogy itt egy adathelyettesítés vagy kizárás történt!

Mindig érdemes a javított adatokat **új adattablába** elmenteni. A mi esetünkben a my_data_corrected nevet adtuk a javított objektumnak. Így a nyers adataink megmaradnak, ami hasznos lehet későbbi műveleteknel.

```
my_data %>%
  select(Exam_Score) %>%
  summary()
```

```
##      Exam_Score
##  Min.   : 55.00
## 1st Qu.: 65.00
##  Median : 67.00
##   Mean  : 67.24
## 3rd Qu.: 69.00
##   Max.  :101.00
```

```
my_data_corrected <- my_data %>%
  mutate(Exam_Score = replace(Exam_Score, Exam_Score=="101", NA))

my_data_corrected %>%
  select(Exam_Score) %>%
  summary()
```

```
##      Exam_Score
##  Min.   : 55.00
## 1st Qu.: 65.00
##  Median : 67.00
##   Mean  : 67.23
## 3rd Qu.: 69.00
##   Max.  :100.00
## NA's    :1
```

Erdemes **megbizonyosodni rola**, hogy az adatcsere sikeres volt, az új javított adat vizualizációjával, vagy a leíró statisztikák lekerdezésével.

Több változó kapcsolatának felterkepezése

Több változó kapcsolatot is felterkepezhetjük tablazatok és ábrák segítségével.

Két kategorikus (csoportosító) változó kapcsolatának felterkepezése

Feltáró elemzés

Most vizsgáljuk meg azt, hogy a család anyagi helyzete (*Family_Income*) milyen összefüggést mutat az interneteléréssel (*Internet_Access*).

A legegyszerűbb módja két csoportosító változó kapcsolatának megvizsgálására a két változó **kereszt-tablázatának (crosstab)** elkészítése a **table()** funkcióval.

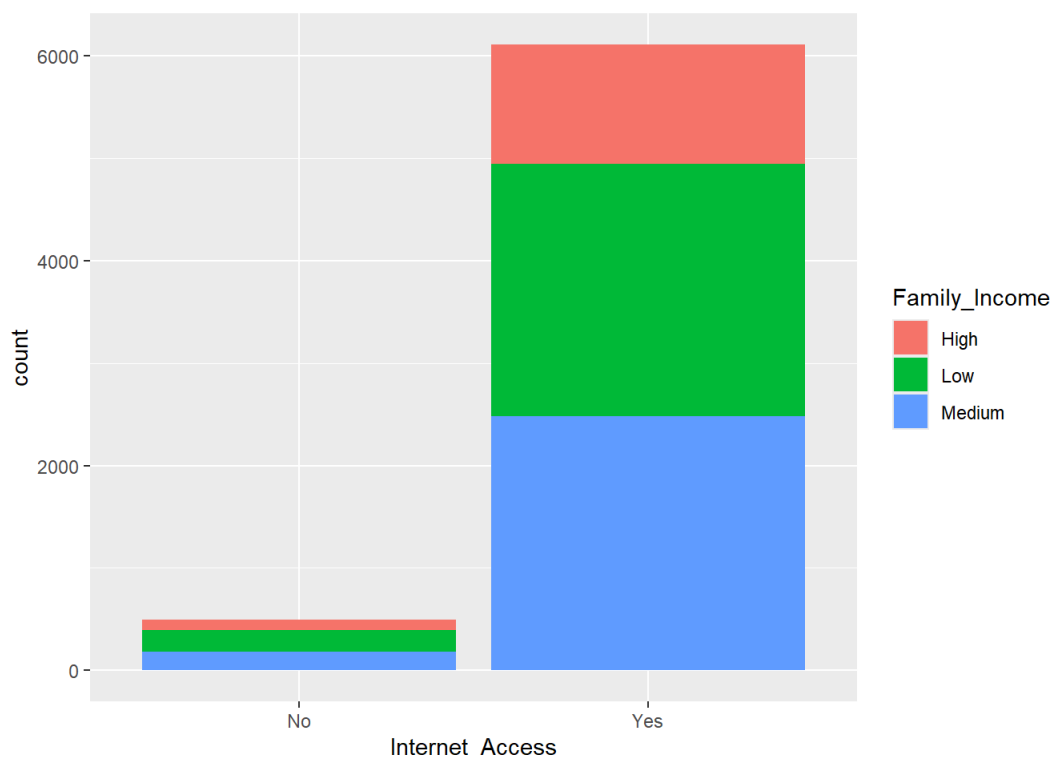
```
table(my_data_corrected$Family_Income, my_data_corrected$Internet_Access)
```

```
##
##           No  Yes
##  High    102 1167
##  Low     211 2461
##  Medium  186 2480
```

Sokszor ennél sokkal **szemleletesebb az ábrák (plot)** használata.

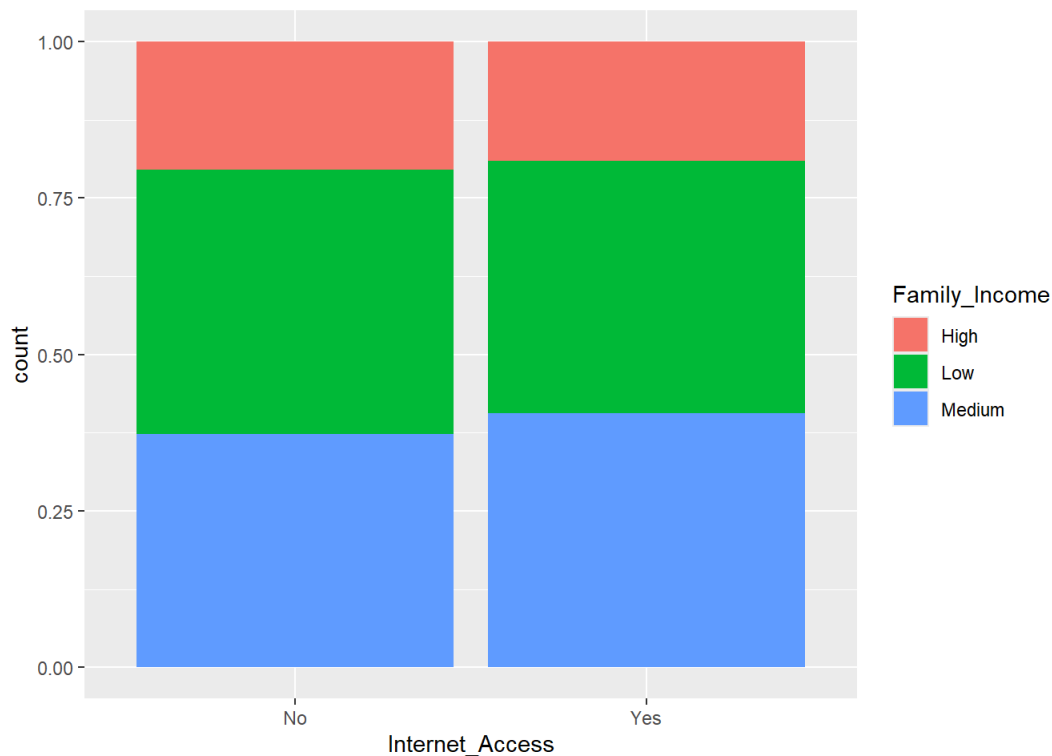
Erre az egyik lehetőség a **stacked bar chart** (egymásra tornyozott oszlopdiagram, a **geom_bar()** geomot használjuk) használata. Itt az egyik változó kategóriái adják meg hány oszlop lesz (ez a változó lesz az x tengelyen reprezentálva, így ezt az "x =" részen adhatjuk meg), a másik változó az oszlopokat színekkel szegmentálja, ezt pedig a "**fill** =" részen adhatjuk meg.

```
my_data_corrected %>%
  drop_na(Family_Income) %>%
  ggplot() +
    aes(x = Internet_Access, fill = Family_Income) +
    geom_bar()
```



Ha az egyes faktorszinteken nagyon **különböző mennyiségű megfigyelés** van, ez a megjelenítés néha felvezető következtetésekhez vezethet, így néha hasznosabb ha az oszlopok nem számosságot (count), hanem **reszarányt (proportion)** jelölnek. Ha ezt szeretnénk, ahelyett hogy üresen hagynánk a **geom_bar()** funkciót, a következőt adjuk meg: **geom_bar(position = "fill")**. Vagy használhatjuk az eltoló oszlopdiagramot (dodged barchart) (a **position = "dodge"** parameter megadásával a **geom_bar()** -on belül)

```
my_data_corrected %>%
  drop_na(Family_Income) %>%
  ggplot() +
    aes(x = Internet_Access, fill = Family_Income) +
    geom_bar(position = "fill")
```



Gyakorlas

Hasznald a fent tanult módszereket, hogy megvizsgald a `my_data_corrected` adatbázisban a **Sleep_Categorical** és a **Distance_from_Home** változók közötti összefüggést. - hasznalj **geom_bar()** geomot a megjelenítéshez - próbald meg mind a **szamossagot**, mind a **reszaranyt** kifejező ábrát megvizsgálni `geom_bar(position = "fill")` - milyen **kovetkeztetést** tudsz levonni az ábrákról?

a fenti gyakorlashoz a `Sleep_Categorical` változót így lehet legeneralni:

```
my_data_corrected = my_data_corrected %>%
  mutate(Sleep_Categorical = factor(
    case_when(Sleep_Hours < 6 ~ "inadequate",
              Sleep_Hours >= 6 ~ "adequate"), ordered = T, levels = c("inadequate",
"adequate")))

levels(my_data_corrected$Sleep_Categorical)
```

```
## [1] "inadequate" "adequate"
```

`geom_bar()` megjelenítésnél fontos hogy ha az egyes megfigyelesek **keves megfigyelesbol allnak**, az ábra megteveszto lehet, mert az ábra nem jelzi a megfigyelesek szamat és így azt, hogy milyen biztosak lehetünk az eredményben. Ilyen esetekben az egyik kategóriát ki lehet venni az ábráról, vagy a **szamossagot** és a **reszaranyt** ábrázoló ábrákat egymás mellett lehet bemutatni, hogy így kiegészítsek egymást. Ehhez használhatjuk a **grid.arrange()** funkciót.

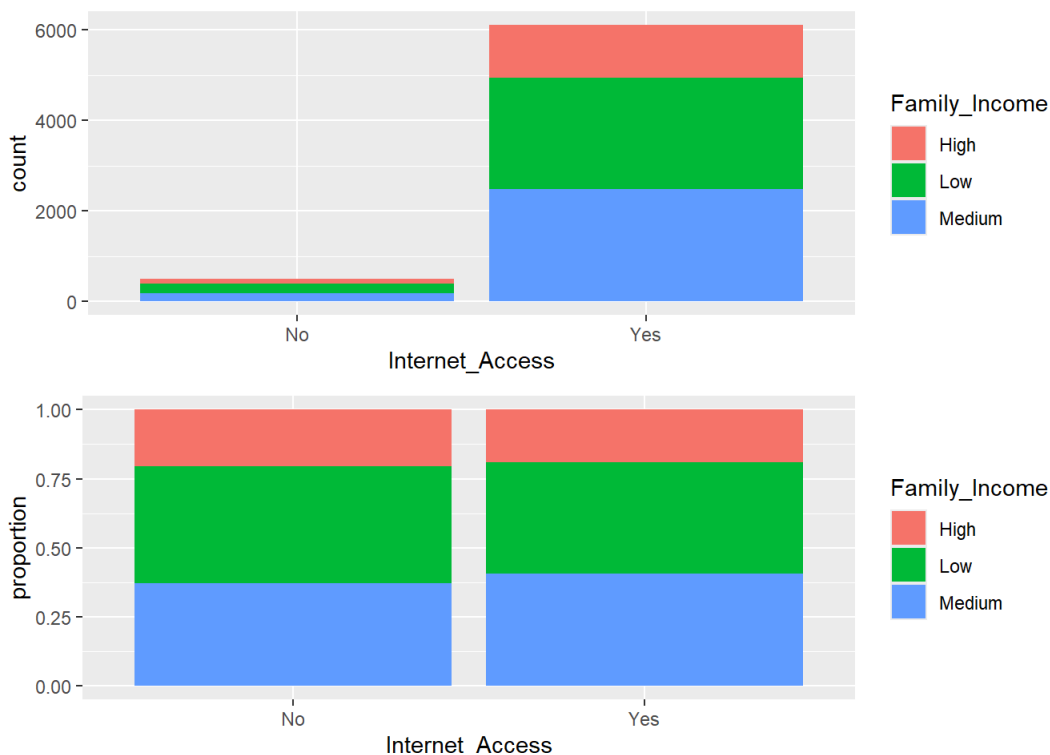
```

szamossag_plot <-
my_data_corrected %>%
  drop_na(Family_Income) %>%
ggplot() +
  aes(x = Internet_Access, fill = Family_Income) +
  geom_bar()

reszarany_plot <-
my_data_corrected %>%
  drop_na(Family_Income) %>%
ggplot() +
  aes(x = Internet_Access, fill = Family_Income) +
  geom_bar(position = "fill") +
  ylab("proportion")

grid.arrange(szamossag_plot, reszarany_plot, nrow=2)

```



A theme(legend.position) és a guides() funciók használatával kontrollálhatjuk hogy hol és hogyan jelenjen meg a **jelmagyarázat** az ábrán. Az ábra **interpretálhatósága** attól függően is **változhat**, hogy melyik változót tesszük az x-tengelyre és melyiket szintenként ábrázolva.

Az alábbi ábrák az egymillió főt vetített új esetek számanak kapcsolatát nézzük meg a gdp-vel. Mindkét változó esetén a csoportosított változót (`_kat`) használjuk.

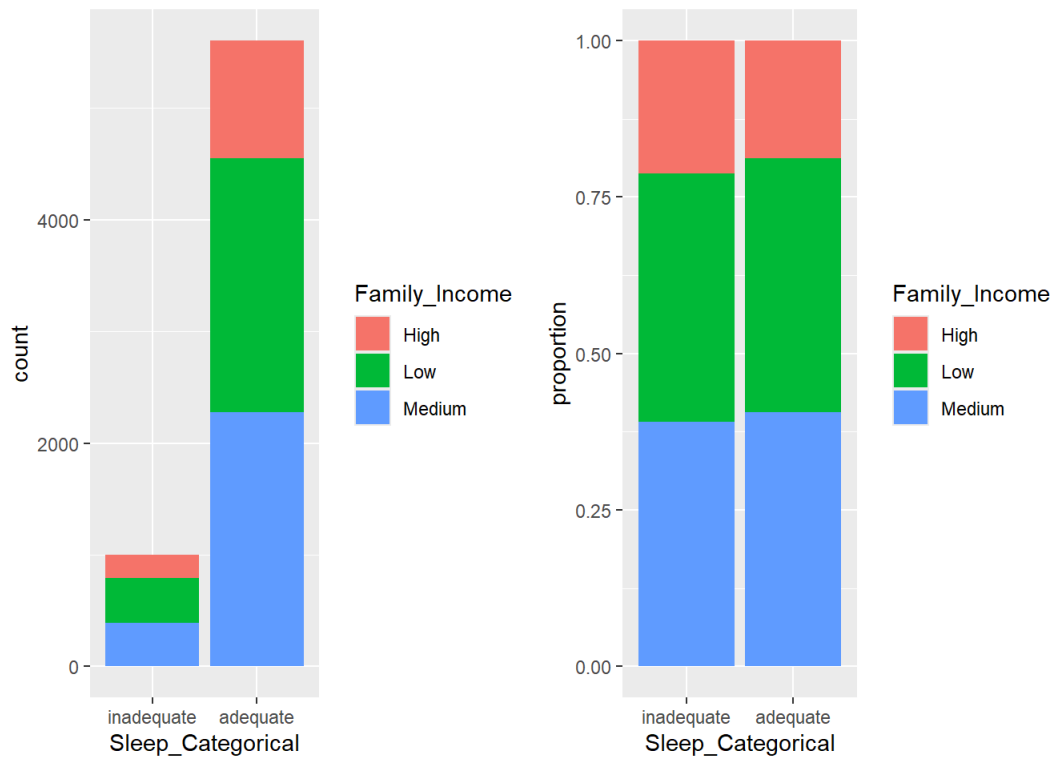
```

barchart_plot_3 <-
my_data_corrected %>%
  select(Sleep_Categorical, Family_Income) %>%
  drop_na() %>%
ggplot() +
  aes(x = Sleep_Categorical, fill = Family_Income) +
  geom_bar()

barchart_plot_4 <-
my_data_corrected %>%
  select(Sleep_Categorical, Family_Income) %>%
  drop_na() %>%
ggplot() +
  aes(x = Sleep_Categorical, fill = Family_Income) +
  geom_bar(position = "fill") +
  ylab("proportion")

grid.arrange(barchart_plot_3, barchart_plot_4, ncol=2)

```

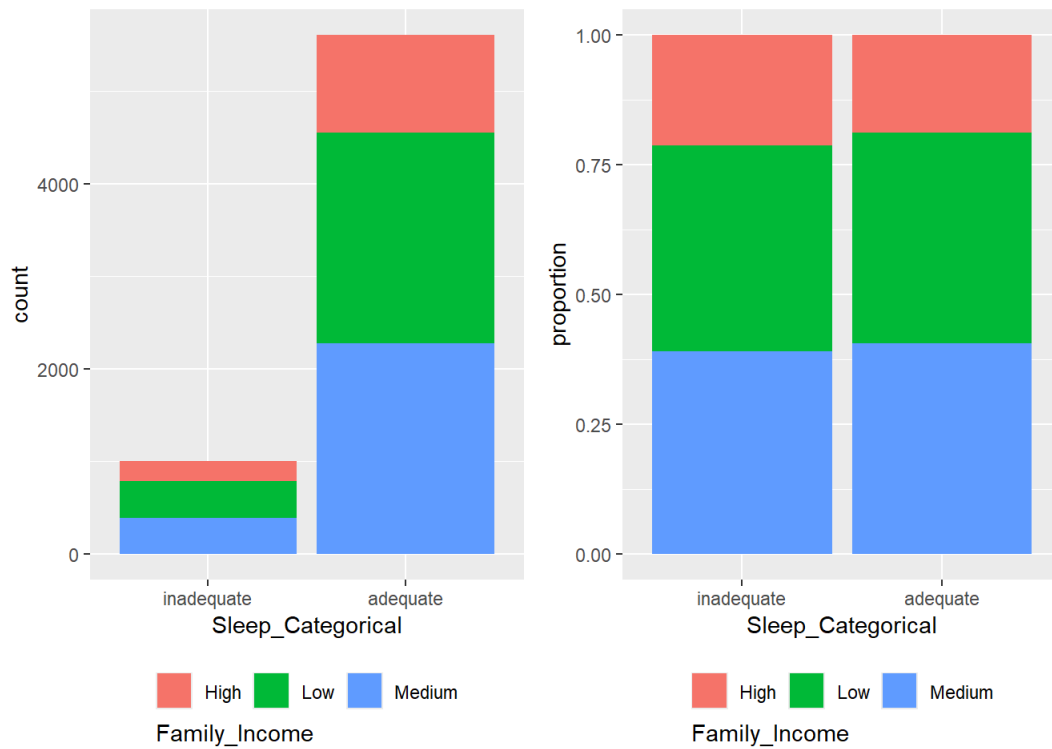


```
# a theme(legend.position) es a guides() funciók
# használatával kontrollálhatjuk hogy hol es hogyan
# jelenjen meg a jelmagyarazat az abran

barchart_plot_3 <-
my_data_corrected %>%
  select(Sleep_Categorical, Family_Income) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Sleep_Categorical, fill = Family_Income) +
    geom_bar() +
    theme(legend.position="bottom") +
    guides(fill = guide_legend(title.position = "bottom"))

barchart_plot_4 <-
my_data_corrected %>%
  select(Sleep_Categorical, Family_Income) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Sleep_Categorical, fill = Family_Income) +
    geom_bar(position = "fill") +
    theme(legend.position="bottom") +
    guides(fill = guide_legend(title.position = "bottom")) +
    ylab("proportion")

grid.arrange(barchart_plot_3, barchart_plot_4, ncol=2)
```

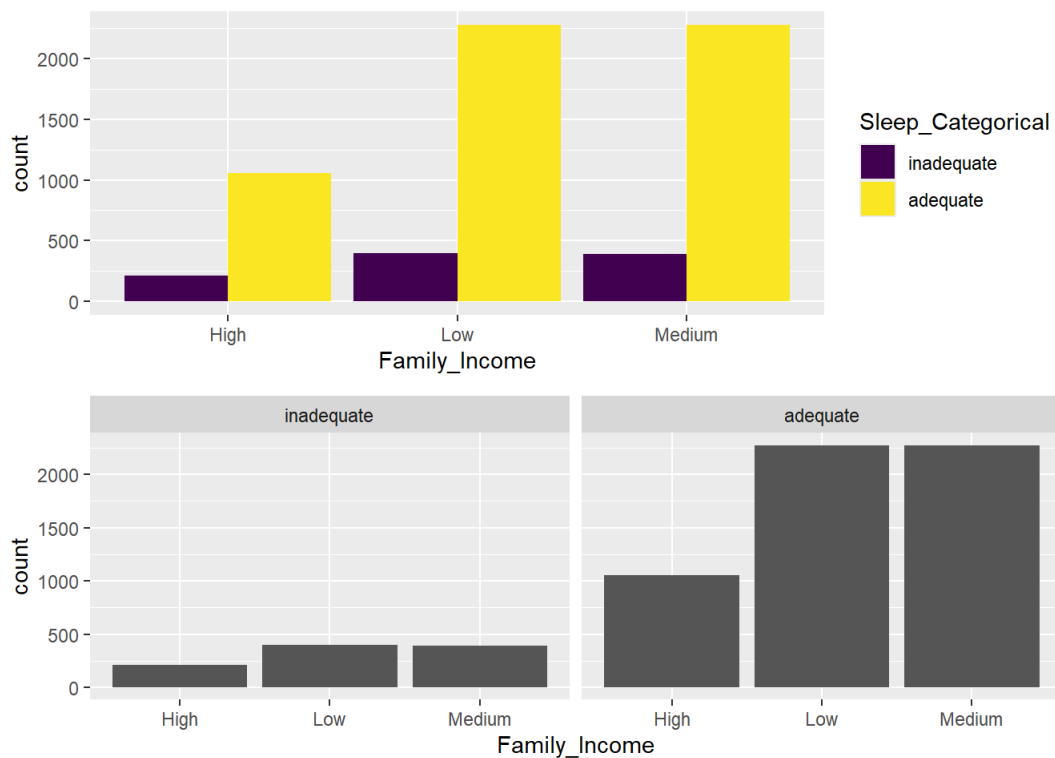


Ujabb módja a barchart segítségével való megjelenítésnek ha az oszlopok nem egymásra tornyozva, hanem **egymás mellett** jelennek meg, vagy ha a második változó szerint **külön paneleken (facet)** jelennek meg.

```
barchart_plot_5 <-
my_data_corrected %>%
  select(Sleep_Categorical, Family_Income) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Family_Income, fill = Sleep_Categorical) +
    geom_bar(position = "dodge")

barchart_plot_6 <-
my_data_corrected %>%
  select(Sleep_Categorical, Family_Income) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Family_Income) +
    geom_bar() +
    facet_wrap(~ Sleep_Categorical)

grid.arrange(barchart_plot_5, barchart_plot_6, nrow=2)
```



Egy kategorikus es egy numerikus valtozo kapcsolata

Vizsgáljuk meg hogy hogyan alakul a vizsgatelijesitmeny (Exam_Score) attol fuggoen hogy milyen a szuloi bevonodas (Parental_Involvement). Az Exam_Score egy folytonos numerikus valtozo, mig a Parental_Involvement kategorikus valtozo.

Az exploraciot kezdhethjuk leiro statisztikak lekerdezesevel csoportonként. Peldaul ha arra vagyunk kivancsiak, milyen a GDP atlaga es szorasa kontinensenként, ezt megvizsgálhatjuk a **group_by()** es a **summarize()** segitsegevel.

```
my_data_corrected = my_data_corrected %>%
  mutate(Parental_Involvement = factor(Parental_Involvement, ordered = T, levels = c("Low", "Medium", "High")))

my_data_corrected %>%
  select(Parental_Involvement, Exam_Score) %>%
  drop_na() %>%
  group_by(Parental_Involvement) %>%
  summarize(mean = mean(Exam_Score),
            sd = sd(Exam_Score))
```

```
## # A tibble: 3 × 3
##   Parental_Involvement mean    sd
##   <ord>                <dbl> <dbl>
## 1 Low                 66.3   3.86
## 2 Medium              67.1   3.73
## 3 High                68.1   3.95
```

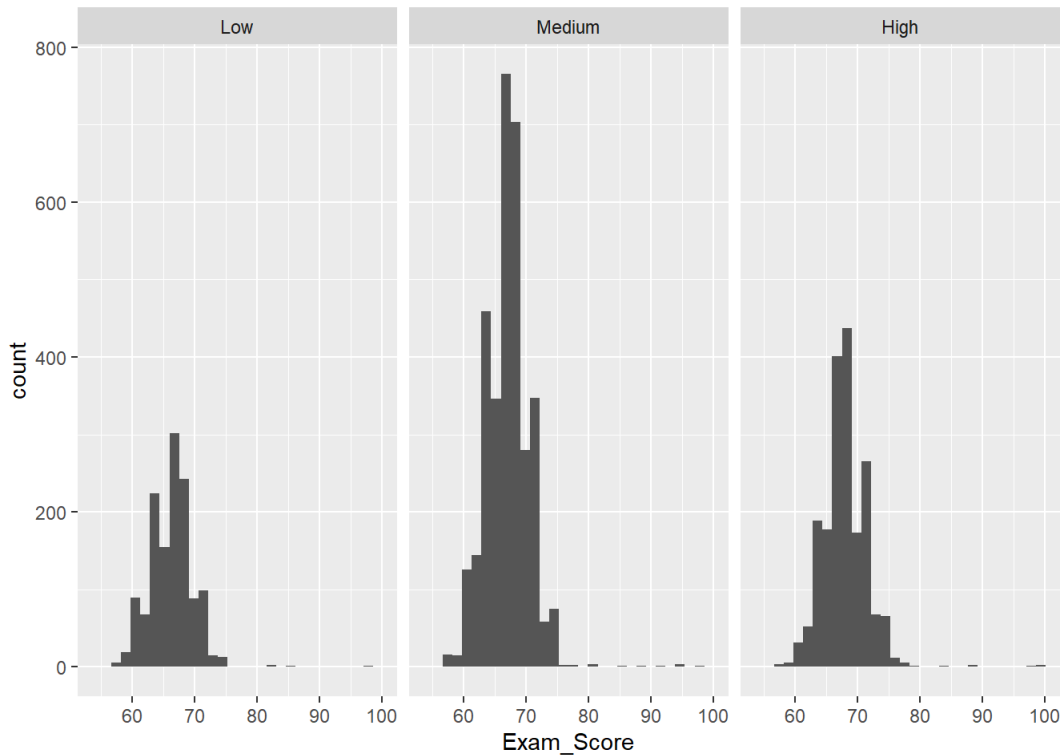
A ket valtozo kapcsolatat megvizsgalhatjuk **abrakkal** is. PI. hasznalhatjuk a

- **facet_wrap()** fuggvenyt egy **geom_histogram()**-al kobinalva
- a **geom_boxplot()** -ot
- esetleg hasznalhatunk egy egymásra illesztett **geom_density()** plot-ot ahol a kategoriak mas mas szinnel vannak jelolve.
- talan ebben az esetben a legtisztabb kepet a **geom_violin()** mutatja, ami a **geom_boxplot()** es a **geom_density()** keverekenek tekintheto. Ezt kiegészithetunk egy **geom_point()** -al, hogy pontosan latsszon, hany megfigyelesen alapulnak az abra adatai.
- az egyik kedvencem a **geom_violin()** a **geom_jitter()**-el valo kombinacioban

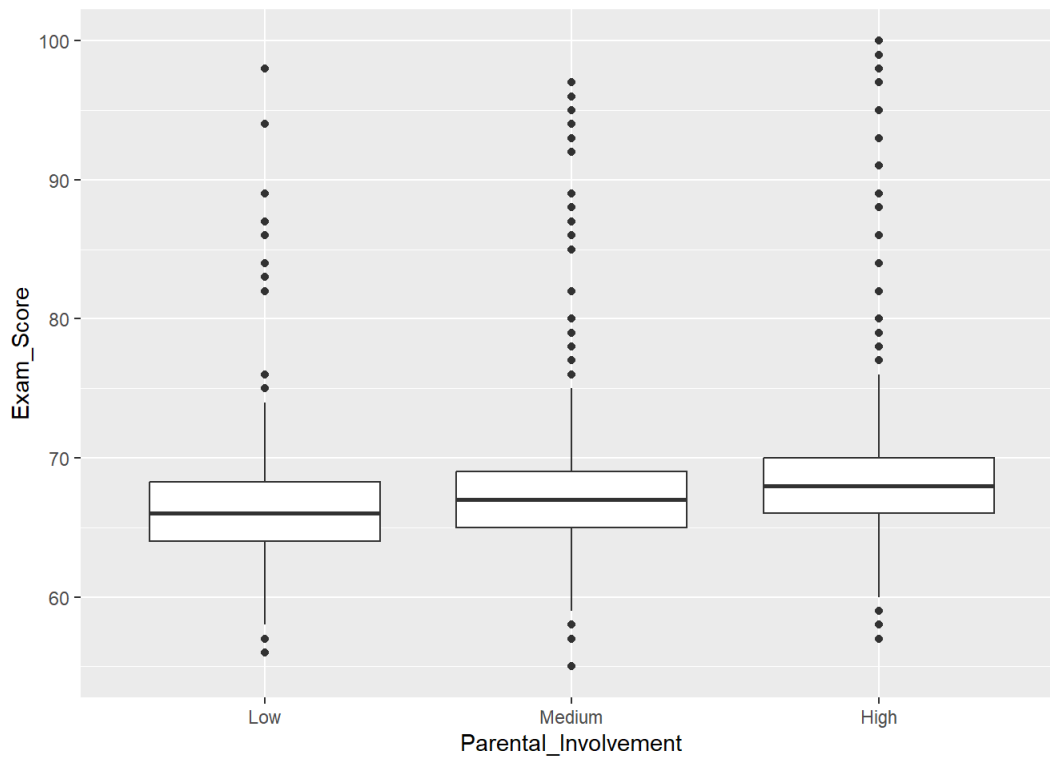
Mindig erdemes **tobb megkozelitest** is hasznalni az adat-exploracio kozben, hogy minel reszletesebb kepet kaphassunk, es csokkentsuk a valoszinuseget hogy egyik vagy masik megkozelites hanyossagai felrevezetnek minket.


```
my_data_corrected %>%
  select(Parental_Involvement, Exam_Score) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Exam_Score) +
    geom_histogram() +
    facet_wrap(~ Parental_Involvement)
```

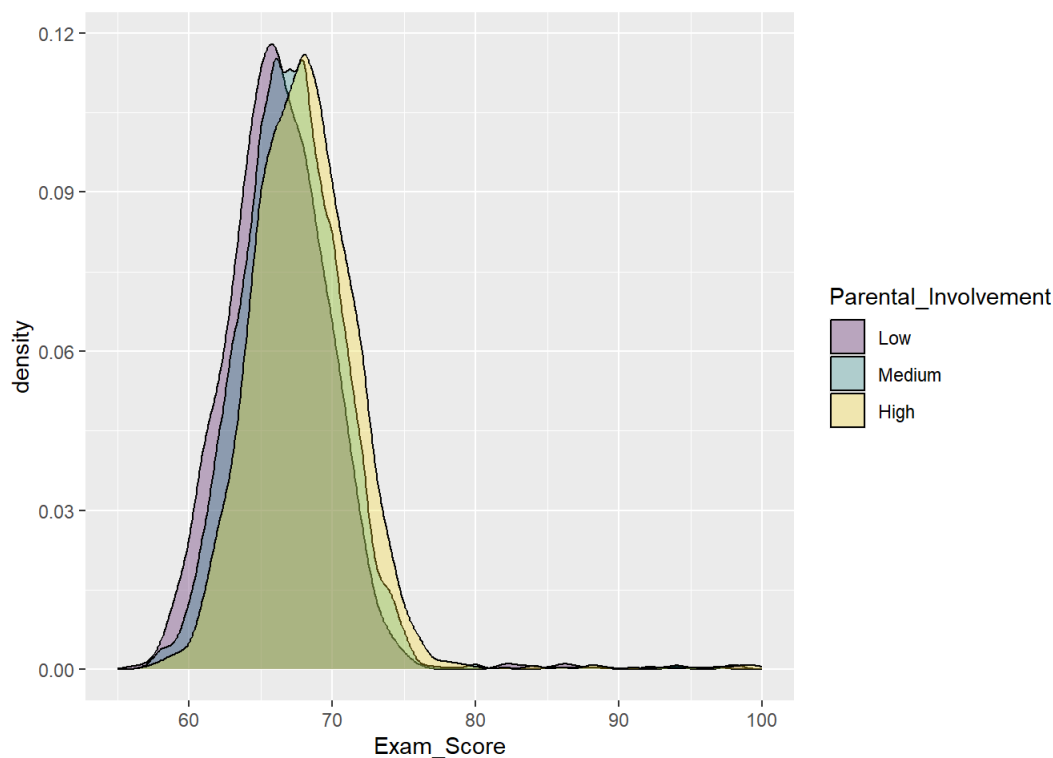
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



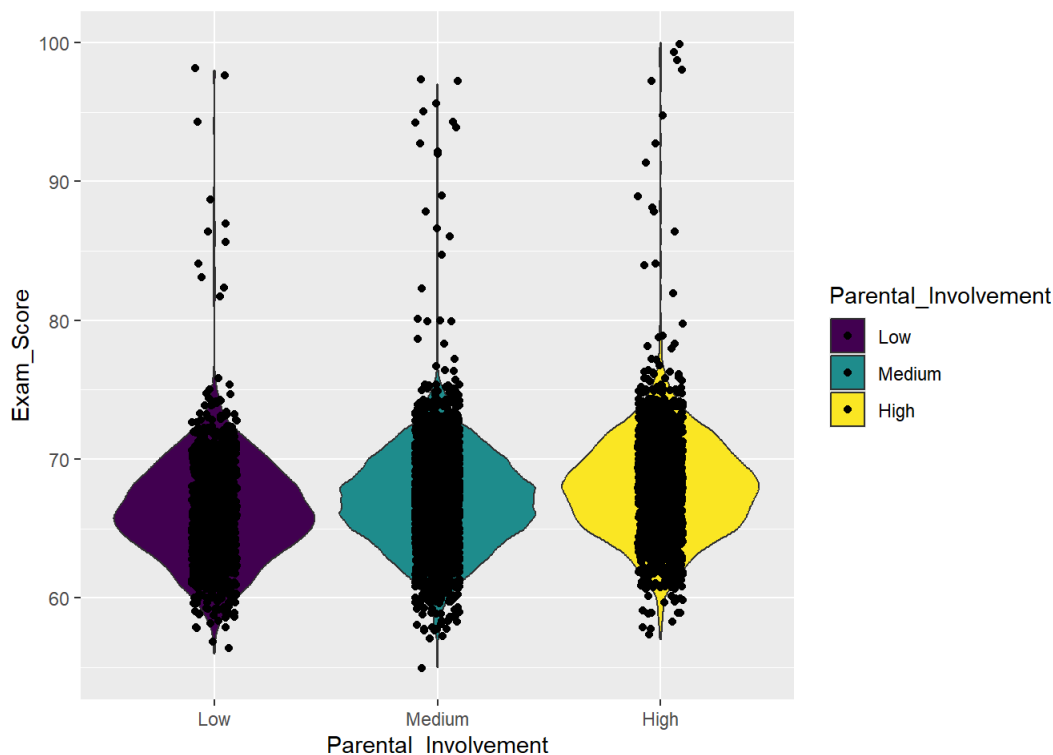
```
my_data_corrected %>%
  select(Parental_Involvement, Exam_Score) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Parental_Involvement, y = Exam_Score) +
    geom_boxplot()
```



```
my_data_corrected %>%
  select(Parental_Involvement, Exam_Score) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Exam_Score, fill = Parental_Involvement) +
    geom_density(alpha = 0.3)
```



```
my_data_corrected %>%
  select(Parental_Involvement, Exam_Score) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Parental_Involvement, y = Exam_Score, fill = Parental_Involvement) +
    geom_violin() +
    geom_jitter(width = 0.1)
```



Gyakorlas

Hasznald a fent tanult módszereket, hogy megvizsgald az **Exam_Score** es a **Learning_Disabilities** változók közötti összefüggést.

- hasznald a fenti geomokat, es készíts legalább két különböző ábrát más-más geomokkal

Két numerikus változó kapcsolata

Két numerikus változó közötti kapcsolat jellemzésére általában a korrelációs együtthatót szoktuk használni (`cor()`). A **cor()** funkciót akár több mint két változó páronkénti korrelációjának meghatározására is lehet használni.

A **drop_na()** funkcióval kiejthetjük azokat a megfigyeléseket, ahol a változók bármelyikében hiányzó adat (NA) van. Ha ezt nem tesszük meg, a `cor()` függvény NA eredményt adhatna ha valamelyik változóban NA-val találkozunk.

```
my_data_corrected %>%
  select(Exam_Score, Hours_Studied) %>%
  drop_na() %>%
  cor()
```

```
##           Exam_Score Hours_Studied
## Exam_Score    1.0000000    0.4465136
## Hours_Studied 0.4465136    1.0000000
```

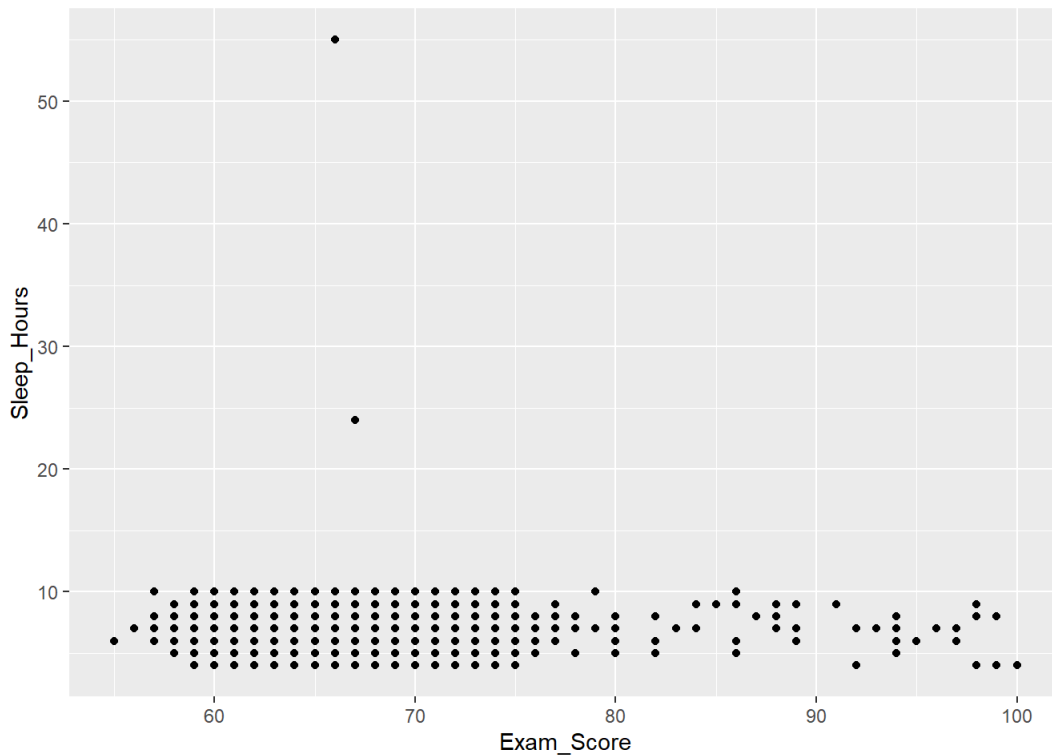
```
my_data_corrected %>%
  select(Exam_Score, Hours_Studied, Previous_Scores) %>%
  drop_na() %>%
  cor()
```

```
##           Exam_Score Hours_Studied Previous_Scores
## Exam_Score    1.0000000    0.44651362    0.17446099
## Hours_Studied 0.4465136    1.00000000    0.02463018
## Previous_Scores 0.1744610    0.02463018    1.00000000
```

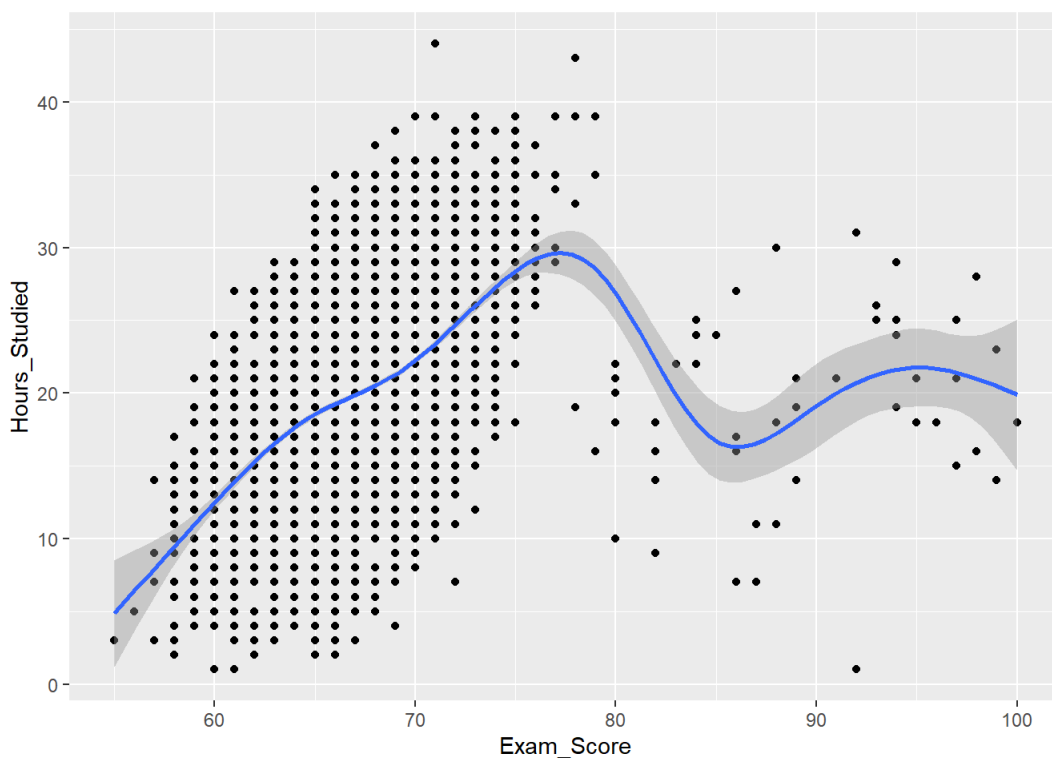
A numerikus változók közötti kapcsolatot általában pont diagrammal szoktuk ábrázolni (**geom_point()**)

A **geom_smooth()** layer hozzáadásával kaphatunk a pontok között meghúzott trendről egy képet. A fekete vonal az úgynevezett trendvonal, a szürke sáv a konfidencia intervallum. Ezekről később még részletesebben beszélünk majd

```
my_data_corrected %>%
  select(Exam_Score, Sleep_Hours) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Exam_Score, y = Sleep_Hours) +
    geom_point()
```



```
my_data_corrected %>%
  select(Exam_Score, Hours_Studied) %>%
  drop_na() %>%
  ggplot() +
    aes(x = Exam_Score, y = Hours_Studied) +
    geom_point() +
    geom_smooth()
```



Gyakorlas

Milyen eros a kapcsolat az Exam_Score es a Sleep_Hours kozott?

- határozd meg a korrelacios egyutthatot a valtozok kozott
- abrazold a valtozok kapcsolatát

Tobb folytonos valtozo kapcsolata megjelenitheto peldaul ugy, hogy az egyik valtozot egy szinskalahoz rendeljuk az alabbi modon.

```
my_data_corrected %>%  
  select(Exam_Score, Hours_Studied, Tutoring_Sessions) %>%  
  drop_na() %>%  
  ggplot() +  
    aes(x = Exam_Score, y = Hours_Studied, col = Tutoring_Sessions) +  
    geom_point()+  
    scale_colour_gradientn(colours=c("green", "black"))
```

