

PSZB17-210 - Seminar_4

Zoltan Kekecs

Marcus 4, 2020

4. Ora - Adatexploracio

Az ora celja az adatexploracios modszerek elsajaitasa.

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 altal 2020.09.28-an feltoltott COVID-19 adatokat a `read_csv()` funkcioval, es elmentjuk egy `COVID_adat` nevű objektumba. A `read_csv()` funkcio a tidyverse resze, es egybol tibble formatumban menti el az adatainkat.

```
COVID_adat <- read_csv("https://raw.githubusercontent.com/owid/covid-19-data/master/public/data/owid-covid-19.csv")
```

Adatok attekintese

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

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

```
COVID_adat
```

```
## # A tibble: 46,902 x 41
##   iso_code continent location date      total_cases new_cases new_cases_smoothed
##   <chr>      <chr>      <chr>  <date>          <dbl>      <dbl>          <dbl>
## 1 AFG      Asia      Afghani~ 2019-12-31          0          0             NA
## 2 AFG      Asia      Afghani~ 2020-01-01          0          0             NA
## 3 AFG      Asia      Afghani~ 2020-01-02          0          0             NA
## 4 AFG      Asia      Afghani~ 2020-01-03          0          0             NA
## 5 AFG      Asia      Afghani~ 2020-01-04          0          0             NA
## 6 AFG      Asia      Afghani~ 2020-01-05          0          0             NA
## 7 AFG      Asia      Afghani~ 2020-01-06          0          0              0
## 8 AFG      Asia      Afghani~ 2020-01-07          0          0              0
## 9 AFG      Asia      Afghani~ 2020-01-08          0          0              0
```

```
## 10 AFG      Asia      Afghani~ 2020-01-09      0      0      0
## # ... with 46,892 more rows, and 34 more variables: total_deaths <dbl>,
## #   new_deaths <dbl>, new_deaths_smoothed <dbl>, total_cases_per_million <dbl>,
## #   new_cases_per_million <dbl>, new_cases_smoothed_per_million <dbl>,
## #   total_deaths_per_million <dbl>, new_deaths_per_million <dbl>,
## #   new_deaths_smoothed_per_million <dbl>, new_tests <lgl>, total_tests <lgl>,
## #   total_tests_per_thousand <lgl>, new_tests_per_thousand <lgl>,
## #   new_tests_smoothed <lgl>, new_tests_smoothed_per_thousand <lgl>,
## #   tests_per_case <lgl>, positive_rate <lgl>, tests_units <lgl>,
## #   stringency_index <dbl>, population <dbl>, population_density <dbl>,
## #   median_age <dbl>, aged_65_older <dbl>, aged_70_older <dbl>,
## #   gdp_per_capita <dbl>, extreme_poverty <dbl>, cardiovasc_death_rate <dbl>,
## #   diabetes_prevalence <dbl>, female_smokers <dbl>, male_smokers <dbl>,
## #   handwashing_facilities <dbl>, hospital_beds_per_thousand <dbl>,
## #   life_expectancy <dbl>, human_development_index <dbl>
```

Leiro statisztikak

Ha az egyes változók **leiro statisztikaira** (descriptive statistics) vagyunk kíváncsiak, kerhetjük ezt a már tanult módon.

Peldaul lekerhetjük a változó alapvető legalacsonyabb és legmagasabb értéket, átlagát, medianját, a kvartiliseket, és hogy hány hiányzó adat van (ha van) a **summary()** funkcióval (miután a select funkcióval kiválasztottuk, melyik változóra vagyunk kíváncsiak)

```
COVID_adat %>%
  select(total_cases) %>%
  summary()
```

```
##   total_cases
##   Min.      :    0
##   1st Qu.:   61
##   Median :  1046
##   Mean   : 103338
##   3rd Qu.: 11200
##   Max.    :33423469
##   NA's    :614
```

Vagy megkaphatjuk ugyanezt az összes változóra, ha ugyanezt az egész adattablára futtatjuk le. Persze a karakter osztályba tartozó változókna mindezeknek a leiro statisztikáknak nincs értelme, ott csak a class információt kaptuk az output-ban.

```
COVID_adat %>%
  summary()
```

Gyakorlas

- Hány regisztrált eset volt összesen Magyarországon a tegnapi napig (*total_cases*)?
 - Mi volt a legmagasabb új eset-szám Magyarországon (*new_cases*)?
-

Meg több leiro statisztika

A **Psych** package segítségével a **describe()** funkció meg több hasznos információt adhat. Ez a funkció elsősorban szám-változók leírására szolgál, és karakter típusú kategorikus változók esetén sok warning

message-et ad, ezért érdemes a funkciót csak a szám-változókra lefuttatni (ezt alább a select() funkcióval érem el.)

```
COVID_adat %>%
  select(-date, -iso_code, -continent, -location, -contains("tests"), -positive_rate) %>%
  describe()
```

##	vars	n	mean	sd	median
## total_cases	1	46288	103338.19	1069717.38	1046.00
## new_cases	2	46078	1450.73	13164.72	10.00
## new_cases_smoothed	3	45296	1438.72	12977.70	15.00
## total_deaths	4	46288	4131.58	38418.41	20.00
## new_deaths	5	46078	43.52	366.48	0.00
## new_deaths_smoothed	6	45296	43.65	357.28	0.14
## total_cases_per_million	7	46014	1934.39	4085.26	284.71
## new_cases_per_million	8	46014	24.86	75.18	1.58
## new_cases_smoothed_per_million	9	45231	24.49	57.96	2.98
## total_deaths_per_million	10	46014	58.27	144.18	5.03
## new_deaths_per_million	11	46014	0.57	3.00	0.00
## new_deaths_smoothed_per_million	12	45231	0.57	1.90	0.02
## stringency_index	13	39192	57.37	27.41	62.96
## population	14	46628	88435794.18	612946888.59	8654618.00
## population_density	15	44499	360.83	1656.04	88.12
## median_age	16	41819	31.32	9.03	31.40
## aged_65_older	17	41197	9.26	6.32	6.98
## aged_70_older	18	41602	5.86	4.32	4.42
## gdp_per_capita	19	41279	20905.50	20433.43	14103.45
## extreme_poverty	20	27542	12.11	19.22	1.80
## cardiovasc_death_rate	21	41827	251.57	117.54	238.34
## diabetes_prevalence	22	43303	8.05	4.15	7.11
## female_smokers	23	32768	10.81	10.48	6.40
## male_smokers	24	32353	32.64	13.42	31.40
## handwashing_facilities	25	19582	52.46	31.60	55.18
## hospital_beds_per_thousand	26	37786	3.11	2.53	2.50
## life_expectancy	27	46040	74.03	7.37	75.49
## human_development_index	28	40355	0.72	0.15	0.75
##	trimmed	mad	min	max	
## total_cases	7697.92	1547.83	0.00	3.342347e+07	
## new_cases	99.65	14.83	-8261.00	3.209380e+05	
## new_cases_smoothed	103.42	22.24	-552.00	2.968079e+05	
## total_deaths	179.19	29.65	0.00	1.002678e+06	
## new_deaths	1.82	0.00	-1918.00	1.049100e+04	
## new_deaths_smoothed	1.99	0.21	-232.14	7.456710e+03	
## total_cases_per_million	967.10	421.88	0.00	4.349475e+04	
## new_cases_per_million	9.44	2.35	-2212.55	4.944380e+03	
## new_cases_smoothed_per_million	10.78	4.42	-269.98	8.829200e+02	
## total_deaths_per_million	21.08	7.46	0.00	1.237550e+03	
## new_deaths_per_million	0.13	0.00	-67.90	2.153800e+02	
## new_deaths_smoothed_per_million	0.17	0.03	-9.68	6.314000e+01	
## stringency_index	59.55	27.46	0.00	1.000000e+02	
## population	15702555.21	12405280.40	809.00	7.794799e+09	
## population_density	124.61	94.65	0.14	1.934750e+04	
## median_age	31.35	12.16	15.10	4.820000e+01	
## aged_65_older	8.69	5.93	1.14	2.705000e+01	
## aged_70_older	5.38	3.91	0.53	1.849000e+01	

## gdp_per_capita	17707.82	15808.23	661.24	1.169356e+05
## extreme_poverty	7.67	2.37	0.10	7.760000e+01
## cardiovasc_death_rate	240.43	121.87	79.37	7.244200e+02
## diabetes_prevalence	7.63	3.68	0.99	2.336000e+01
## female_smokers	9.49	8.01	0.10	4.400000e+01
## male_smokers	31.98	14.38	7.70	7.810000e+01
## handwashing_facilities	53.01	45.28	1.19	9.900000e+01
## hospital_beds_per_thousand	2.73	1.93	0.10	1.380000e+01
## life_expectancy	74.72	6.98	53.28	8.675000e+01
## human_development_index	0.73	0.16	0.35	9.500000e-01
##	range	skew	kurtosis	se
## total_cases	3.342347e+07	20.81	501.25	4972.04
## new_cases	3.291990e+05	16.64	315.47	61.33
## new_cases_smoothed	2.973599e+05	16.51	308.42	60.98
## total_deaths	1.002678e+06	17.96	366.58	178.57
## new_deaths	1.240900e+04	14.66	245.57	1.71
## new_deaths_smoothed	7.688860e+03	13.91	211.66	1.68
## total_cases_per_million	4.349475e+04	4.31	25.51	19.04
## new_cases_per_million	7.156920e+03	12.49	507.59	0.35
## new_cases_smoothed_per_million	1.152900e+03	5.23	41.71	0.27
## total_deaths_per_million	1.237550e+03	4.25	22.41	0.67
## new_deaths_per_million	2.832800e+02	30.75	1637.65	0.01
## new_deaths_smoothed_per_million	7.282000e+01	9.70	153.80	0.01
## stringency_index	1.000000e+02	-0.60	-0.66	0.14
## population	7.794798e+09	11.80	144.08	2838568.97
## population_density	1.934736e+04	9.93	106.26	7.85
## median_age	3.310000e+01	-0.03	-1.22	0.04
## aged_65_older	2.591000e+01	0.65	-0.87	0.03
## aged_70_older	1.797000e+01	0.79	-0.55	0.02
## gdp_per_capita	1.162744e+05	1.65	3.46	100.57
## extreme_poverty	7.750000e+01	1.81	2.32	0.12
## cardiovasc_death_rate	6.450500e+02	0.91	0.86	0.57
## diabetes_prevalence	2.237000e+01	1.09	1.42	0.02
## female_smokers	4.390000e+01	0.89	-0.31	0.06
## male_smokers	7.040000e+01	0.55	0.34	0.07
## handwashing_facilities	9.781000e+01	-0.13	-1.45	0.23
## hospital_beds_per_thousand	1.370000e+01	1.77	3.95	0.01
## life_expectancy	3.347000e+01	-0.75	-0.11	0.03
## human_development_index	6.000000e-01	-0.50	-0.74	0.00

Gyakorlas

- Mi az egy millio fore eso uj esetek (*new_cases_per_million*) ferdesegi mutatoja (skew/skewness)?
- Hany valid (nem NA) adat szerepel az adatbazisban az egy fore eso gdp-rol (*gdp_per_capita*)?

Faktorok

Nehany karaktervaltozonak csak **korlatozott mennyisegu eleme** lehet, mint peldaul a continent (North America, Asia, Africa, Europe, South America, Oceania). Ezeket megjelolhetjuk faktor (factor) osztalyu valtozokent, es akkor az R tobb informaciot fog adni rola.

A `levels()` funkcio megmutatja mik a faktorunk szintjei, de lathato ez akkor is ha csak meghivjuk a valtozot

magat.

A `table()` funkció pedig táblázatot készít arról, hogy az egyes csoportokban hány megfigyelés található

Amikor klistázzuk a faktor változót, akkor is kiírja az R a lista aljára, hogy milyen faktorszintek vannak. (Alább csinálunk egy `COVID_adat_tegnap` változót, amivel csak a tegnapi adatokat nézzük, hogy kisebb legyen az adattábla amivel dolgozunk.)

```
COVID_adat <- COVID_adat %>%
  mutate(continent = factor(continent),
         location = factor(location))

levels(COVID_adat$continent)

## [1] "Africa"      "Asia"      "Europe"    "North America"
## [5] "Oceania"    "South America"

table(COVID_adat$continent)

##
##      Africa      Asia      Europe North America      Oceania
##      10942     11224     12320      7325      1751
## South America
##      2792

COVID_adat_tegnap = COVID_adat %>%
  filter(date == "2020-09-28")

COVID_adat_tegnap$continent

## [1] Asia      Europe    Africa    Europe    Africa
## [6] North America North America South America Asia      North America
## [11] Oceania    Europe    Asia      North America Asia
## [16] Asia      North America Europe    Europe    North America
## [21] Africa    North America Asia      South America North America
## [26] Europe    Africa    South America North America Asia
## [31] Europe    Africa    Africa    Asia      Africa
## [36] North America Africa    North America Africa    Africa
## [41] South America Asia      South America Africa    Africa
## [46] North America Africa    Europe    North America North America
## [51] Europe    Europe    Africa    Europe    Africa
## [56] North America North America South America Africa    North America
## [61] Africa    Africa    Europe    Africa    Europe
## [66] South America Oceania    Europe    Europe    Oceania
## [71] Africa    Africa    Asia      Europe    Africa
## [76] Europe    Europe    North America North America Oceania
## [81] North America Europe    Africa    Africa    South America
## [86] North America North America Europe    Europe    Asia
## [91] Asia      Asia      Asia      Europe    Europe
## [96] Asia      Europe    North America Asia      Europe
## [101] Asia      Asia      Africa    Europe    Asia
## [106] Asia      Asia      Europe    Asia      Africa
## [111] Africa    Africa    Europe    Europe    Europe
## [116] Europe    Africa    Africa    Asia      Asia
## [121] Africa    Europe    Africa    Africa    North America
## [126] Europe    Europe    Asia      Europe    North America
## [131] Africa    Africa    Asia      Africa    Asia
```

```
## [136] Europe      Oceania      Oceania      North America Africa
## [141] Africa      Oceania      Europe      Asia      Asia
## [146] Asia      North America Oceania      South America South America
## [151] Asia      Europe      Europe      North America Asia
## [156] Europe      Europe      Africa      North America North America
## [161] North America Europe      Africa      Asia      Africa
## [166] Europe      Africa      Africa      Asia      North America
## [171] Europe      Europe      Africa      Africa      Asia
## [176] Africa      Europe      Asia      Africa      South America
## [181] Africa      Europe      Asia      Asia      Asia
## [186] Africa      Asia      Asia      Africa      North America
## [191] Africa      Asia      North America Africa      Europe
## [196] Asia      Europe      North America North America South America
## [201] Asia      Europe      South America Asia      Africa
## [206] Asia      Africa      Africa      <NA>      <NA>
## Levels: Africa Asia Europe North America Oceania South America
```

Igy mar a fenti `summary()` funkcio is kiadja az egyes faktorszintekrol hogy hanyan tartoznak oda.

```
COVID_adat_tegnap %>%
  select(continent) %>%
  summary()
```

```
##      continent
## Africa      :55
## Asia        :46
## Europe      :50
## North America:36
## Oceania     : 8
## South America:13
## NA's        : 2
```

Van, hogy szeretnenk **kizarni** bizonyos **faktorszinteket** az elemzesbol. Pl. ha valamelyik faktor szintbol nagyon keves megfigyeles van, mondjuk Oceaniat, mondjuk mert ugy gondoljuk hogy az tulsagosan “elszigetelt” a világ többi reszetol, oket lehet hogy szeretnenk kizarni a kesobbi elemzesekbol hogy egyszerusitsuk az eredményeink értelmezését. Ezt a mar korabban tanult `filter()` funkcio segitsegevel konnyeden megtehetjuk, azonban arra figyelniunk kell, hogy az R megjegyzi a faktorszinteket, es azt azt kovetoen is a **valtozohoz rendelve tartja**, miutan mar az adott faktorszintbol nincs egy megfigyeles sem az adattablában.

```
COVID_adat_tegnap %>%
  filter(continent != "Oceania") %>%
  select(total_cases, continent) %>%
  summary()
```

```
##   total_cases      continent
## Min.      :    3 Africa      :55
## 1st Qu.: 1731 Asia         :46
## Median : 9664 Europe      :50
## Mean   : 165381 North America:36
## 3rd Qu.: 72210 Oceania     : 0
## Max.    :7115046 South America:13
```

Igy ezeket a szinteket ejthetjuk a `droplevels()` funkcioval.

```
COVID_adat_tegnap_noOceania = COVID_adat_tegnap %>%
  filter(continent != "Oceania") %>%
  mutate(continent = droplevels(continent))
```

```
COVID_adat_tegnap_noOceania %>%
  select(continent) %>%
  summary()
```

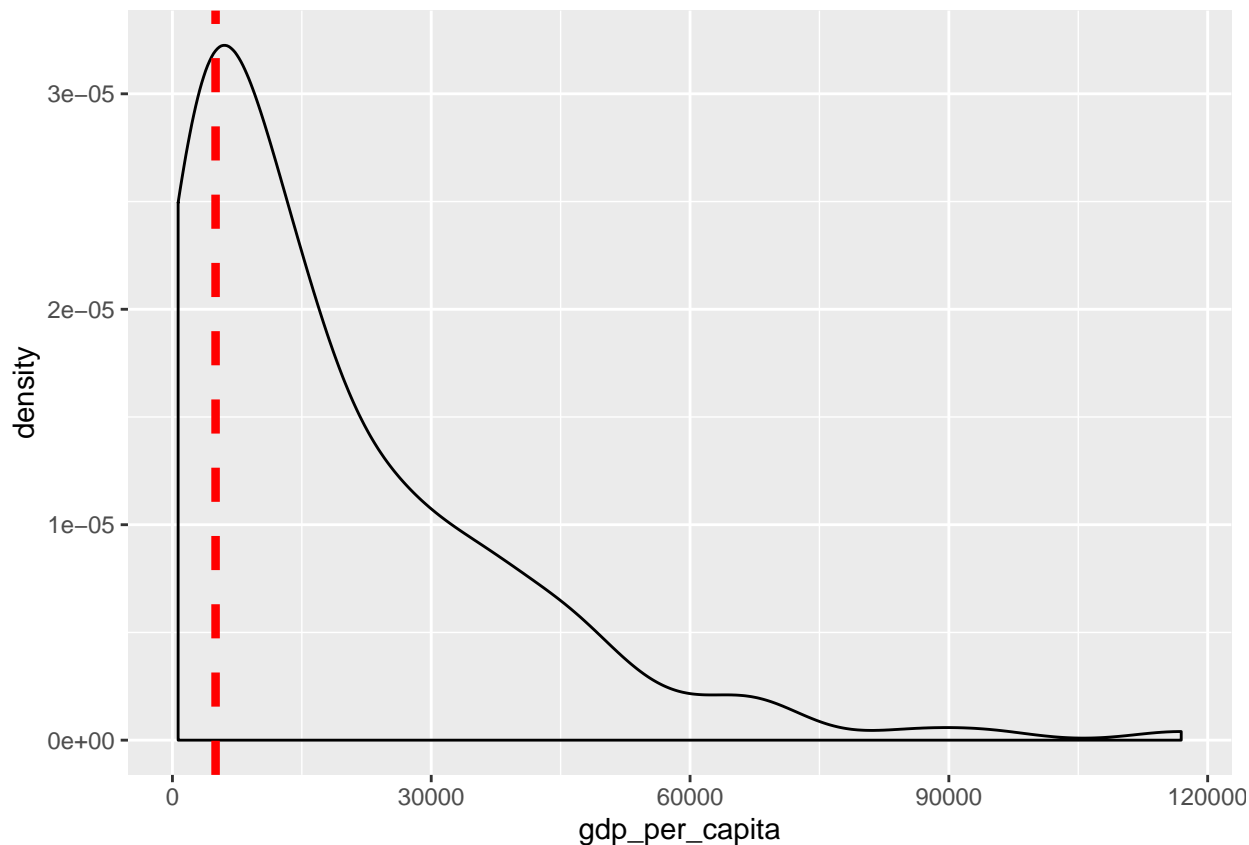
```
##           continent
## Africa           :55
## Asia             :46
## Europe           :50
## North America:36
## South America:13
```

Elofordul, hogy egy **numerikus változot akarunk atalakítani faktorra**, pl. elkepzelheto hogy ossze akarjuk hasonlítani azokat az orszagokat ahol 5000 alatti a gdp_per_capita azokkal akinek e feletti, hogy hogyan kulonboznak a COVID adatok.

```
COVID_adat_tegnap %>%
  select(gdp_per_capita, continent) %>%
  drop_na() %>%
  group_by(continent) %>%
  summarize(mean_gdp = mean(gdp_per_capita))
```

```
## # A tibble: 6 x 2
##   continent    mean_gdp
##   <fct>         <dbl>
## 1 Africa         5444.
## 2 Asia          22185.
## 3 Europe         33029.
## 4 North America  21655.
## 5 Oceania        23315.
## 6 South America  13841.
```

```
COVID_adat_tegnap %>%
  select(gdp_per_capita) %>%
  drop_na() %>%
  ggplot() +
  aes(x = gdp_per_capita) +
  geom_density() +
  geom_vline(xintercept = 5000, linetype="dashed",
            color = "red", size=1.5)
```



Folytonos változók atkódolása kategorikus változóra

Ilyenkor használhatjuk a `mutate()` és `case_when()` funkciók kombinációját hogy csináljunk egy új változót. Ebbe a kódba beleépítettem a `factor()` funkciót is, hogy azonnal meghatározzuk, hogy ez az új változó egy faktor, és nem egy egyszerű karaktervektor. A `factor()` funkció nélkül is lefut a kód, de akkor még kellene egy külön sor ahol megadjuk hogy ez egy faktorváltozó.

```
COVID_adat = COVID_adat %>%
  mutate(gdp_per_capita_kat = factor(
    case_when(gdp_per_capita < 5000 ~ "small",
              gdp_per_capita >= 5000 & gdp_per_capita < 10000 ~ "medium",
              gdp_per_capita > 10000 ~ "large")))
levels(COVID_adat$gdp_per_capita_kat)

## [1] "large" "medium" "small"

# ugyanez a COVID_adat_tegnap -al

COVID_adat_tegnap = COVID_adat_tegnap %>%
  mutate(gdp_per_capita_kat = factor(
    case_when(gdp_per_capita < 5000 ~ "small",
              gdp_per_capita >= 5000 & gdp_per_capita < 10000 ~ "medium",
              gdp_per_capita > 10000 ~ "large")))
```


Kategorikus változó újrakodolása

Hasonló eset ha kategorikus változókat szeretnénk atkódolni. Mondjuk ha szeretnénk a déli felteket az északi feltekeivel összehasonlítani. Ezt a **recode()** funkcióval lehet megoldani.

```
COVID_adat = COVID_adat %>%  
  mutate(continent_south_north = factor(recode(continent,  
                                                "Oceania" = "South",  
                                                "South America" = "South",  
                                                "Africa" = "South",  
                                                "Asia" = "North",  
                                                "Europe" = "North",  
                                                "North America" = "North"))))  
  
levels(COVID_adat$continent_south_north)
```

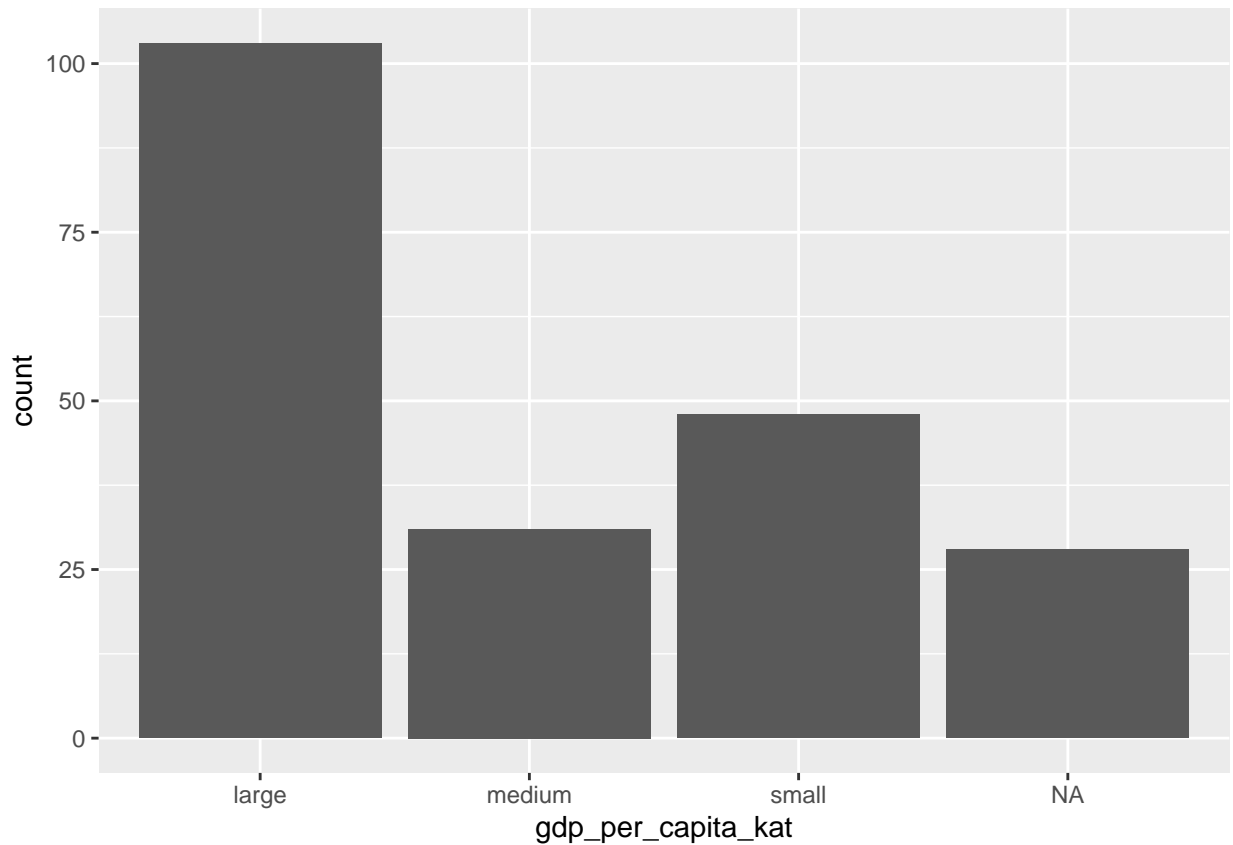
```
## [1] "South" "North"
```

```
COVID_adat_tegnap = COVID_adat_tegnap %>%  
  mutate(continent_south_north = factor(recode(continent,  
                                                "Oceania" = "South",  
                                                "South America" = "South",  
                                                "Africa" = "South",  
                                                "Asia" = "North",  
                                                "Europe" = "North",  
                                                "North America" = "North"))))
```

Faktorszintek sorrendje, ordinalis változók

Amikor van értelme a **sorrendiségnek** a faktorszintek között, **ordinalis változókrol** beszélünk (vagyis az egyik faktorszint alacsonyabb, vagy kisebb “értéku” mint a másik). Arra figyelünk kell, hogy amikor faktorokat hozunk létre, az R automatikusan a faktorszintek neveinek **ABC sorrendje** alapján rakja őket sorba, és az ábrák is így szemlélteti majd őket.

```
COVID_adat_tegnap %>%  
  ggplot() +  
  aes(x = gdp_per_capita_kat) +  
  geom_bar()
```

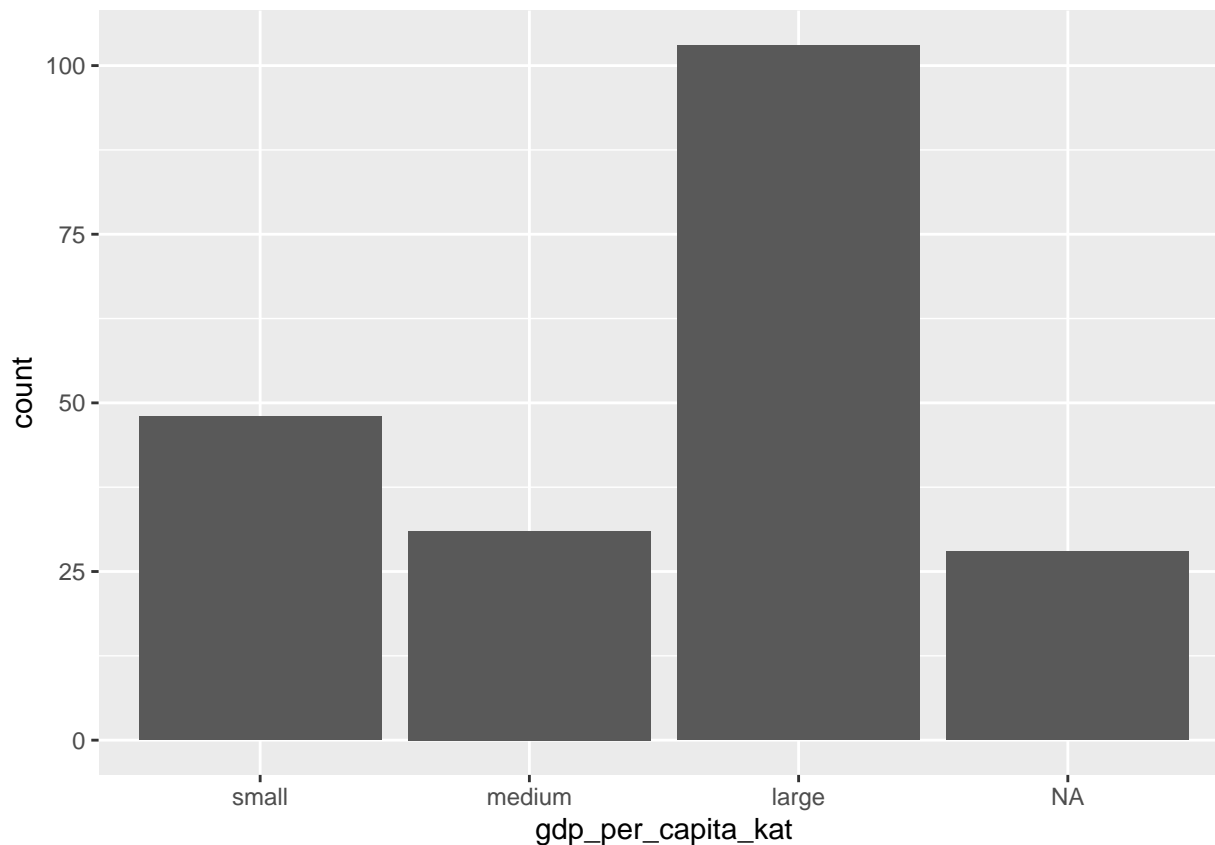


Ilyenkor érdemes meghatározni a faktorszintek sorrendjét (**order**). Ezt legegyszerűbben a `factor()` funkcion belül tehetjük meg, az **ordered** = **T** beállításával, és a **levels** = resznel a szintek sorrendjének meghatározásával.

```
COVID_adat_tegnap = COVID_adat_tegnap %>%
mutate(gdp_per_capita_kat = factor(gdp_per_capita_kat, ordered = T, levels = c(
  "small",
  "medium",
  "large"))))
```

Igy már az R minden funkciója tudni fogja, hogy egy ordinalis változóról van szó, ahol fontos a sorrend, és tudni fogja a sorrendet is.

```
COVID_adat_tegnap %>%
  ggplot() +
  aes(x = gdp_per_capita_kat) +
  geom_bar()
```



Gyakorlas

- szurd az adatokat ugy hogy csak a 2020-09-28-ai adatokkal dolgozzunk csak.
- csinalj egy uj kategorikus valtozot (nevezzuk ezt *new_cases_per_million_kat*-nak) a `mutate()` funkcio hasznalataval amiben azok az orszagok ahol a *new_cases_per_million* valtozo 20 alatt van “small”, ahol 20 vagy a felett van “large” kategoriaba keruljenek.
- figyelj oda hogy faktorkent jelold meg ezt az uj valtozot (Ezt lehet az elozi lepesben a `mutate()` funkcion belül, vagy egy kulon lepesben, de mindenképpen a `factor()` vagy az `as.factor()` funkciokat erdemes hozza hasznalni)
- mentsd el ezt a valtozot az eredeti adatobjektumban ugy hogy kesobb is lehessen vele dolgozni
- keszits egy tablazatot arrol, hogy hanyan esnek a *new_cases_per_million_kat* egyes kategoriaba.
- Add meg a faktorszintek helyes sorrendjet: small, large (Ird felul a *new_cases_per_million_kat* korabbi valtozatot ezzel a valtozattal ahol a szintek mar helyes sorrendben vannak, vagy ezt a sorrendezest is bele vonhatod az eredeti funkcioba, amivel a valtozot generaltad)
- Ellenorizd, hogy valoban helyes sorrendben szerepelnek-e a faktor szintjei.

Exploracio vizualizacion keresztul

Egyes valtozok vizualizacioja

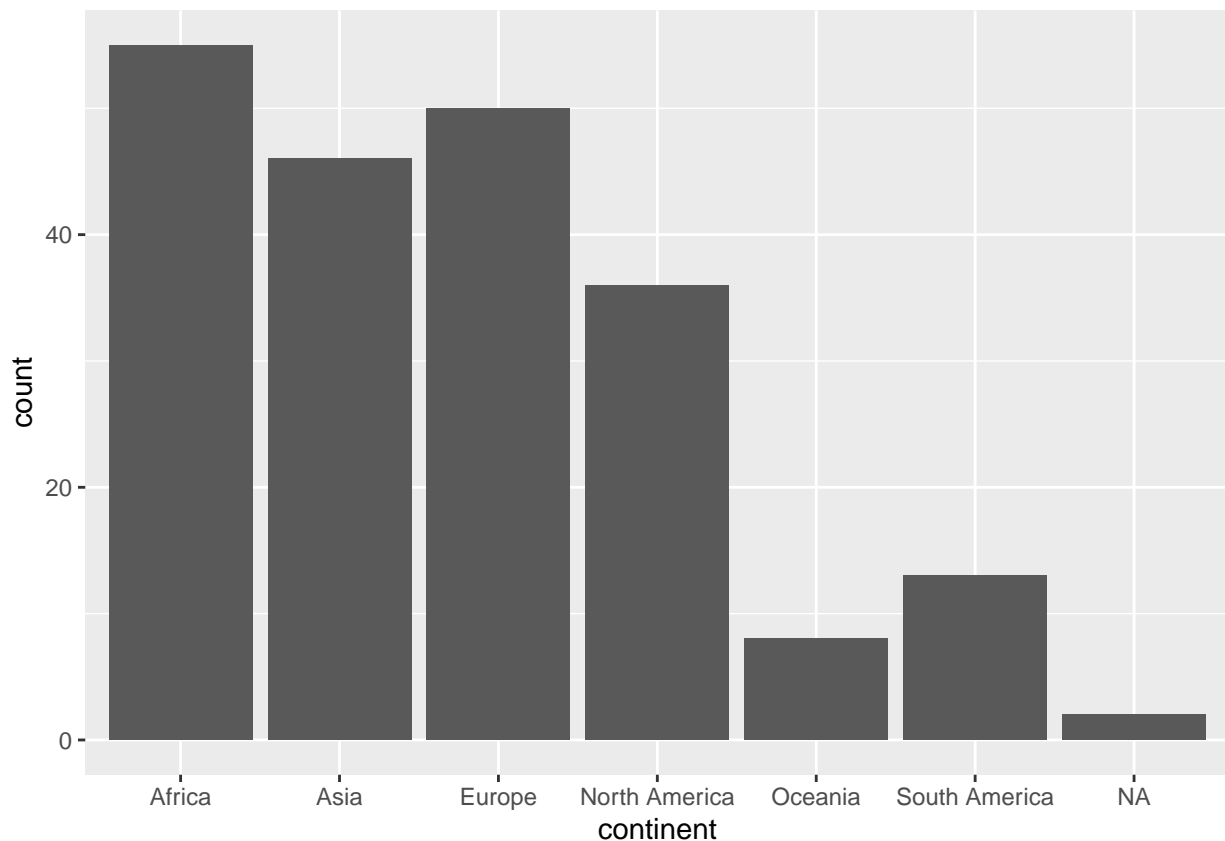
Az egyes valtozok **abrak** (plot) segitsegevel is megvizsgalhatok. A **kategorikus** valtozokat gyakran oszlopdiagrammal (**geom_bar**) abrazoljuk,

Mig a **numerikus** valtozokat inkabb **dotplot**, **histogram**, vagy **density plot** segitsegevel szoktuk abrazolni.

Az egyes valtozok vizualizacioja es a leiro statisztikak atvizsgalasa elengedhetetlen hogy azonositsuk az esetleges adatbeviteli **hibakat es egyeb nemvart furcsasagokat** az adataink kozott.

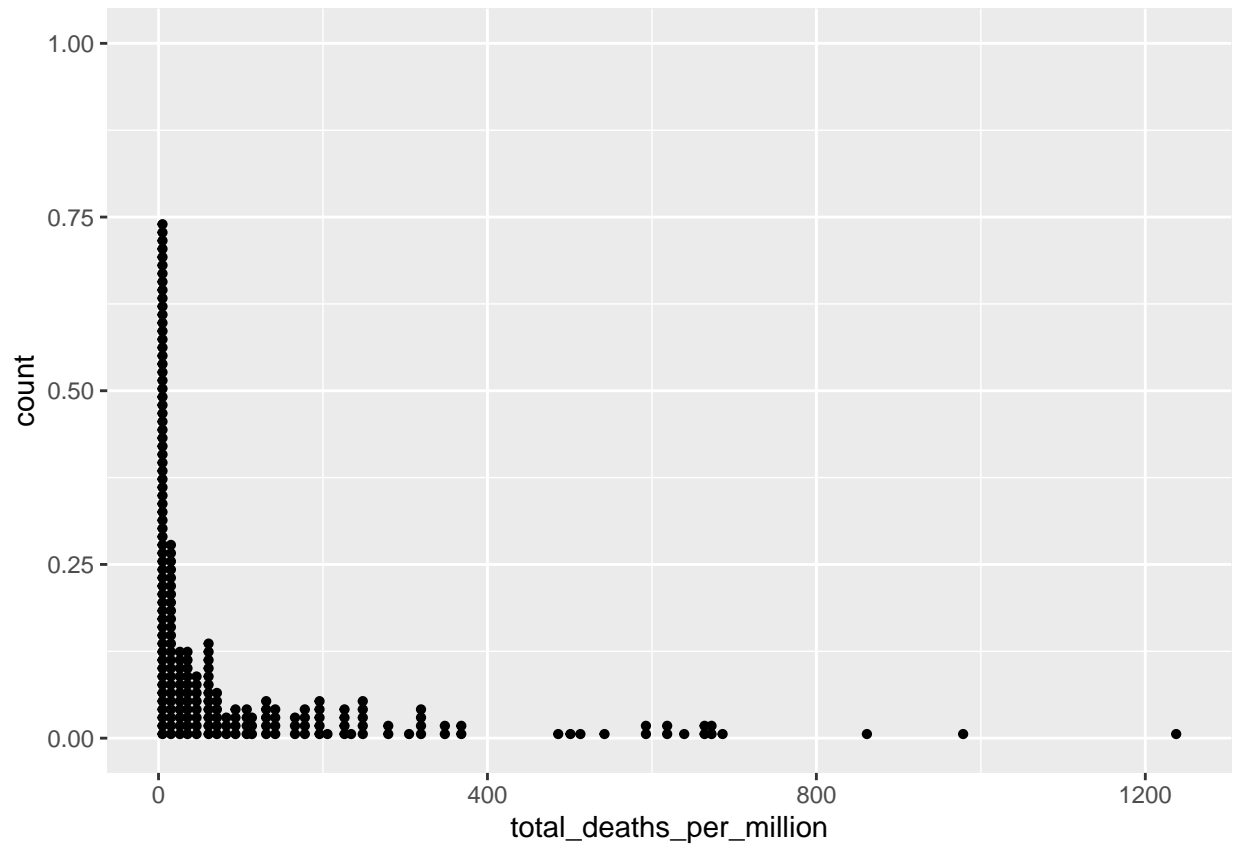
MINDING ellenorizd az adataidat ezekkel a modszerekkel mielőtt komolyabb adatelemzesbe kezdesz, hogy meggyozodj rola, hogy az adatok tisztak es megfelelnek az elvarasaidnak.

```
COVID_adat_tegnap %>%  
ggplot() +  
  aes(x = continent) +  
  geom_bar()
```



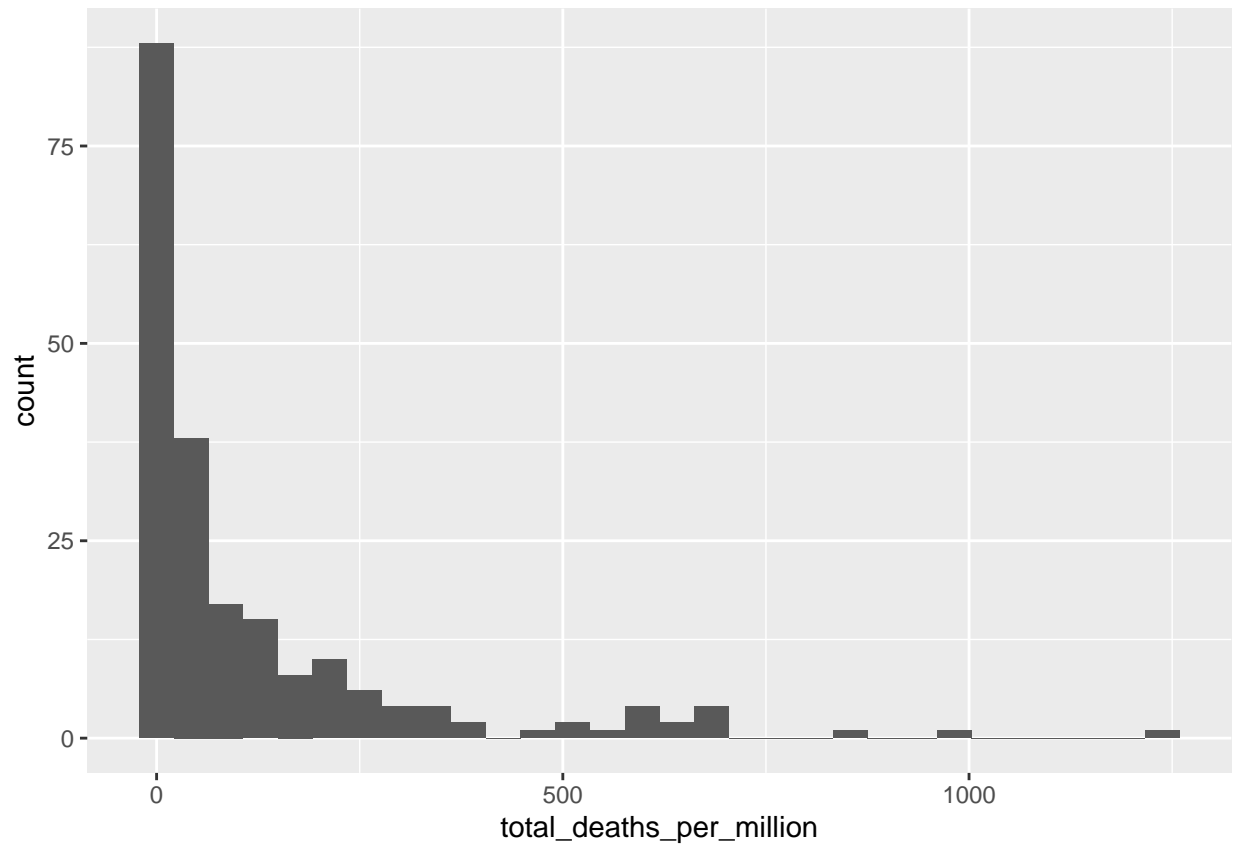
```
COVID_adat_tegnap %>%  
ggplot() +  
  aes(x = total_deaths_per_million) +  
  geom_dotplot(binwidth = 10)
```

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



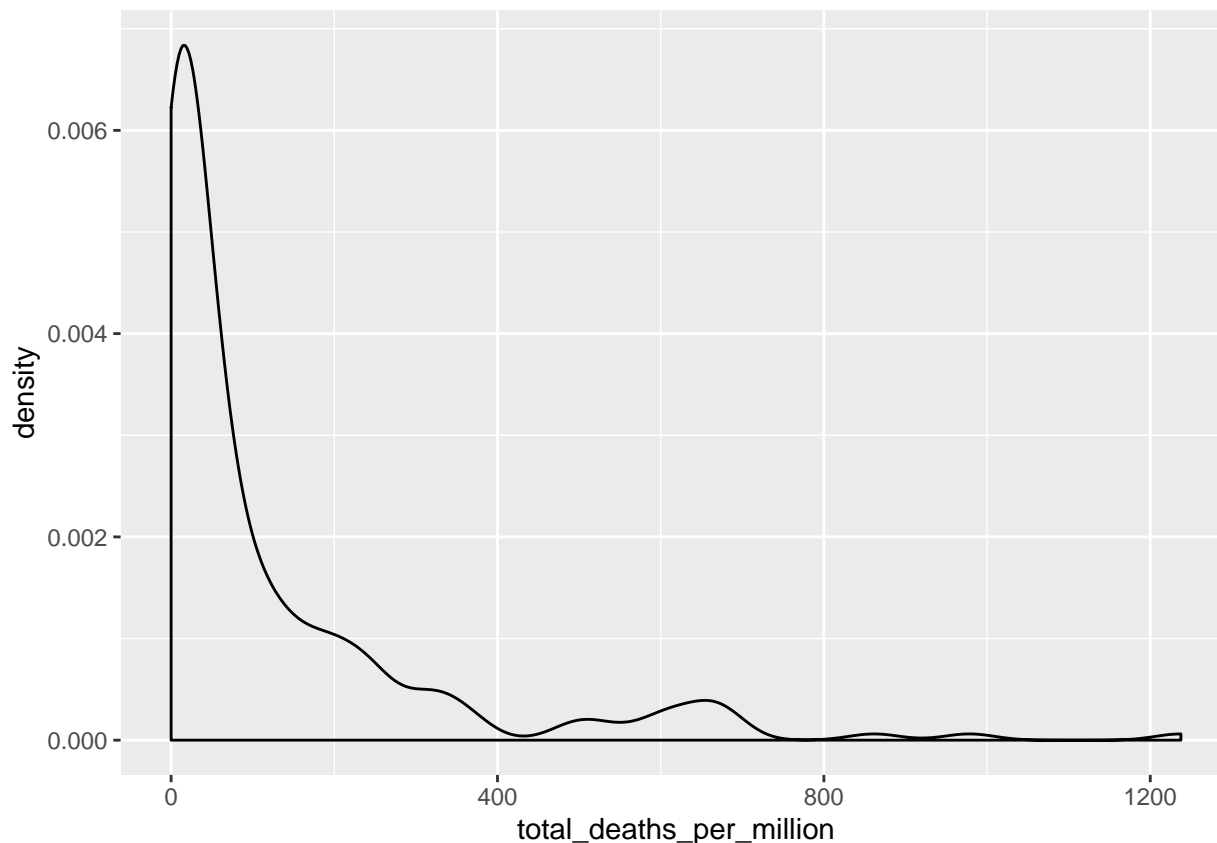
```
COVID_adat_tegnap %>%
  ggplot() +
    aes(x = total_deaths_per_million) +
    geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 1 rows containing non-finite values (stat_bin).
```



```
COVID_adat_tegnap %>%  
ggplot() +  
  aes(x = total_deaths_per_million) +  
  geom_density()
```

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



Gyakorlas

Szurd az adatokat úgy hogy csak a 2020-09-07-en jeletett adatokkal dolgozzunk

Hasznald a fent tanult módszereket, hogy **azonosítsd az COVID_adat adattáblában lévő hibákat** vagy nem várt furcsaságokat.

- A vizualizáción túl a View(), describe(), és summary() funciókat érdemes használni az adatok első betekintésére
- A numerikus (vagy éppen folytonos) változókna vizsgald meg a minimum és maximum értéket és a hiányzó adatok mennyiségét, valamint az eloszlást.
- A kategorikus változókna vizsgald meg az összes faktorszintet és az egyes szintekhez tartozó megfigyelések mennyiségét.

A hibákat a következőképpen javíthatjuk.

A **mutate()** és a **replace()** funkciók használatával **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 válasz, beírhatjuk, **DE minden javítást fel kell tüntetni** a kutatási jelentésben (és a ZH során is), hogy az olvasó számára tiszta legyen, hogy itt egy adathelyettesítés vagy kizárás történt!

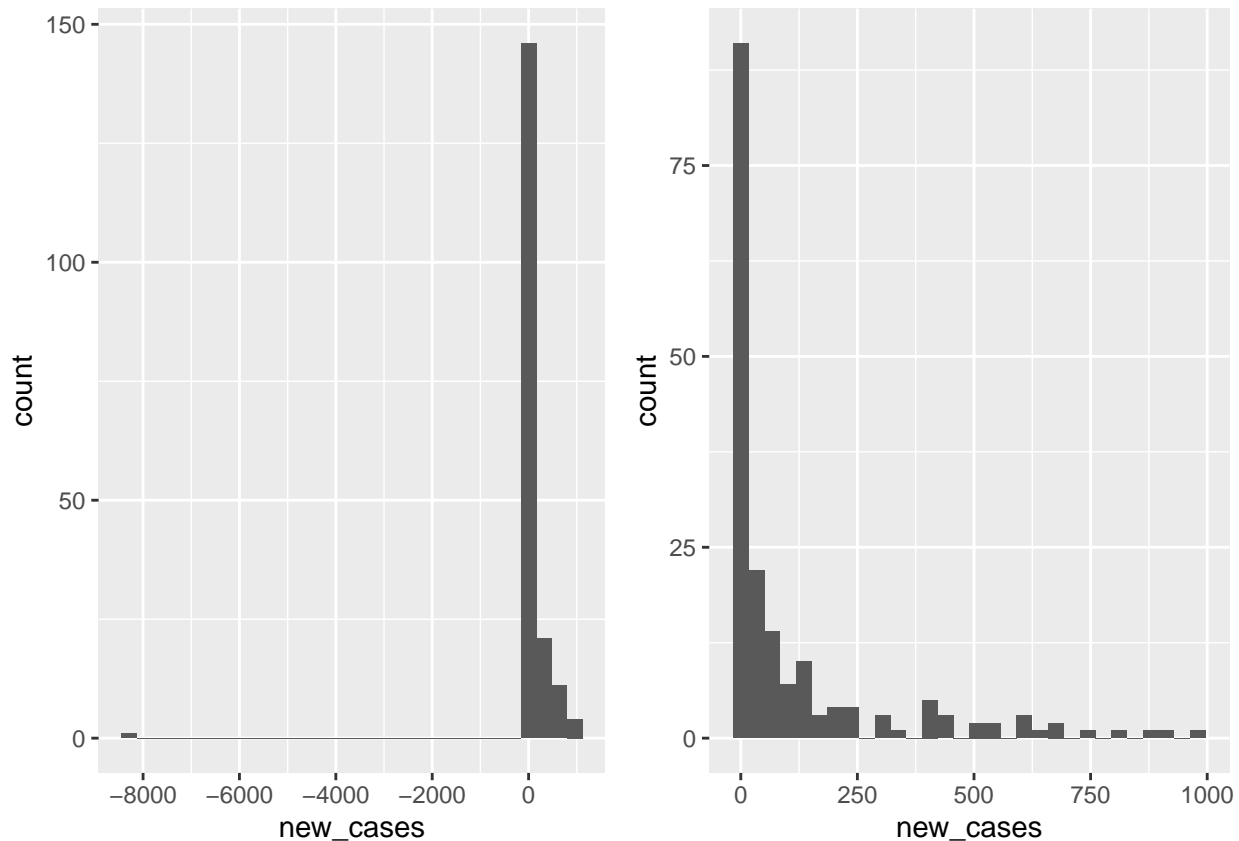
Mindig érdemes a javított adatokat **új adattáblába** elmenteni. A mi esetünkben az COVID_adat_corrected nevet adtuk a javított objektumnak. Így a nyers adataink megmaradnak, ami hasznos lehet későbbi

muveleteknel.

```
COVID_adat_corrected <- COVID_adat %>%  
  mutate(new_cases = replace(new_cases, new_cases=="-8261", NA))
```

Erdemes **megbizonyosodni** rola, hogy az adatcsere sikeres volt. Alabb az adatok vizualizaciojaval gyzodunk meg errol, de az adatok megjelenitesevel, vagy a leiro statisztikak lekerdezesevel is megtehető ez, ha az informatív.

```
# hasznalhatnak meg az alabbiakat is arra,  
# hogy megbizonyosodjunk abban, hogy sikeres volt a csere  
# View(COVID_adat_corrected)  
# describe(COVID_adat_corrected)  
# summary(COVID_adat_corrected$szocmedia_3)  
# COVID_adat_corrected$szocmedia_3  
  
old_plot <-  
  COVID_adat %>%  
  filter(date == "2020-09-07", new_cases < 1000) %>%  
  ggplot()+  
    aes(x = new_cases)+  
    geom_histogram()  
  
new_plot <-  
  COVID_adat_corrected %>%  
  filter(date == "2020-09-07", new_cases < 1000) %>%  
  ggplot()+  
    aes(x = new_cases)+  
    geom_histogram()  
  
grid.arrange(old_plot, new_plot, ncol=2)
```

Tobb változó kapcsolatának felterkepezése

Több változó kapcsolatát is felterkepezhetjük táblázatok és ábrák segítségével.

Két kategorikus (csoportosított) változó kapcsolatának felterkepezése

Feltáró elemzés

Most vizsgáljuk meg azt, hogy 2020-09-28-an mi az összefüggése a gdp kategorianak (*gdp_per_capita_kat*) a kontinenssel (*continent*) ahol az ország elhelyezkedik.

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

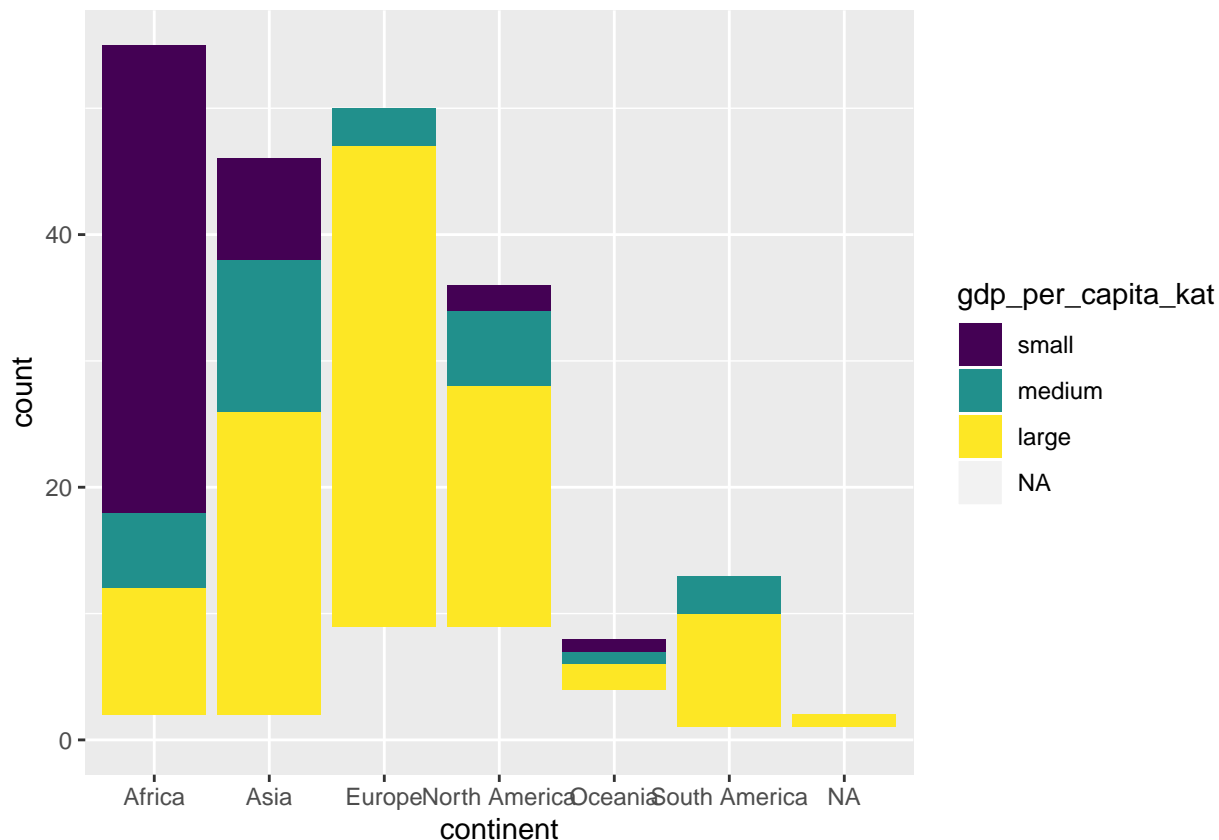
```
table(COVID_adat_tegnap$gdp_per_capita_kat, COVID_adat_tegnap$continent)
```

```
##
##           Africa Asia Europe North America Oceania South America
##  small         37    8      0              2      1              0
##  medium         6   12      3              6      1              3
##  large         10   24     38             19      2              9
```

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

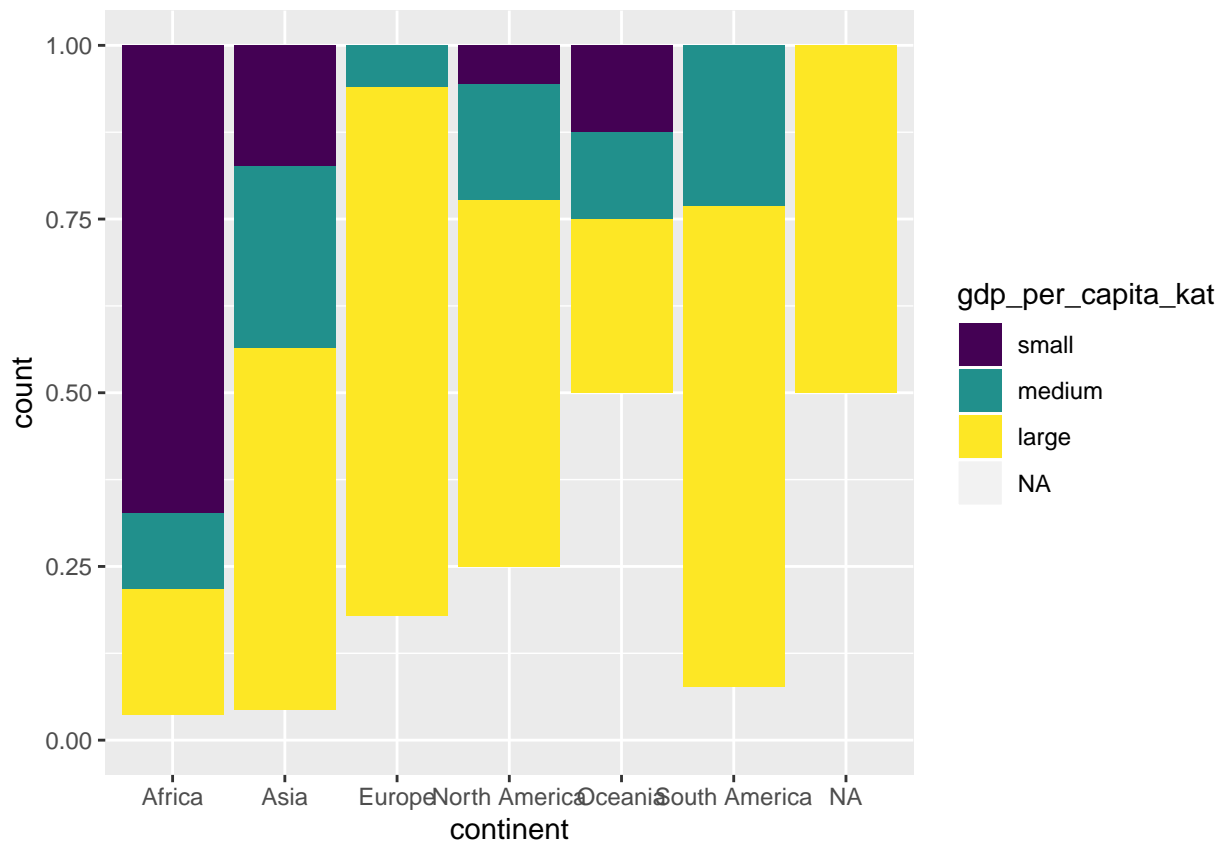
Erre az egyik lehetőség a **stacked bar chart** (egymásra tornyozott oszlopdiagram, a **geom_bar()** geomot használjuk) használata. Itt az egyik változó kategóriái adják meg hány oszlop lesz (ez a változó lesz az x tengelyen reprezentálva, így ezt az "x =" részen adhatjuk meg), a másik változó az oszlopokat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

```
COVID_adat_tegnap %>%
  ggplot() +
    aes(x = continent, fill = gdp_per_capita_kat) +
    geom_bar()
```



Ha az egyes faktorszinteken nagyon **különbozo mennyisegu megfigyeles** van, ez a megjelenites neha felrevezeto kovetkeztetesekekhez vezethet, így neha hasznosabb ha az oszlopok nem szamossagot (count), hanem **reszaranyt (proportion)** jelolnek. Ha ezt szeretnenk, ahelyett hogy uresen hagynank a `geom_bar()` funkciot, a kovetkezozt adjuk meg: `geom_bar(position = "fill")`.

```
COVID_adat_tegnap %>%
  ggplot() +
    aes(x = continent, fill = gdp_per_capita_kat) +
    geom_bar(position = "fill")
```



```
COVID_adat = COVID_adat %>%
  mutate(new_cases_per_million_kat = factor(
    case_when(new_cases_per_million < 20 ~ "small",
              new_cases_per_million >= 20 ~ "large"), ordered = T, levels =
  levels(COVID_adat$new_cases_per_million_kat)

## [1] "small" "large"
```

```
# ugyanez a COVID_adat_tegnap -al

COVID_adat_tegnap = COVID_adat_tegnap %>%
  mutate(new_cases_per_million_kat = factor(
    case_when(new_cases_per_million < 20 ~ "small",
              new_cases_per_million >= 20 ~ "large"), ordered = T, levels =
```

Gyakorlas

Hasznald a fent tanult modszereket, hogy megvizsgald a COVID_adat_tegnap adatbazisban a **new_cases_per_million_kat** es a **continent** 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?

Ennel a megjelenítésnél fontos hogy ha az egyes megfigyelek **keves megfigyelesbol allnak**, az ábra megteveszto lehet, mert az ábra nem jelzi a megfigyelek szamat es így azt, hogy milyen biztosak lehetünk

az eredményben. Ilyen esetekben az egyik kategóriát ki lehet venni az ábráról, vagy a **szamossagot es a reszaranyt abrazolo abrakat egymás mellet** lehet bemutatni, hogy így kiegészítsek egymást. Ehhez használhatjuk a **grid.arrange()** funkciót.