

PSZB17-210 - Seminar_4

Zoltan Kekecs

March 2, 2021

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

```
COVID_data_raw <- read_csv("https://raw.githubusercontent.com/owid/covid-19-data/master/public/data/owid-covid-data.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_data_raw
```

```
## # A tibble: 72,096 x 59
##   iso_code continent location date      total_cases new_cases new_cases_smooth~
##   <chr>      <chr>      <chr>  <date>          <dbl>      <dbl>          <dbl>
## 1 AFG      Asia      Afghani~ 2020-02-24          1          1            NA
## 2 AFG      Asia      Afghani~ 2020-02-25          1          0            NA
## 3 AFG      Asia      Afghani~ 2020-02-26          1          0            NA
## 4 AFG      Asia      Afghani~ 2020-02-27          1          0            NA
## 5 AFG      Asia      Afghani~ 2020-02-28          1          0            NA
## 6 AFG      Asia      Afghani~ 2020-02-29          1          0            0.143
## 7 AFG      Asia      Afghani~ 2020-03-01          1          0            0.143
## 8 AFG      Asia      Afghani~ 2020-03-02          1          0            0
## 9 AFG      Asia      Afghani~ 2020-03-03          2          1            0.143
```

```
## 10 AFG      Asia      Afghani~ 2020-03-04          4          2          0.429
## # ... with 72,086 more rows, and 52 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>, reproduction_rate <dbl>,
## #   icu_patients <lgl>, icu_patients_per_million <lgl>, hosp_patients <lgl>,
## #   hosp_patients_per_million <lgl>, weekly_icu_admissions <lgl>,
## #   weekly_icu_admissions_per_million <lgl>, weekly_hosp_admissions <lgl>,
## #   weekly_hosp_admissions_per_million <lgl>, new_tests <dbl>,
## #   total_tests <dbl>, total_tests_per_thousand <dbl>,
## #   new_tests_per_thousand <dbl>, new_tests_smoothed <dbl>,
## #   new_tests_smoothed_per_thousand <dbl>, positive_rate <dbl>,
## #   tests_per_case <dbl>, tests_units <chr>, total_vaccinations <lgl>,
## #   people_vaccinated <lgl>, people_fully_vaccinated <lgl>,
## #   new_vaccinations <lgl>, new_vaccinations_smoothed <lgl>,
## #   total_vaccinations_per_hundred <lgl>, people_vaccinated_per_hundred <lgl>,
## #   people_fully_vaccinated_per_hundred <lgl>,
## #   new_vaccinations_smoothed_per_million <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_data_raw %>%
  select(total_cases) %>%
  summary()
```

```
##   total_cases
##   Min.      :    1
##   1st Qu.:   706
##   Median :  7577
##   Mean    : 610304
##   3rd Qu.: 82603
##   Max.    :114442646
##   NA's    :964
```

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 kaptjuk az output-ban.

```
COVID_data_raw %>%
  summary()
```

Az exploráció megmutatta hogy van néhány irreálisztikus adat. Ennek az az oka hogy kontinensekre és régiókra lebontott összefoglaló adatokat is tartalmaz a táblázat. Ezeket úgy tudjuk legkönnyebben kivenni

hogy kivesszük azokat a sorokat, ahol a continent változó NA értéket vesz fel. (Vedd észre hogy ezt “!” és az is.na() funktciók kombinációjával oldjuk meg. A ! jelentése “NOT”.)

```
COVID_data <- COVID_data_raw %>%  
  filter(!is.na(continent))
```

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

```
##   total_cases  
##   Min.      :    1  
##   1st Qu.:   628  
##   Median :  6565  
##   Mean   : 203025  
##   3rd Qu.: 67334  
##   Max.    :28664481  
##   NA's    :950
```

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

```
##   total_cases  
##   Min.      :    1  
##   1st Qu.:   706  
##   Median :  7577  
##   Mean   : 610304  
##   3rd Qu.: 82603  
##   Max.    :114442646  
##   NA's    :964
```

Gyakorlas

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

Megtöbb leíró statisztika

A **Psych** csomag segítségével a **describe()** funkció még 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_data %>%  
  select(-date, -iso_code, -continent, -location, -contains("tests"), -positive_rate) %>%  
  describe()
```

##		vars	n	mean	sd
##	total_cases	1	67566	203024.59	1139284.95
##	new_cases	2	67559	1681.76	9159.12
##	new_cases_smoothed	3	66611	1689.43	9008.48
##	total_deaths	4	58643	6075.65	25639.29

## new_deaths	5	58642	43.29	185.58
## new_deaths_smoothed	6	66611	37.76	163.03
## total_cases_per_million	7	67566	7153.58	14314.81
## new_cases_per_million	8	67559	64.78	168.79
## new_cases_smoothed_per_million	9	66611	64.77	141.33
## total_deaths_per_million	10	58643	168.05	306.43
## new_deaths_per_million	11	58642	1.36	3.93
## new_deaths_smoothed_per_million	12	66611	1.19	2.80
## reproduction_rate	13	57602	1.02	0.35
## icu_patients	14	0	NaN	NA
## icu_patients_per_million	15	0	NaN	NA
## hosp_patients	16	0	NaN	NA
## hosp_patients_per_million	17	0	NaN	NA
## weekly_icu_admissions	18	0	NaN	NA
## weekly_icu_admissions_per_million	19	0	NaN	NA
## weekly_hosp_admissions	20	0	NaN	NA
## weekly_hosp_admissions_per_million	21	0	NaN	NA
## total_vaccinations	22	0	NaN	NA
## people_vaccinated	23	0	NaN	NA
## people_fully_vaccinated	24	0	NaN	NA
## new_vaccinations	25	0	NaN	NA
## new_vaccinations_smoothed	26	0	NaN	NA
## total_vaccinations_per_hundred	27	0	NaN	NA
## people_vaccinated_per_hundred	28	0	NaN	NA
## people_fully_vaccinated_per_hundred	29	0	NaN	NA
## new_vaccinations_smoothed_per_million	30	0	NaN	NA
## stringency_index	31	61934	58.92	22.15
## population	32	68516	43599575.60	156432943.82
## population_density	33	66980	328.83	1596.52
## median_age	34	65184	30.57	9.15
## aged_65_older	35	64428	8.81	6.27
## aged_70_older	36	64814	5.59	4.28
## gdp_per_capita	37	65362	19156.64	19740.78
## extreme_poverty	38	44519	13.28	20.00
## cardiovasc_death_rate	39	65973	257.23	118.73
## diabetes_prevalence	40	66800	7.79	3.94
## female_smokers	41	51834	10.59	10.45
## male_smokers	42	51117	32.63	13.51
## handwashing_facilities	43	32902	50.91	31.95
## hospital_beds_per_thousand	44	60413	3.04	2.48
## life_expectancy	45	68103	73.14	7.57
## human_development_index	46	65951	0.73	0.15
##		median	trimmed	mad
## total_cases		6565.00	35972.93	9699.17
## new_cases		50.00	289.35	74.13
## new_cases_smoothed		60.71	306.87	89.80
## total_deaths		202.00	1118.22	292.07
## new_deaths		1.00	7.54	1.48
## new_deaths_smoothed		0.86	6.03	1.27
## total_cases_per_million		996.97	3495.34	1463.16
## new_cases_per_million		5.55	26.75	8.23
## new_cases_smoothed_per_million		7.62	29.95	11.27
## total_deaths_per_million		28.54	89.88	41.08
## new_deaths_per_million		0.09	0.54	0.13

## new_deaths_smoothed_per_million	0.10	0.50	0.15	
## reproduction_rate	1.03	1.02	0.24	
## icu_patients	NA	NaN	NA	
## icu_patients_per_million	NA	NaN	NA	
## hosp_patients	NA	NaN	NA	
## hosp_patients_per_million	NA	NaN	NA	
## weekly_icu_admissions	NA	NaN	NA	
## weekly_icu_admissions_per_million	NA	NaN	NA	
## weekly_hosp_admissions	NA	NaN	NA	
## weekly_hosp_admissions_per_million	NA	NaN	NA	
## total_vaccinations	NA	NaN	NA	
## people_vaccinated	NA	NaN	NA	
## people_fully_vaccinated	NA	NaN	NA	
## new_vaccinations	NA	NaN	NA	
## new_vaccinations_smoothed	NA	NaN	NA	
## total_vaccinations_per_hundred	NA	NaN	NA	
## people_vaccinated_per_hundred	NA	NaN	NA	
## people_fully_vaccinated_per_hundred	NA	NaN	NA	
## new_vaccinations_smoothed_per_million	NA	NaN	NA	
## stringency_index	61.11	60.24	23.34	
## population	9660350.00	16582763.24	13023827.05	
## population_density	83.48	113.57	88.35	
## median_age	29.70	30.44	12.75	
## aged_65_older	6.38	8.14	5.06	
## aged_70_older	3.86	5.04	3.12	
## gdp_per_capita	12951.84	15810.25	14552.01	
## extreme_poverty	2.00	8.89	2.67	
## cardiovasc_death_rate	242.65	246.67	121.48	
## diabetes_prevalence	7.11	7.42	3.45	
## female_smokers	6.30	9.21	7.86	
## male_smokers	31.40	32.00	14.53	
## handwashing_facilities	49.84	51.14	45.66	
## hospital_beds_per_thousand	2.40	2.67	1.93	
## life_expectancy	74.62	73.69	7.06	
## human_development_index	0.75	0.74	0.17	
##	min	max	range	skew
## total_cases	1.00	2.866448e+07	2.866448e+07	14.41
## new_cases	-46076.00	2.997860e+05	3.458620e+05	15.23
## new_cases_smoothed	-1121.71	2.497266e+05	2.508483e+05	15.12
## total_deaths	1.00	5.146570e+05	5.146560e+05	9.49
## new_deaths	-1918.00	4.398000e+03	6.316000e+03	10.01
## new_deaths_smoothed	-232.14	3.354140e+03	3.586290e+03	9.64
## total_cases_per_million	0.00	1.409306e+05	1.409306e+05	3.30
## new_cases_per_million	-2153.44	8.652660e+03	1.080609e+04	9.15
## new_cases_smoothed_per_million	-276.82	2.648770e+03	2.925600e+03	4.07
## total_deaths_per_million	0.00	2.180450e+03	2.180450e+03	2.61
## new_deaths_per_million	-76.44	2.183300e+02	2.947700e+02	12.66
## new_deaths_smoothed_per_million	-10.92	6.314000e+01	7.406000e+01	4.79
## reproduction_rate	0.00	6.740000e+00	6.740000e+00	0.96
## icu_patients	Inf	-Inf	-Inf	NA
## icu_patients_per_million	Inf	-Inf	-Inf	NA
## hosp_patients	Inf	-Inf	-Inf	NA
## hosp_patients_per_million	Inf	-Inf	-Inf	NA
## weekly_icu_admissions	Inf	-Inf	-Inf	NA

## weekly_icu_admissions_per_million	Inf	-Inf	-Inf	NA
## weekly_hosp_admissions	Inf	-Inf	-Inf	NA
## weekly_hosp_admissions_per_million	Inf	-Inf	-Inf	NA
## total_vaccinations	Inf	-Inf	-Inf	NA
## people_vaccinated	Inf	-Inf	-Inf	NA
## people_fully_vaccinated	Inf	-Inf	-Inf	NA
## new_vaccinations	Inf	-Inf	-Inf	NA
## new_vaccinations_smoothed	Inf	-Inf	-Inf	NA
## total_vaccinations_per_hundred	Inf	-Inf	-Inf	NA
## people_vaccinated_per_hundred	Inf	-Inf	-Inf	NA
## people_fully_vaccinated_per_hundred	Inf	-Inf	-Inf	NA
## new_vaccinations_smoothed_per_million	Inf	-Inf	-Inf	NA
## stringency_index	0.00	1.000000e+02	1.000000e+02	-0.48
## population	809.00	1.439324e+09	1.439323e+09	7.88
## population_density	0.14	2.054677e+04	2.054663e+04	10.35
## median_age	15.10	4.820000e+01	3.310000e+01	0.10
## aged_65_older	1.14	2.705000e+01	2.591000e+01	0.78
## aged_70_older	0.53	1.849000e+01	1.797000e+01	0.91
## gdp_per_capita	661.24	1.169356e+05	1.162744e+05	1.81
## extreme_poverty	0.10	7.760000e+01	7.750000e+01	1.65
## cardiovasc_death_rate	79.37	7.244200e+02	6.450500e+02	0.87
## diabetes_prevalence	0.99	3.053000e+01	2.954000e+01	1.16
## female_smokers	0.10	4.400000e+01	4.390000e+01	0.95
## male_smokers	7.70	7.810000e+01	7.040000e+01	0.54
## handwashing_facilities	1.19	9.900000e+01	9.781000e+01	-0.05
## hospital_beds_per_thousand	0.10	1.380000e+01	1.370000e+01	1.75
## life_expectancy	53.28	8.675000e+01	3.347000e+01	-0.59
## human_development_index	0.39	9.600000e-01	5.600000e-01	-0.38
##	kurtosis	se		
## total_cases	271.15	4382.97		
## new_cases	314.39	35.24		
## new_cases_smoothed	303.00	34.90		
## total_deaths	123.39	105.88		
## new_deaths	145.38	0.77		
## new_deaths_smoothed	130.45	0.63		
## total_cases_per_million	13.84	55.07		
## new_cases_per_million	250.56	0.65		
## new_cases_smoothed_per_million	23.55	0.55		
## total_deaths_per_million	7.15	1.27		
## new_deaths_per_million	400.14	0.02		
## new_deaths_smoothed_per_million	38.91	0.01		
## reproduction_rate	11.41	0.00		
## icu_patients	NA	NA		
## icu_patients_per_million	NA	NA		
## hosp_patients	NA	NA		
## hosp_patients_per_million	NA	NA		
## weekly_icu_admissions	NA	NA		
## weekly_icu_admissions_per_million	NA	NA		
## weekly_hosp_admissions	NA	NA		
## weekly_hosp_admissions_per_million	NA	NA		
## total_vaccinations	NA	NA		
## people_vaccinated	NA	NA		
## people_fully_vaccinated	NA	NA		
## new_vaccinations	NA	NA		

```
## new_vaccinations_smoothed      NA      NA
## total_vaccinations_per_hundred NA      NA
## people_vaccinated_per_hundred  NA      NA
## people_fully_vaccinated_per_hundred NA      NA
## new_vaccinations_smoothed_per_million NA      NA
## stringency_index               -0.42    0.09
## population                     65.44 597629.76
## population_density             115.04    6.17
## median_age                    -1.25    0.04
## aged_65_older                 -0.69    0.02
## aged_70_older                 -0.36    0.02
## gdp_per_capita                 4.14    77.22
## extreme_poverty               1.68    0.09
## cardiovasc_death_rate         0.77    0.46
## diabetes_prevalence           2.42    0.02
## female_smokers                 -0.17    0.05
## male_smokers                   0.26    0.06
## handwashing_facilities        -1.50    0.18
## hospital_beds_per_thousand     3.97    0.01
## life_expectancy               -0.40    0.03
## human_development_index       -0.89    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.

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

```
levels(COVID_data$continent)
```

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

```
table(COVID_data$continent)
```

```
##
##      Africa      Asia      Europe North America      Oceania
##      18950     17244     17439      8444      2035
## South America
##      4404
```

```
COVID_data <- COVID_data %>%
  mutate(continent = factor(continent))
```

A `levels()` funkció megmutatja mik a faktorunk szintjei, de látható ez akkor is ha csak meghívjuk a változót magát.

A `table()` funkció pedig táblázatot készít arról, hogy az egyes csoportokban hány megfigyelés található. Amikor klistazzuk a faktor változót, akkor is kiírja az R a lista aljára, hogy milyen faktorszintek vannak.

```
levels(COVID_data$continent)

table(COVID_data$continent)

COVID_data$continent
```

Alább csinálunk egy `COVID_data_latest` változót, amivel csak az adatbázisban szereplő legutolsó napra vonatkozó adatok szerepelnek, hogy kisebb legyen az adattábla amivel dolgozunk.

```
COVID_data_latest = COVID_data %>%
  filter(date == max(COVID_data$date))
```

Miután egy változót faktorként azonosítottunk, bizonyos funkciók képesek felhasználni ezt az információt. Például a `summary()` function így már a fenti `summary()` funkció is kiadja az **egyes faktorszintekről** hogy hány megfigyelés tartozik az egyes kategóriákba (faktorszintekbe).

```
COVID_data %>%
  mutate(continent = as.character(continent)) %>%
  select(continent) %>%
  summary()
```

```
##   continent
## Length:68516
## Class :character
## Mode  :character
```

```
# continent is already recognized as a factor variable
COVID_data_latest %>%
  select(continent) %>%
  summary()
```

```
##           continent
## Africa           :54
## Asia             :46
## Europe           :47
## North America:23
## Oceania          : 9
## South America:12
```

Van, hogy szeretnénk **kizárni** bizonyos **faktorszinteket** az elemzésből. Pl. ha valamelyik faktor szintből nagyon keves megfigyelés van, mondjuk Oceaniát, mondjuk mert úgy gondoljuk hogy az tulságosan “elszigetelt” a világ többi részétől, okét lehet hogy szeretnénk kizárni a későbbi elemzésekből hogy egyszerűsítsuk az eredményeink értelmezését. Ezt a már korábban tanult `filter()` funkció segítségével könnyedén megtehetjük, azonban arra figyelni kell, hogy az R megjegyzi a faktorszinteket, és azt azt követően is a **változéhoz rendelve tartja**. A **faktorszintek meg akkor is megmaradnak ha nem marad egy megfigyelés sem** az adott faktorszinten az adattáblában.

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



```
##   total_cases      continent
##   Min.      :    27   Africa      :54
##   1st Qu.:   8595   Asia         :46
##   Median :  79002   Europe      :47
##   Mean    : 632095   North America:23
##   3rd Qu.: 288267   Oceania      : 0
##   Max.    :28664481   South America:12
##   NA's     :1
```

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

```
COVID_data_latest_noOceania = COVID_data_latest %>%
  filter(continent != "Oceania") %>%
  mutate(continent = droplevels(continent))
```

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

```
##           continent
##   Africa      :54
##   Asia        :46
##   Europe      :47
##   North America:23
##   South America:12
```

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 végzettség nélküli < általános iskolai < közepiskolai < felsőfokú ... Ittfaktorszinteknek van egy meghatározott hierarchiaja, vagy sorrendje. Ilyen változó nincs ebben az adatbázisban, de könnyedén csinálhatunk ilyen faktor változót.

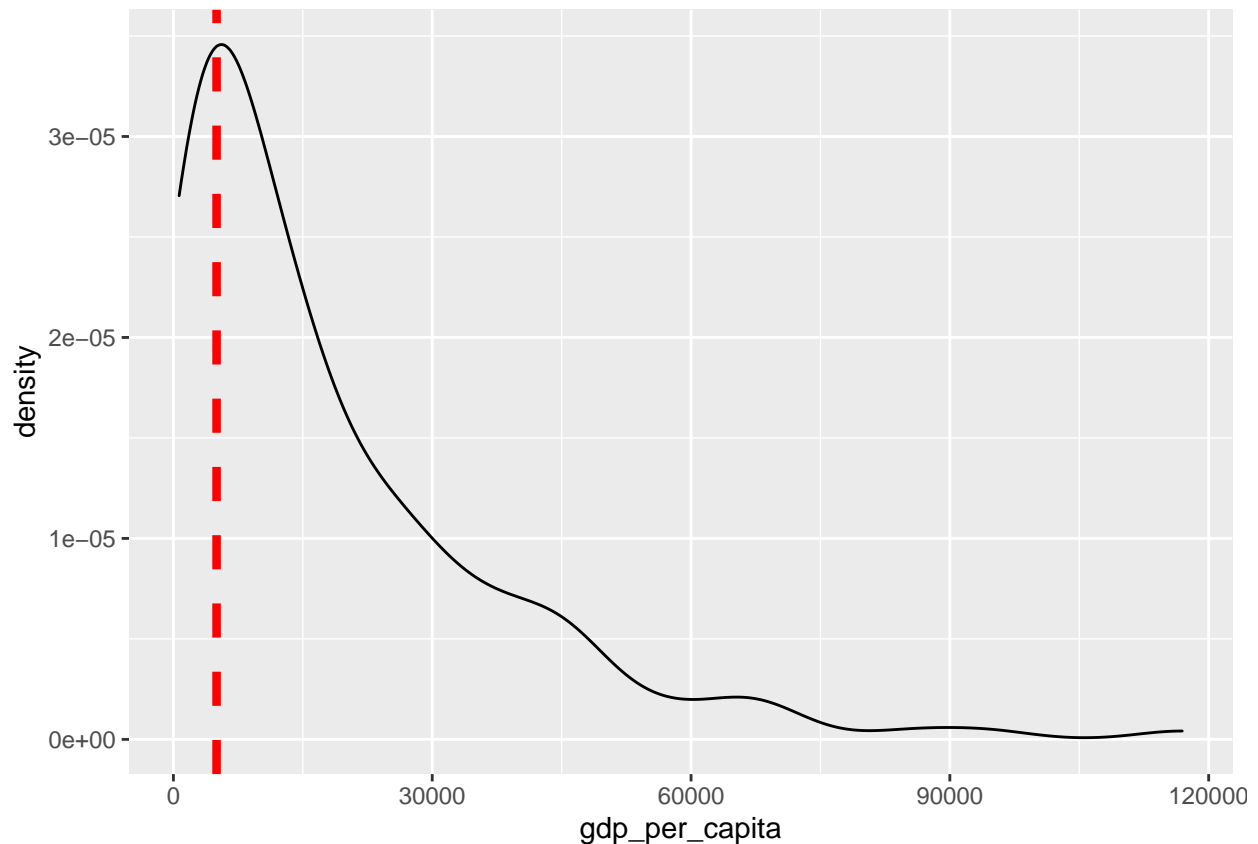
Ehhez arra van szükségünk, hogy egy **numerikus változót alakítsunk faktorra**, pl. elképzelhető hogy össze akarjuk hasonlítani azokat az országokat ahol 5000 alatti a `gdp_per_capita` azokkal akinek e feletti, hogy hogyan különböznek a COVID adatok.

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

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

```
## # A tibble: 6 x 2
##   continent    mean_gdp
##   <fct>         <dbl>
## 1 Africa         5444.
## 2 Asia          22185.
## 3 Europe        33361.
## 4 North America 17126.
## 5 Oceania       12392.
## 6 South America 13841.
```

```
COVID_data_latest %>%
  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)
```



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 meg kellene egy külön sor ahol megadjuk hogy ez egy faktorváltozó.

```
COVID_data = COVID_data %>%
  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_data$gdp_per_capita_kat)

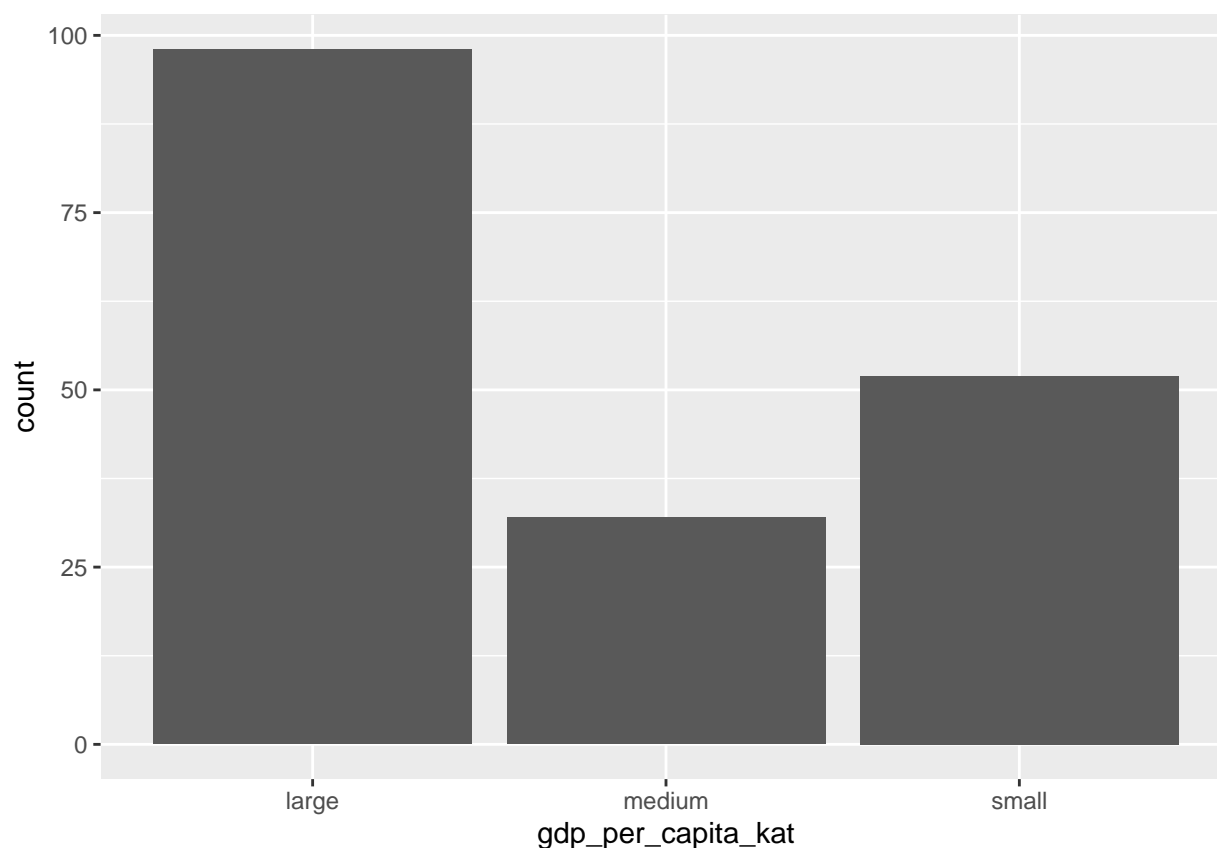
## [1] "large" "medium" "small"
# ugyanez a COVID_data_latest -al

COVID_data_latest = COVID_data_latest %>%
  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"))
```

Amikor abrat rajzolunk erreol a változóra, láthatjuk hogy a faktorszintek sorrendje “large”, “medium”, és “small” az abran.

```
COVID_data_latest %>%
  select(gdp_per_capita_kat) %>%
  drop_na() %>%
  ggplot() +
  aes(x = gdp_per_capita_kat) +
  geom_bar()
```



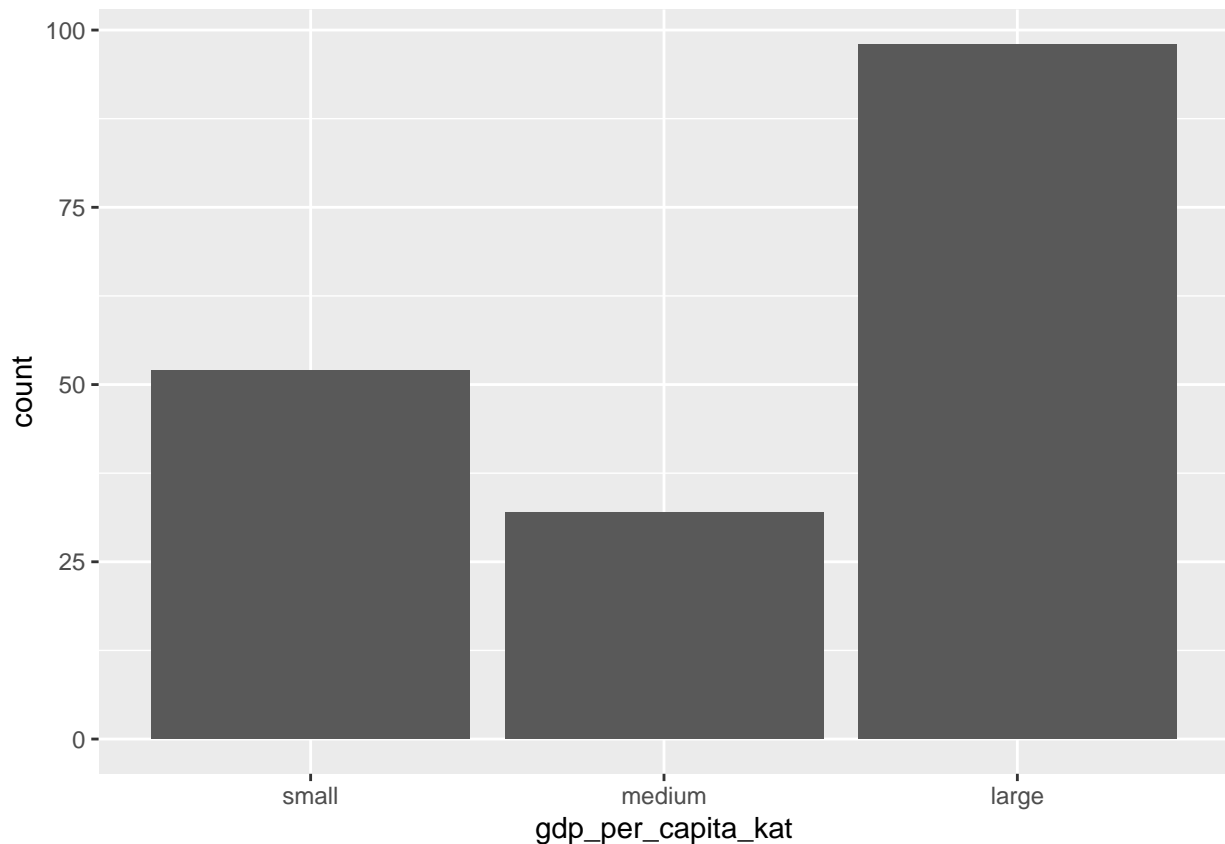
Ez nem feltétlenül legikus ábrázolás, hiszen általában a kisebbtől a nagyobbig szoktunk haladni balról jobbra. De az R nem tudja mit jelentenek a faktorszintek nevei. A faktorszintek sorrendjének meghatározásánál ezért alapértelmezett módon **abc sorrendet** használ.

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

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

```
COVID_data_latest %>%
```

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



Attól meg hogy megadjuk a levels-el a faktorszintek listazasi sorrendjet, az R meg mindig egyenrangukent kezeli a faktorszinteket. Ha azt szeretnénk ha az R ugy ertelmeze hogy a faktorszintek valamilyen hierarchikus sorrendben van, vagyis **ordinalis változók**, akkor ezt a factor() funkcion belül az **ordered = T** parameter beallitasaval tehetjuk meg.

Ha ezt teszük, a faktor változó kilistazasakor relacio-jelek kerülnek a faktorszintek köze, és más funkiók is fel tudják majd használni ezt az információt.

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

```
## [1] small large large <NA> medium large large medium large large
## [11] large large large small large large large medium small medium
## [21] medium large large large large large large small small small
## [31] large medium small small large large large small small large
## [41] small large <NA> large large small large small medium large
## [51] large large medium large small large medium small medium large
```

```
## [61] large large small medium large small large large medium small
## [71] small medium small small large large medium large large large
## [81] large <NA> large large medium large medium large small medium
## [91] large small medium large large small small large <NA> large
## [101] large small small large large small large small small large
## [111] large small medium <NA> large large medium small medium medium
## [121] small large large medium small medium large large large medium
## [131] small large small medium large medium large large large large
## [141] large small large large large medium large small large small
## [151] large large small large large large small <NA> large large
## [161] small large large small large large large <NA> <NA> small
## [171] small large medium small large large large small medium large
## [181] large large large medium small <NA> large medium small small
## [191] small
## Levels: small < medium < large
```

Kategorikus változó ujrakodolása

Egy másik funkció amivel manipulálhatjuk a faktorszinteket, a **recode()**. Ha kategorikus változókat szeretnénk átkodolni, mondjuk ha szeretnénk a déli felteket az északi feltekeivel összehasonlítani, ezt a következőképpen tehetjük:

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

levels(COVID_data$continent_south_north)
```

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

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

Gyakorlás

- szurd az adatokat úgy hogy csak a tegnapi adatokkal dolgozzunk.
- csinálj egy új kategorikus változót (nevezzük ezt *new_cases_per_million_kat*-nak) a `mutate()` funkció használatával amiben azok az országok ahol a *new_cases_per_million* változó 20 alatt van “small”, ahol 20 vagy a felett van “large” kategóriába kerüljenek.
- figyelj 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()` funkciót érdemes hozzá használni)
- mentsd el ezt a változót az eredeti adatobjektumban úgy hogy később is lehessen vele dolgozni

- készíts egy táblázatot arról, hogy hányan esnek a *new_cases_per_million_kat* egyes kategóriaiba.
- Add meg a faktorszintek helyes sorrendjét: small, large (Írd felül a *new_cases_per_million_kat* korábbi változatát ezzel a változattal ahol a szintek már helyes sorrendben vannak, vagy ezt a sorrendezést is bele vonhatod az eredeti funkcioba, amivel a változót generáltad)
- Ellenorizd, hogy valóban helyes sorrendben szerepelnek-e a faktor szintjei.

Exploracio vizualizacion keresztul

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

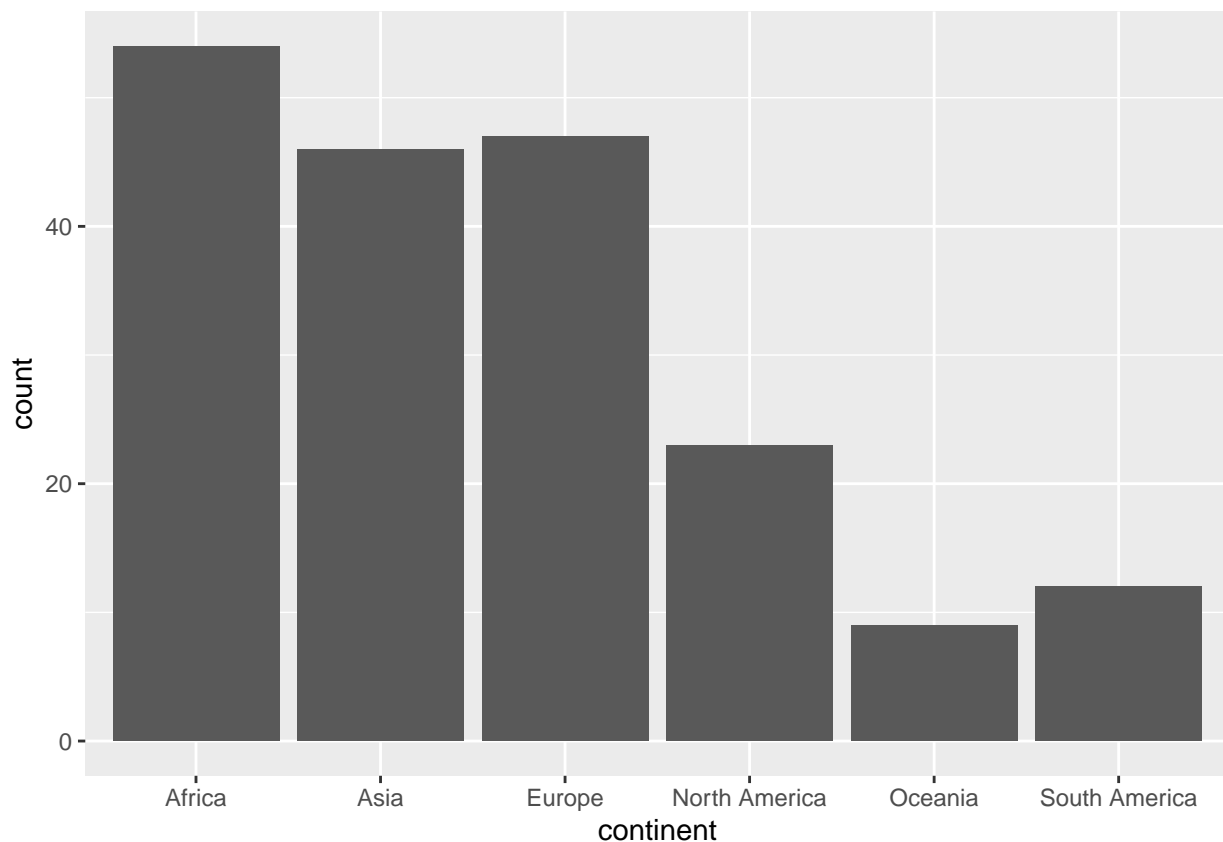
MINDING ellenorizd az adataidat ezekkel a módszerekkel mielőtt komolyabb adatelemzésbe kezdesz, hogy meggyőződj róla, hogy az adatok tiszták és megfelelnek az elvárásaidnak.

Egyes változók vizualizacioja

Az egyes változók például **abrák** (plot) segítségével megvizsgálhatók.

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

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

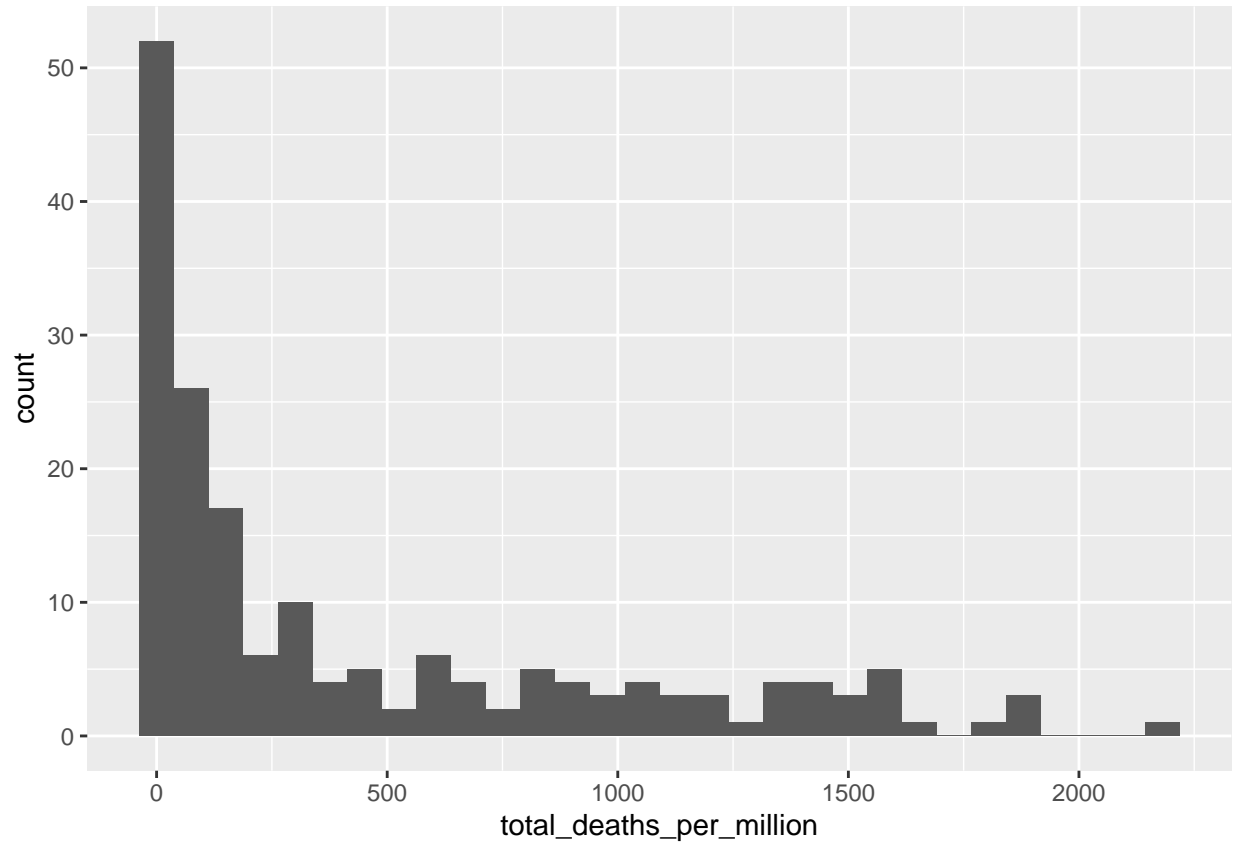


```
COVID_data_latest %>%
  ggplot() +
```

```
aes(x = total_deaths_per_million) +  
geom_histogram()
```

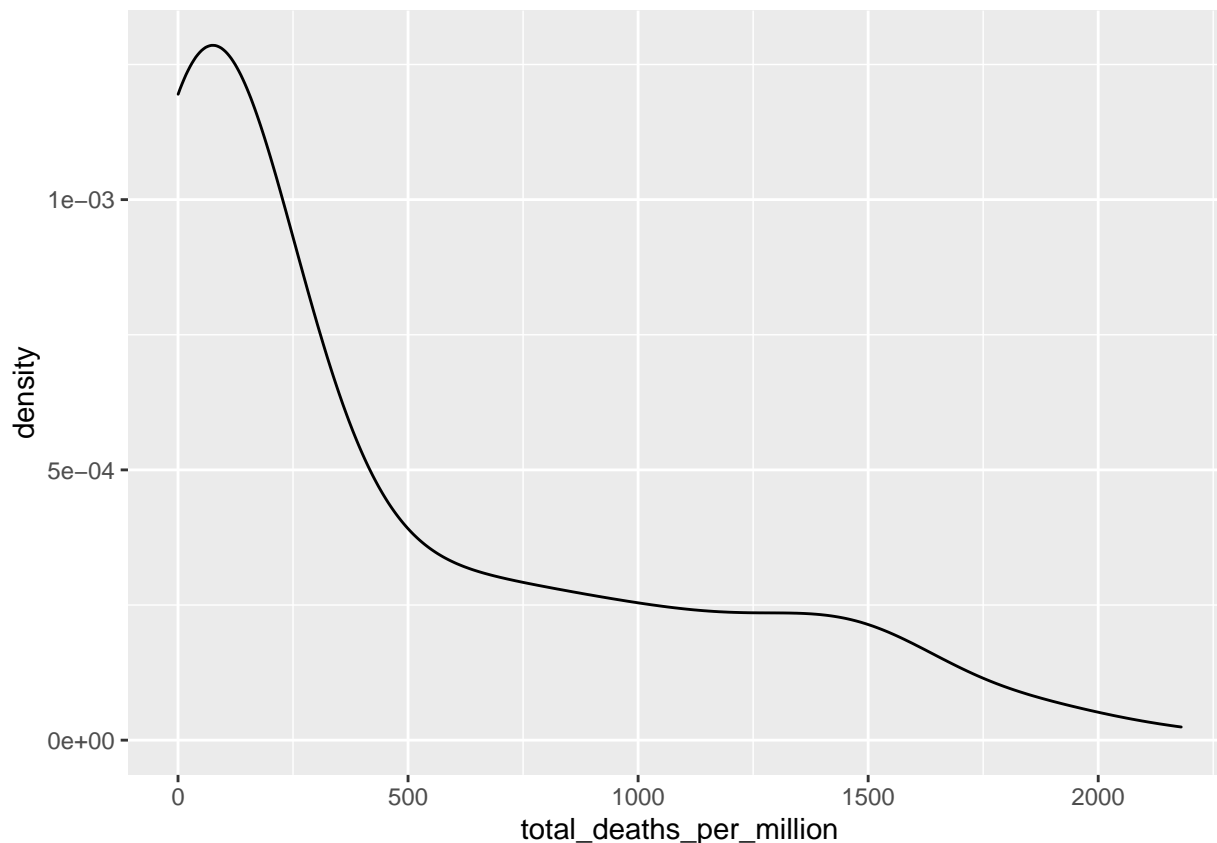
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

```
## Warning: Removed 12 rows containing non-finite values (stat_bin).
```



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

```
## Warning: Removed 12 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_data 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!

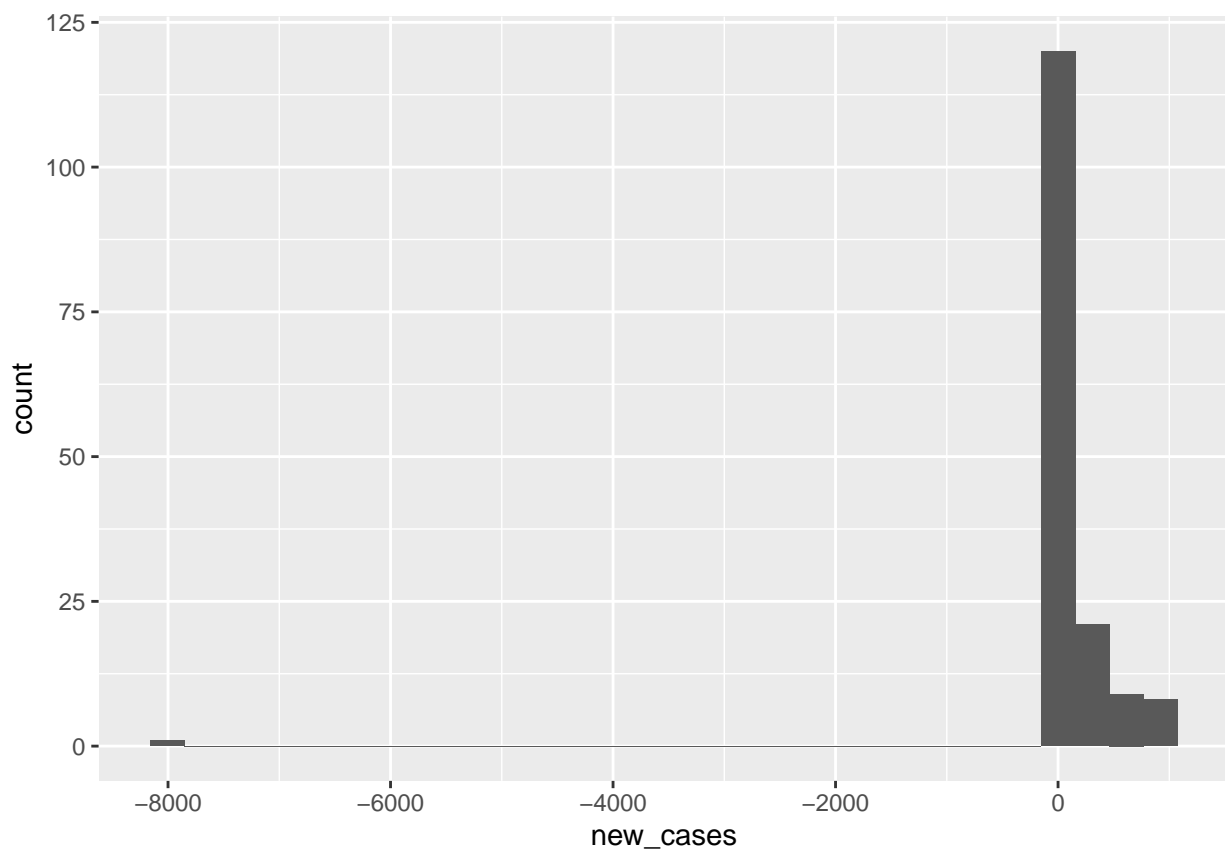
Mindig érdemes a javított adatokat **új adattáblába** elmenteni. A mi esetünkben az COVID_data_corrected nevet adtuk a javított objektumnak. Így a nyers adataink megmaradnak, ami hasznos lehet későbbi

muveleteknel.

```
COVID_data %>%  
  filter(date == "2020-09-07") %>%  
  select(new_cases) %>%  
  summary()
```

```
##    new_cases  
##  Min.     :-7953  
## 1st Qu.:    4  
## Median :   62  
## Mean   : 1177  
## 3rd Qu.:  400  
## Max.   :75809  
## NA's   :1
```

```
COVID_data %>%  
  filter(date == "2020-09-07", new_cases < 1000) %>%  
  ggplot()+  
    aes(x = new_cases)+  
    geom_histogram()
```



```
COVID_data_corrected <- COVID_data %>%  
  mutate(new_cases = replace(new_cases, new_cases=="-7953", NA))
```

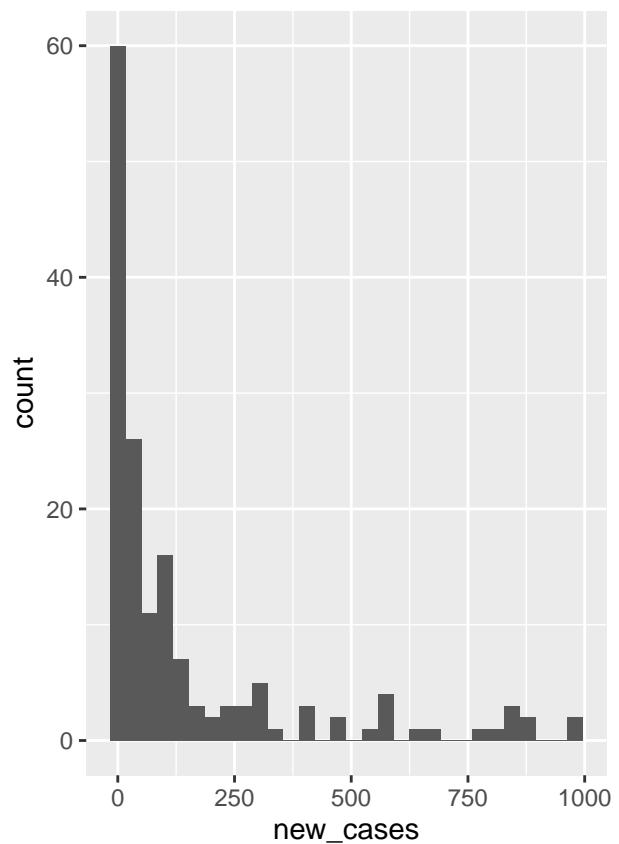
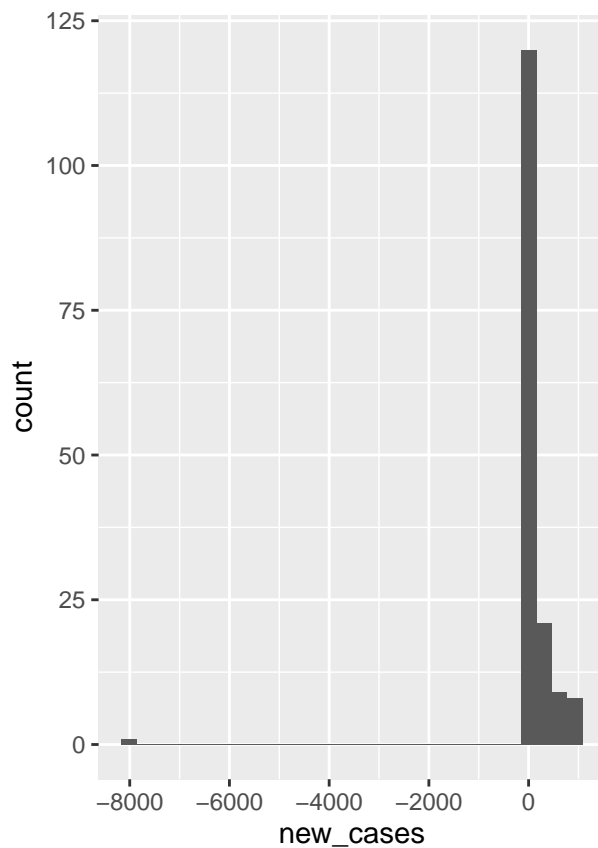
Erdemes **megbizonyosodni rola**, hogy az adatsere sikeres volt. Alabb az adatok vizualizaciojaval gyozodunk meg errol, de az adatok megjelenitesevel, vagy a leiro statisztikak lekerdezesevel is megtehető ez, ha az informatív.

```
# hasznalhatnak meg az alabbiakat is arra,
# hogy megbizonyosodjunk abban, hogy sikeres volt a csere
# View(COVID_data_corrected)
# describe(COVID_data_corrected)
# summary(COVID_data_corrected$szocmedia_3)
# COVID_data_corrected$szocmedia_3
```

```
old_plot <-
  COVID_data %>%
  filter(date == "2020-09-07", new_cases < 1000) %>%
  ggplot()+
  aes(x = new_cases)+
  geom_histogram()
```

```
new_plot <-
  COVID_data_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 kapcsolatot 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)** elkezdése a **table()** funkcióval.

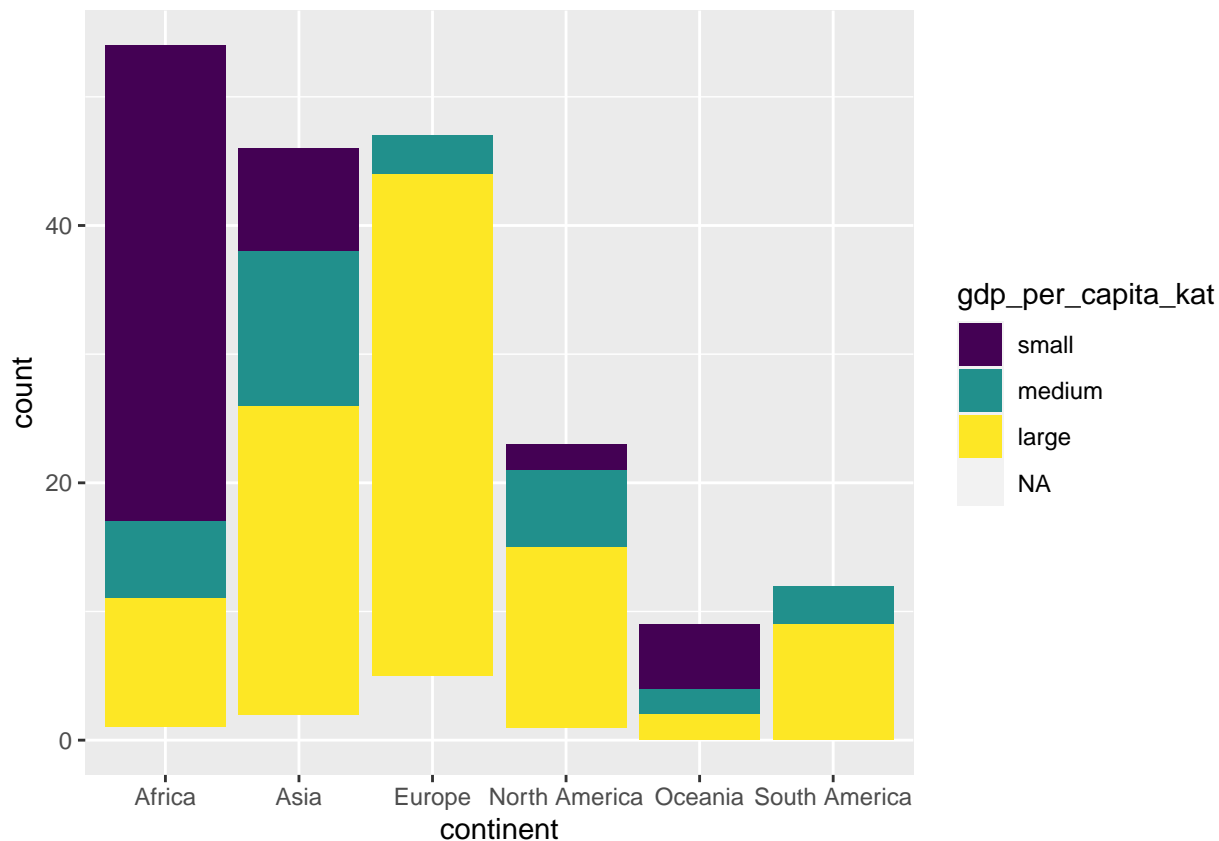
```
table(COVID_data_latest$gdp_per_capita_kat, COVID_data_latest$continent)
```

```
##
##           Africa Asia Europe North America Oceania South America
##   small         37    8      0             2         5             0
##   medium         6   12      3             6         2             3
##   large         10   24     39            14         2             9
```

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

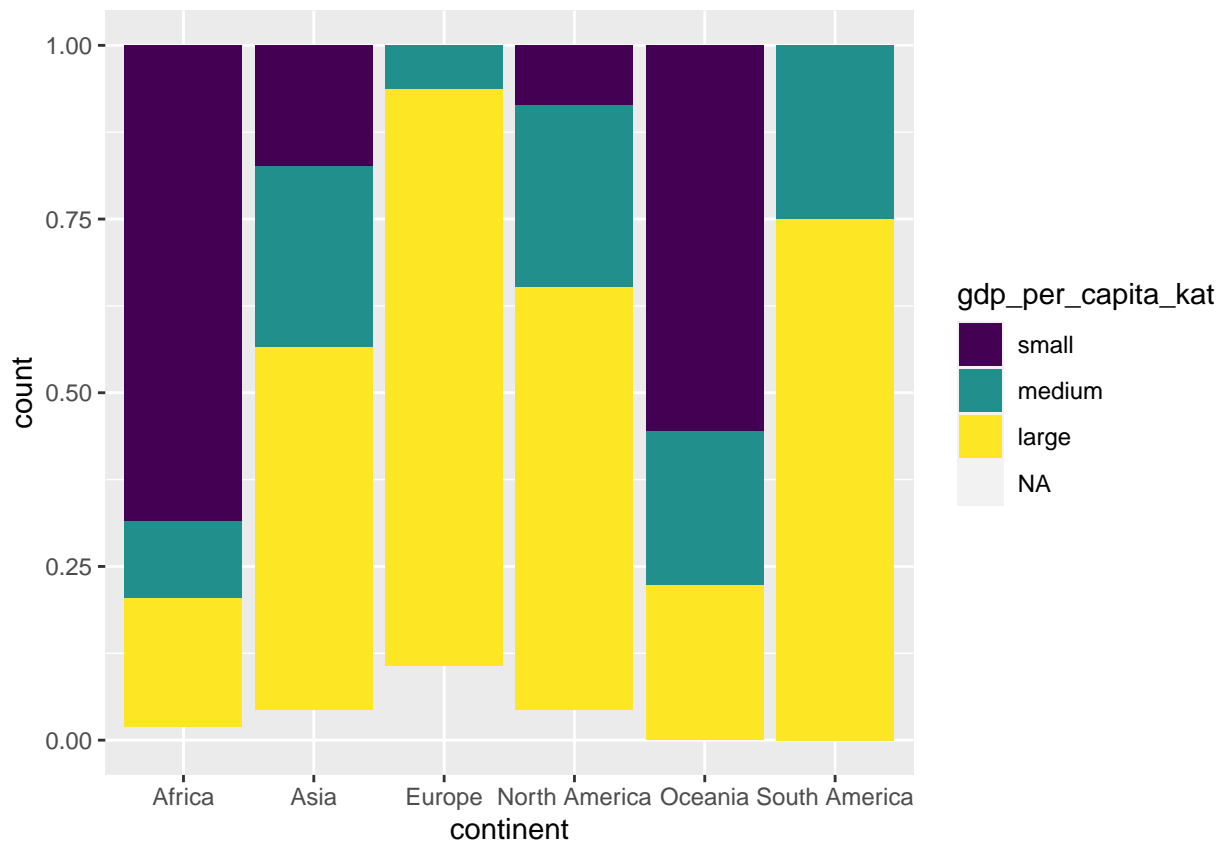
Erre az egyik lehetőség a **stacked bar chart** (egymásra tornyozott oszlopdiagram, a **geom_bar()** geomot használjuk) használata. Itt az egyik változó kategoriai adják meg hány oszlop lesz (ez a változó lesz az x tengelyen reprezentálva, így ezt az “x =” részen adhatjuk meg), a másik változó az oszlopokat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

```
COVID_data_latest %>%
  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 kovetkeztetesekehez 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 kovetkezo adjuk meg: `geom_bar(position = "fill")`. Vagy hasznalhatjuk az eltolt oszlopdiagramot (dodged barchart) (a `position = "dodge"` parameter megadasaval a `geom_bar()` -on belül)

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



Gyakorlas

Hasznald a fent tanult módszereket, hogy megvizsgald a COVID_data_latest adatbázisban a **new_cases_per_million_kat** és 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?

a fenti gyakorlashoz a new_cases_per_million_kat változót így lehet legeneralni:

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

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

ugyanez a COVID_data_latest -al

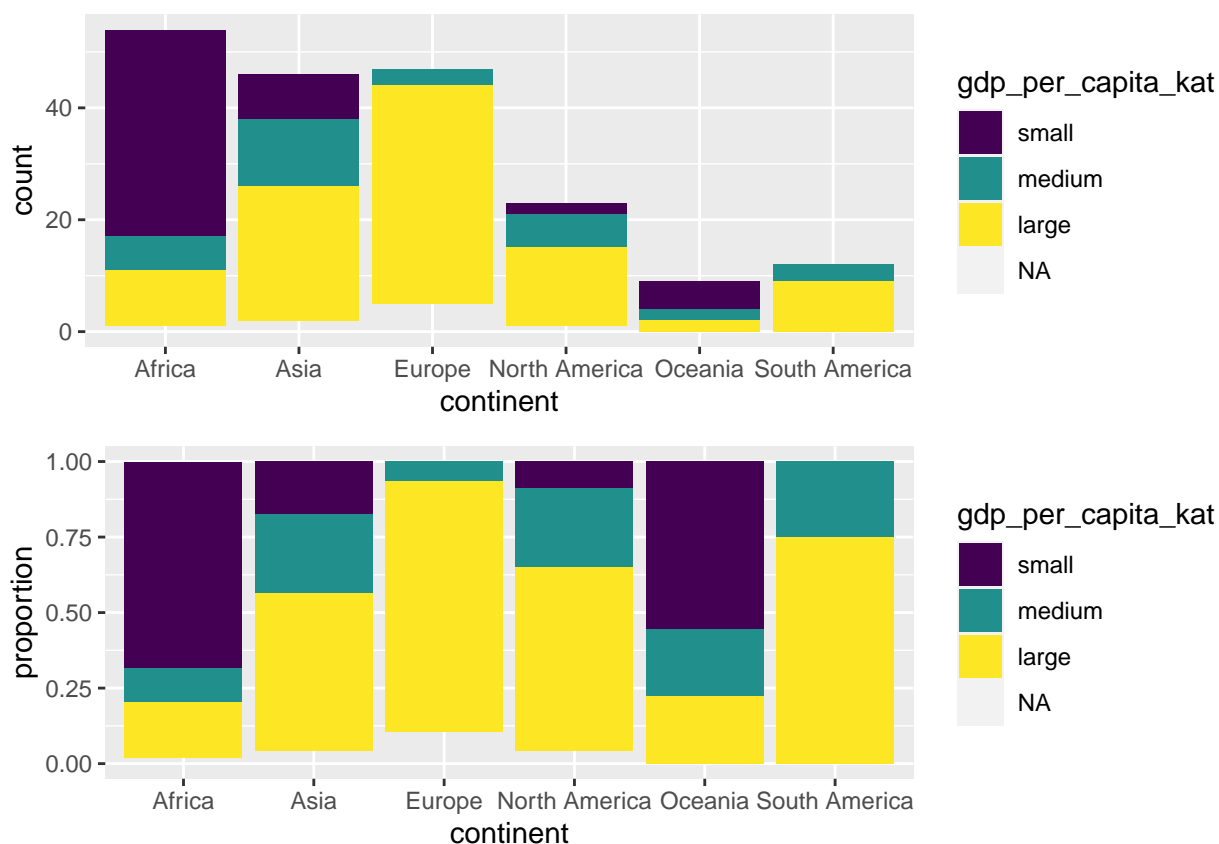
```
COVID_data_latest = COVID_data_latest %>%
  mutate(new_cases_per_million_kat = factor(
    case_when(new_cases_per_million < 20 ~ "small",
              new_cases_per_million >= 20 ~ "large"), ordered = T, levels(COVID_data_latest$new_cases_per_million_kat))
```

`geom_bar()` megjelenítésnél fontos hogy ha az egyes megfigyelesek **keves megfigyelesbol allnak**, az abra megteveszto lehet, mert az abra nem jelzi a megfigyelesek szamat es ilyet, hogy milyen biztosak lehetunk az eredményben. Ilyen esetekben az egyik kategoriat ki lehet venni az abrarol, vagy a **szamossagot es a reszaranyt abrazolo abrakat egymás mellet** lehet bemutatni, hogy ilyet kiegeszitse egymast. Ehhez hasznalhatjuk a `grid.arrange()` funkciot.

```
szamossag_plot <-
COVID_data_latest %>%
ggplot() +
  aes(x = continent, fill = gdp_per_capita_kat) +
  geom_bar()

reszarany_plot <-
COVID_data_latest %>%
ggplot() +
  aes(x = continent, fill = gdp_per_capita_kat) +
  geom_bar(position = "fill") +
  ylab("proportion")

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



A `theme(legend.position)` es a `guides()` funckiok hasznalataval kontrollalhatjuk hogy hol es hogyan jelenjen meg a **jelmagyazat** az abran. Az abra **interpretalhatosaga** attol fuggoen is **valtozhat**, hogy melyik valtozot tesszuk az x-tengelyre es melyiket szinkent abrazolva.

Az alabbi abrakon az egymillio fore vetitett uj esetek szamanak kapcsolatát nezzuk meg a gdp-vel. Mindket valtozo eseten a csoportosított valtozot (`_kat`) hasznaljuk.

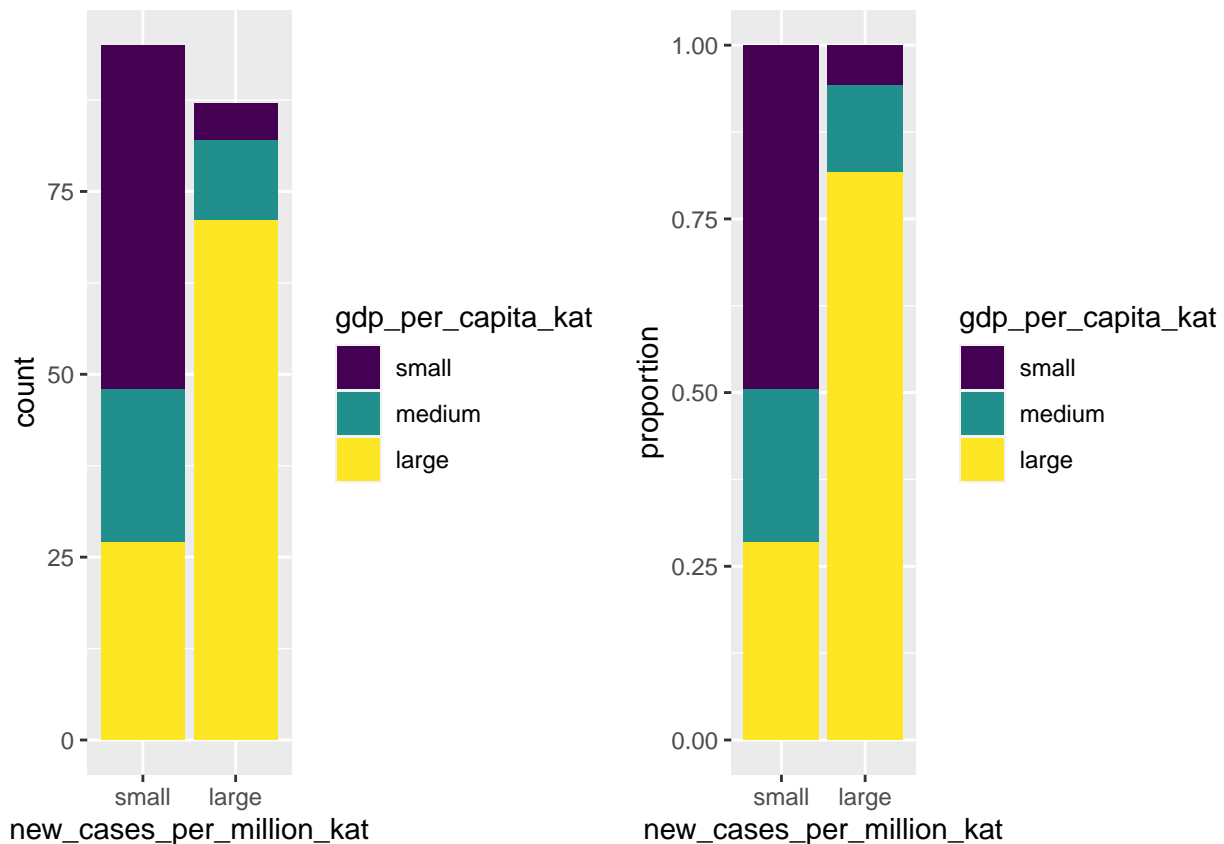
```

barchart_plot_3 <-
COVID_data_latest %>%
  select(new_cases_per_million_kat, gdp_per_capita_kat) %>%
  drop_na() %>%
ggplot() +
  aes(x = new_cases_per_million_kat, fill = gdp_per_capita_kat) +
  geom_bar()

barchart_plot_4 <-
COVID_data_latest %>%
  select(new_cases_per_million_kat, gdp_per_capita_kat) %>%
  drop_na() %>%
ggplot() +
  aes(x = new_cases_per_million_kat, fill = gdp_per_capita_kat) +
  geom_bar(position = "fill") +
  ylab("proportion")

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

```



*# a theme(legend.position) es a guides() funckiok
 # hasznalataval kontrollalhatjuk hogy hol es hogyan
 # jelenjen meg a jelmagyarazat az abran*

```

barchart_plot_3 <-
COVID_data_latest %>%

```

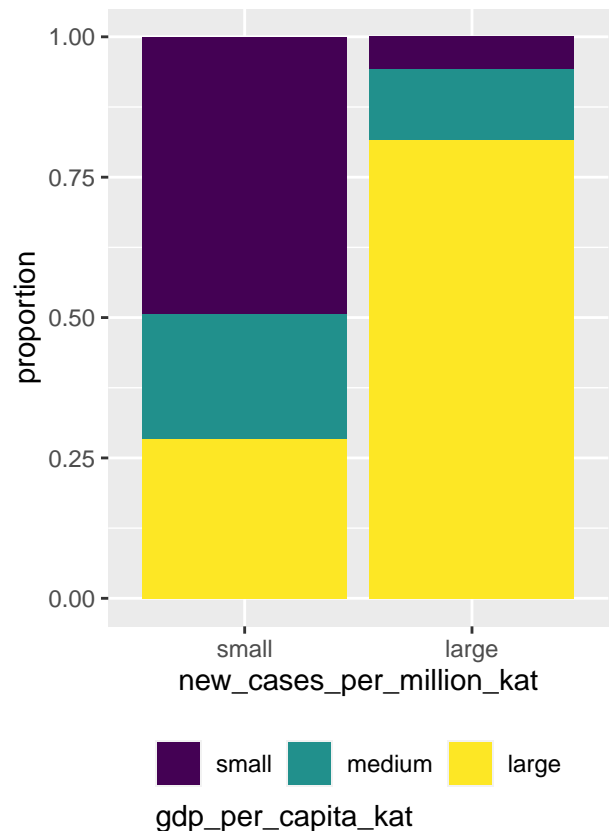
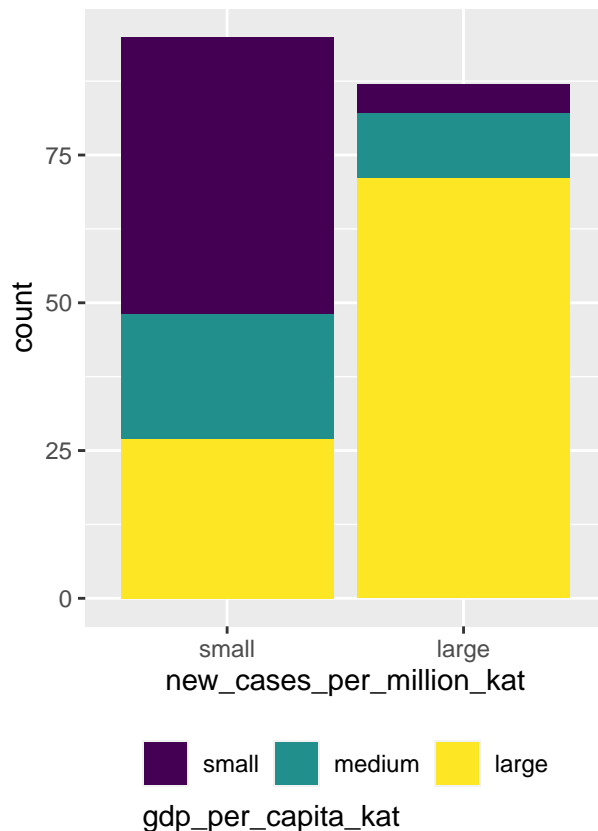
```

select(new_cases_per_million_kat, gdp_per_capita_kat) %>%
drop_na() %>%
ggplot() +
  aes(x = new_cases_per_million_kat, fill = gdp_per_capita_kat) +
  geom_bar() +
  theme(legend.position="bottom") +
  guides(fill = guide_legend(title.position = "bottom"))

barchart_plot_4 <-
COVID_data_latest %>%
  select(new_cases_per_million_kat, gdp_per_capita_kat) %>%
  drop_na() %>%
  ggplot() +
    aes(x = new_cases_per_million_kat, fill = gdp_per_capita_kat) +
    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 **egyedül** jelennek meg, vagy ha a második változó szerint **külön paneleken (facet)** jelennek meg.

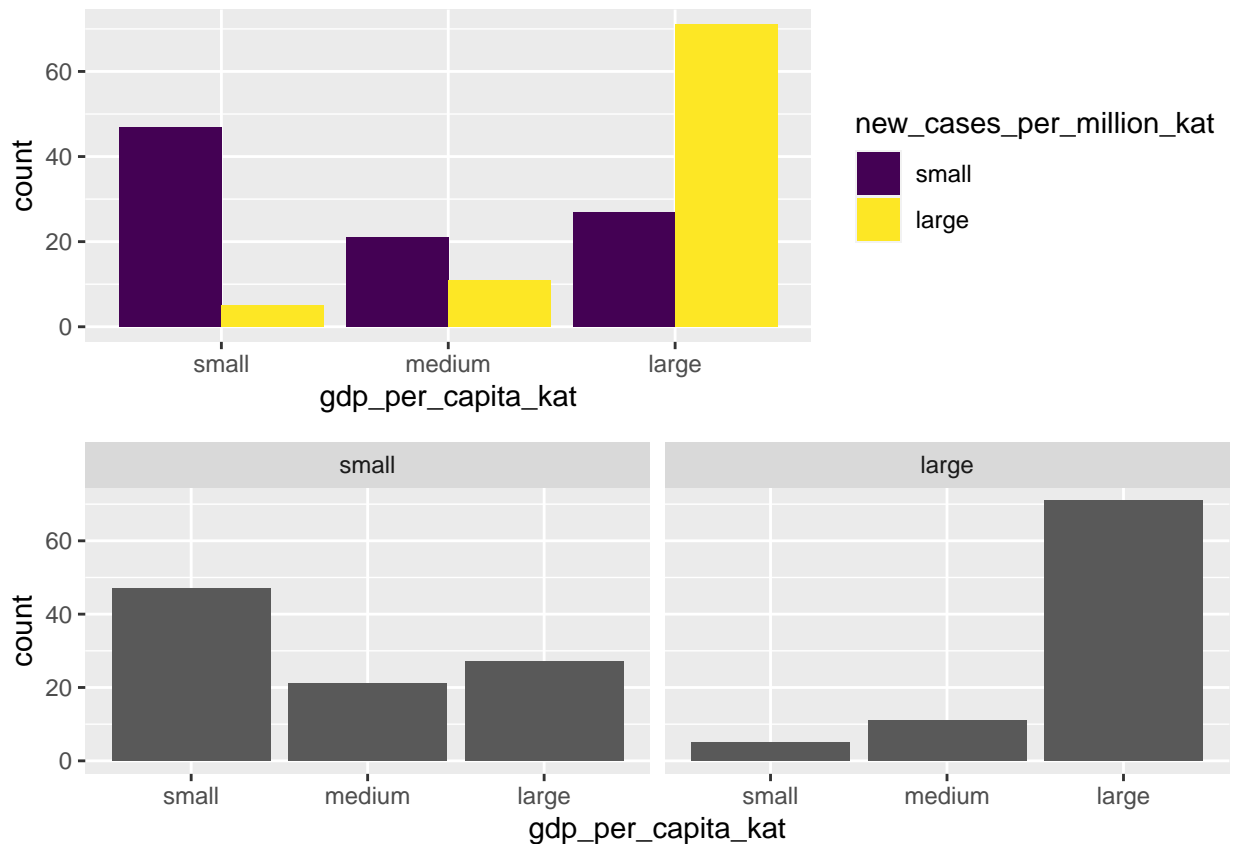

```

barchart_plot_5 <-
COVID_data_latest %>%
  select(new_cases_per_million_kat, gdp_per_capita_kat) %>%
  drop_na() %>%
ggplot() +
  aes(x = gdp_per_capita_kat, fill = new_cases_per_million_kat) +
  geom_bar(position = "dodge")

barchart_plot_6 <-
COVID_data_latest %>%
  select(new_cases_per_million_kat, gdp_per_capita_kat) %>%
  drop_na() %>%
ggplot() +
  aes(x = gdp_per_capita_kat) +
  geom_bar() +
  facet_wrap(~ new_cases_per_million_kat)

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

```



Egy kategorikus es egy numerikus valtozo kapcsolata

Vizsgáljuk meg hogy hogyan alakul az egy fore juto GDP kontinensenként. A GDP ebben az esetben egy folytonos változó (`gdp_per_capita`), es ennek az összefüggést szeretnénk megvizsgálni egy kategorikus változóval (`continent`).

Az explorációt kezdetben leíró statisztikák lekerdezésével csoportonként. Például ha arra vagyunk kíváncsiak,

milyen a GDP atlaga es szorasa kontinensenként, ezt megvizsgálhatjuk a `group_by()` es a `summarize()` segítségével.

```
COVID_data_latest %>%
  select(continent, gdp_per_capita) %>%
  drop_na() %>%
  group_by(continent) %>%
  summarize(mean = mean(gdp_per_capita),
            sd = sd(gdp_per_capita))
```

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

```
## # A tibble: 6 x 3
##   continent      mean      sd
##   <fct>         <dbl> <dbl>
## 1 Africa         5444.  6183.
## 2 Asia          22185. 25406.
## 3 Europe         33361. 18030.
## 4 North America 17126. 12871.
## 5 Oceania        12392. 16121.
## 6 South America 13841.  5110.
```

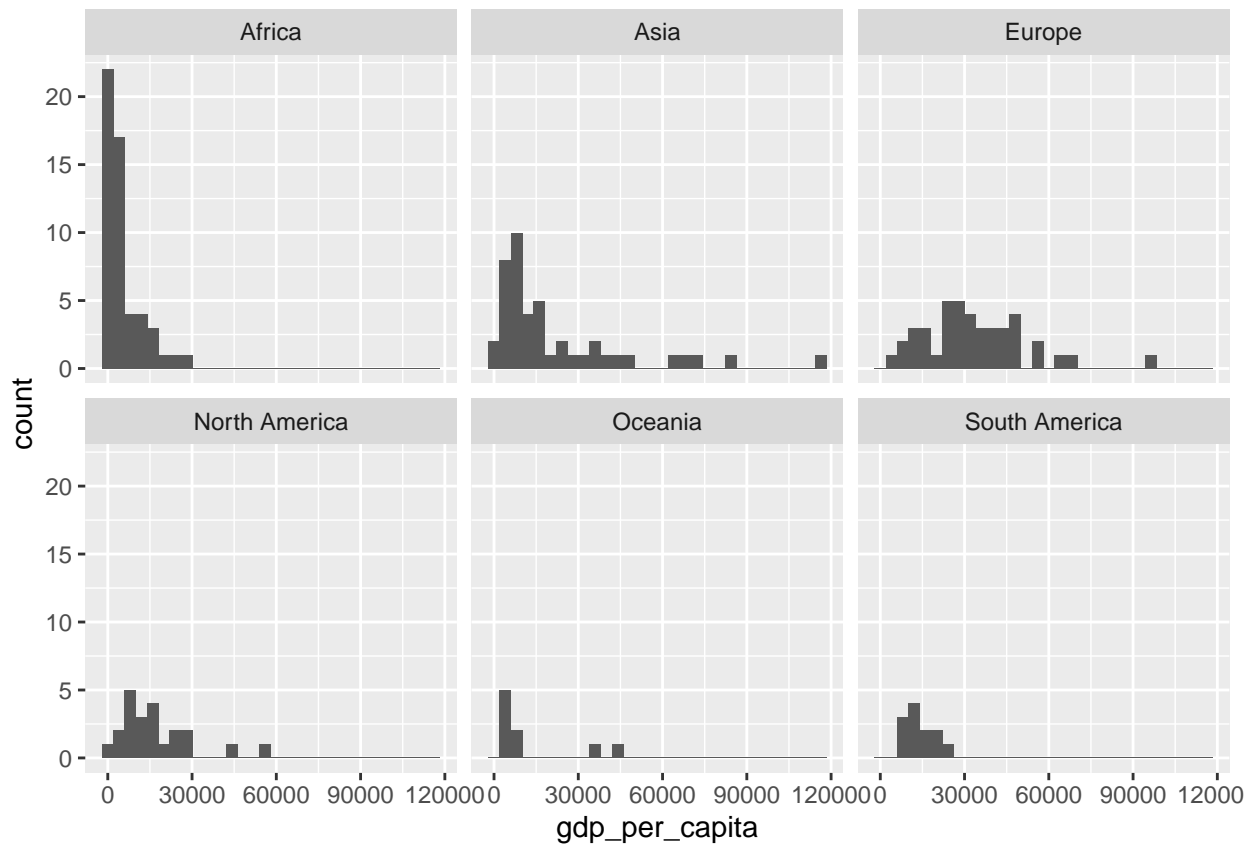
A ket változo kapcsolatát megvizsgálhatjuk **abrakkal** is. Pl. használhatjuk a

- `facet_wrap()` függvényt egy `geom_histogram()`-al kobinalva
- a `geom_boxplot()` -ot
- esetleg használhatunk egy egymásra illesztett `geom_density()` plot-ot ahol a kategoriak mas mas szinnel vannak jelolve.
- talan ebben az esetben a legtisztább kepet a `geom_violin()` mutatja, ami a `geom_boxplot()` es a `geom_density()` keverekenek tekintheto. Ezt kiegeszithetunk 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.

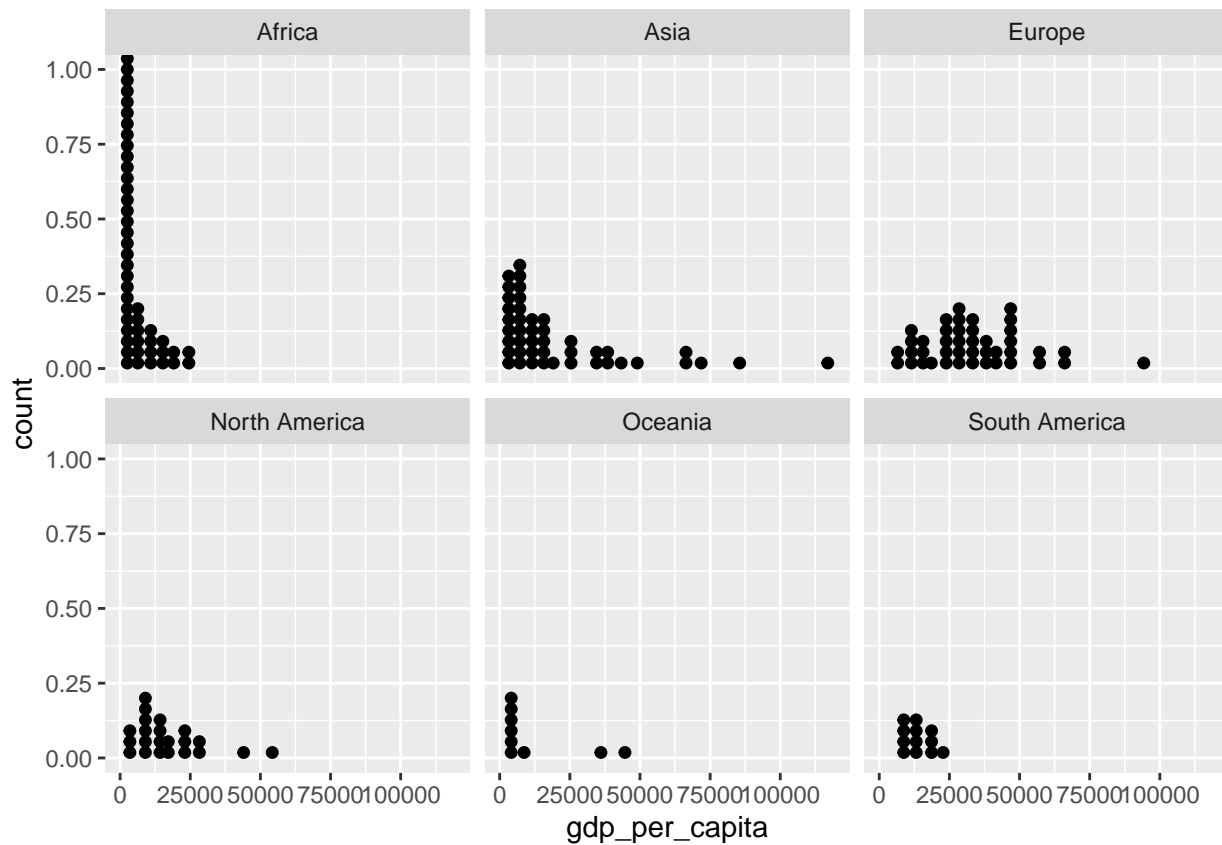
```
COVID_data_latest %>%
  select(continent, gdp_per_capita) %>%
  drop_na() %>%
  ggplot() +
    aes(x = gdp_per_capita) +
    geom_histogram() +
    facet_wrap(~ continent)
```

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

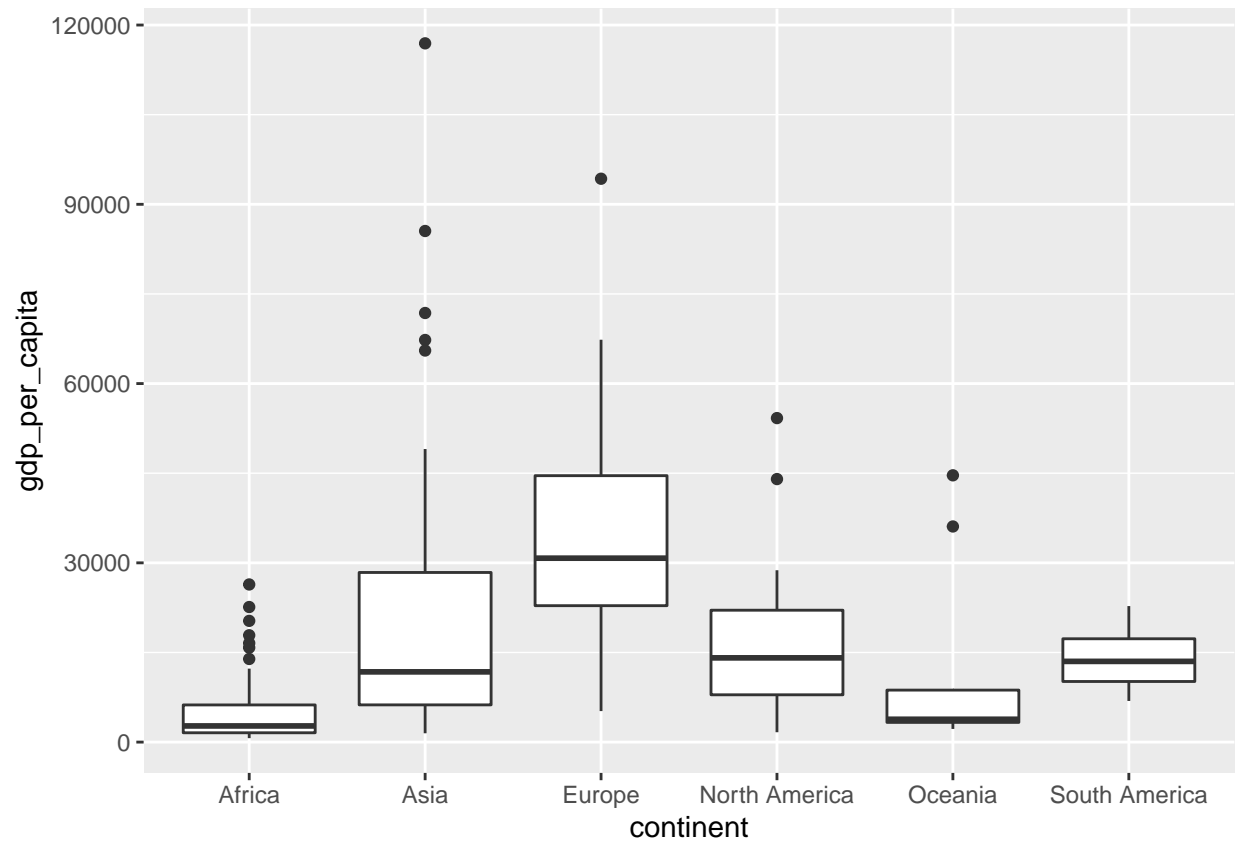


```
COVID_data_latest %>%
  select(continent, gdp_per_capita) %>%
  drop_na() %>%
  ggplot() +
    aes(x = gdp_per_capita) +
    geom_dotplot() +
    facet_wrap(~ continent)
```

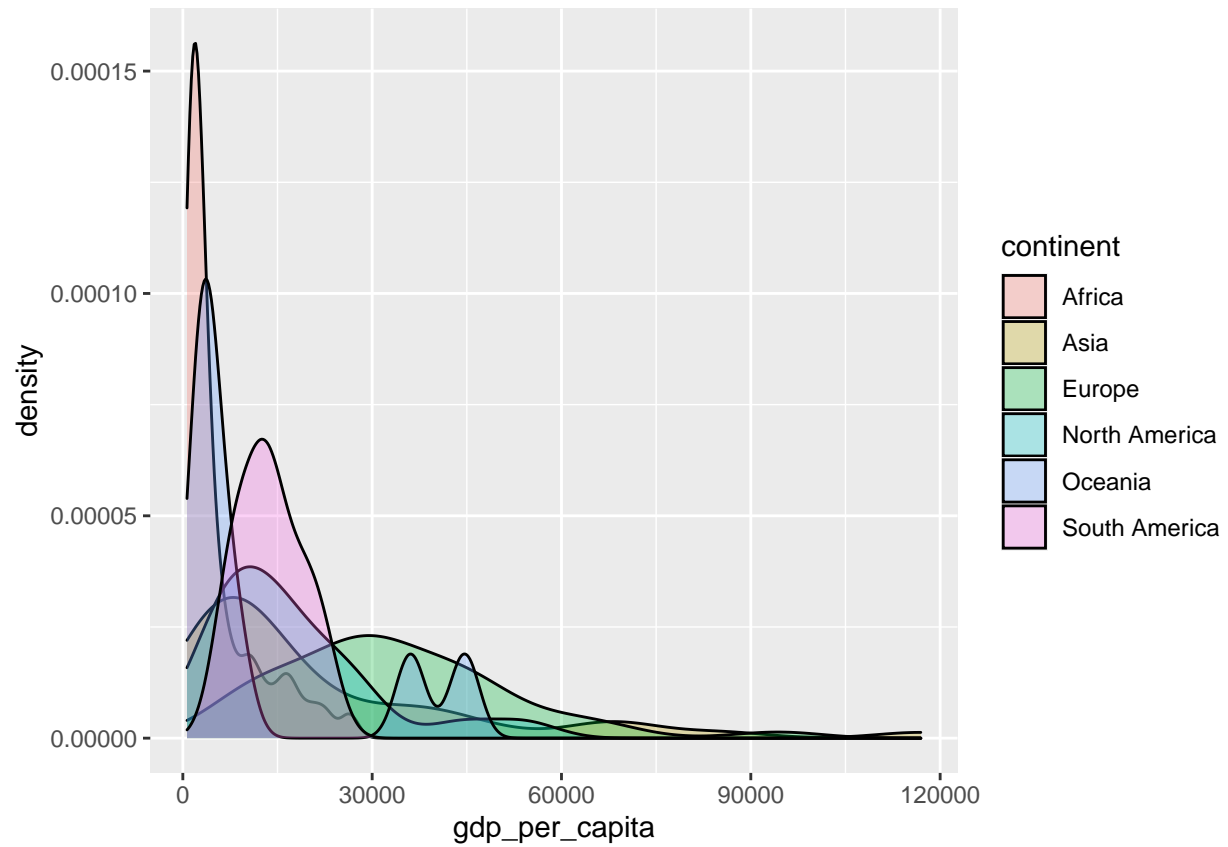
```
## `stat_bindot()` using `bins = 30`. Pick better value with `binwidth`.
```



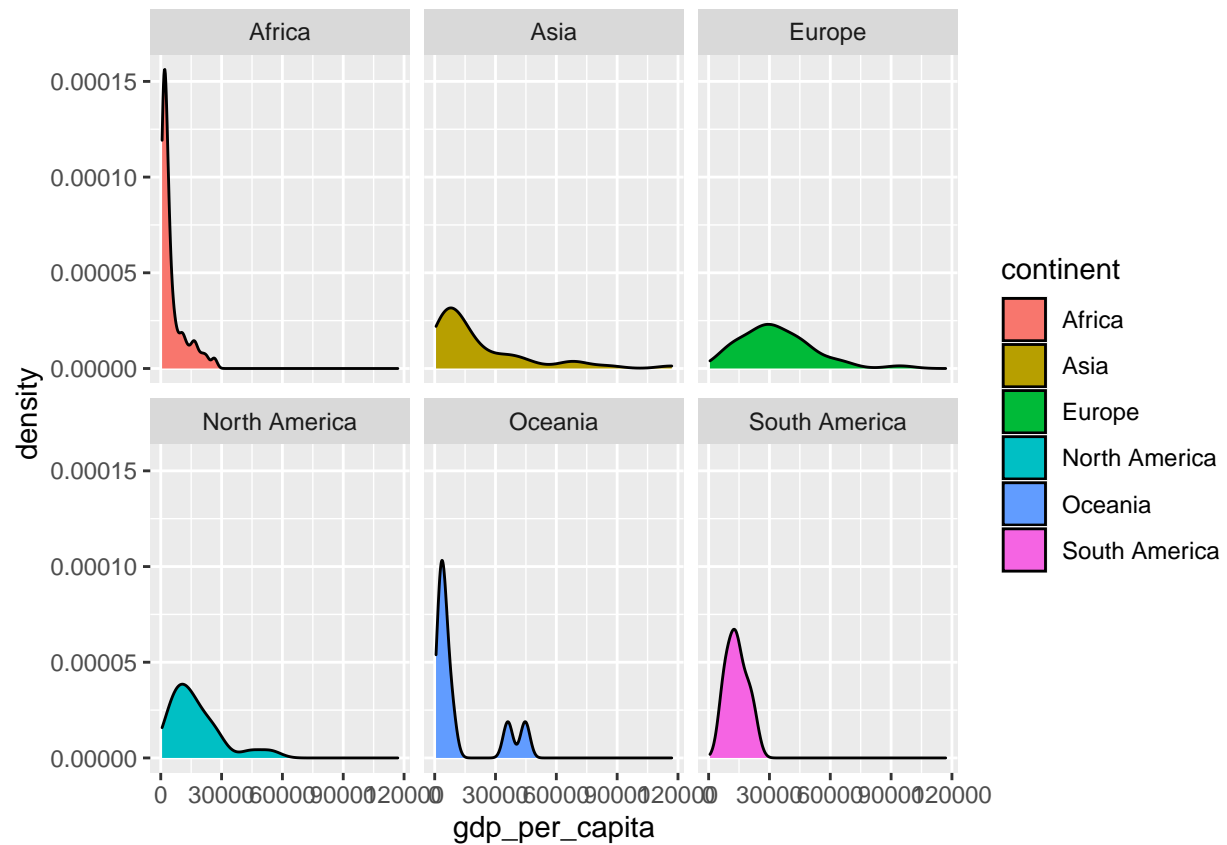
```
COVID_data_latest %>%
  select(continent, gdp_per_capita) %>%
  drop_na() %>%
  ggplot() +
    aes(x = continent, y = gdp_per_capita) +
    geom_boxplot()
```



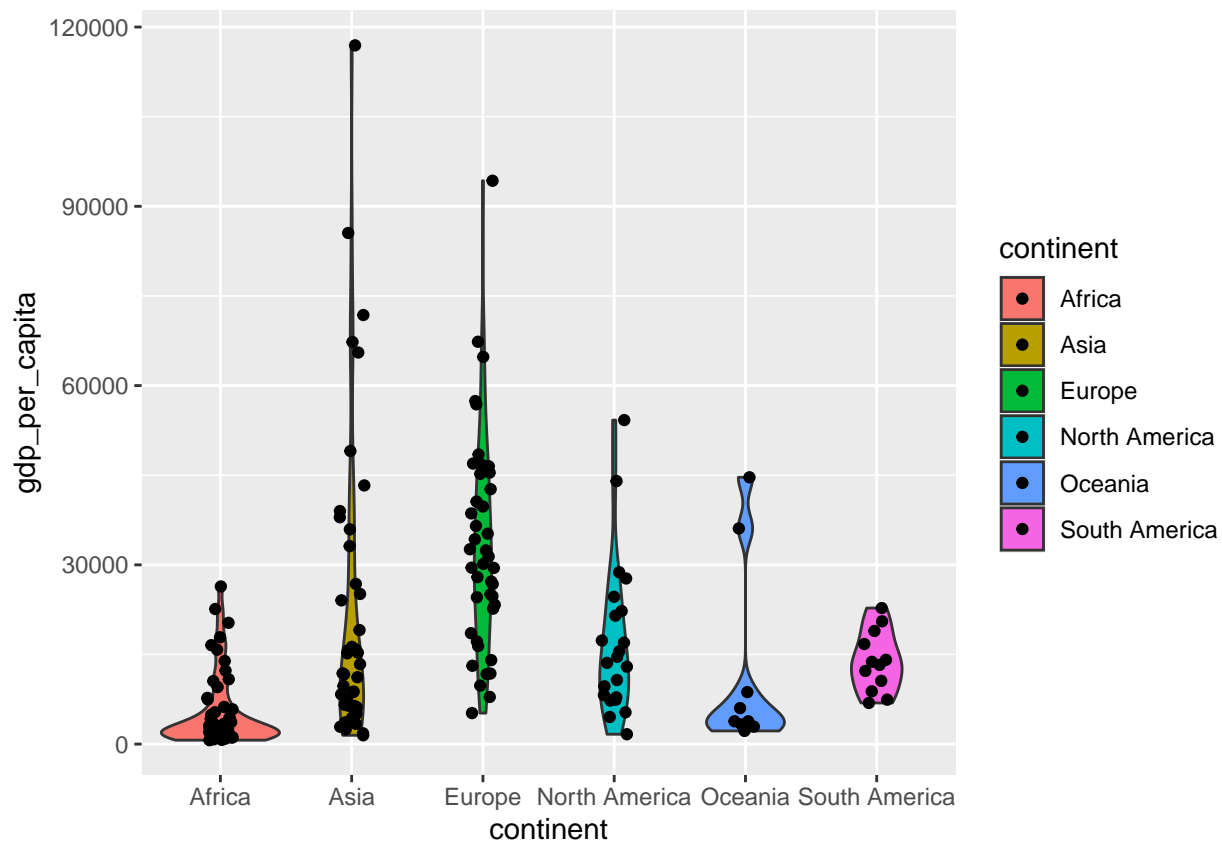
```
COVID_data_latest %>%  
  select(continent, gdp_per_capita) %>%  
  drop_na() %>%  
  ggplot() +  
    aes(x = gdp_per_capita, fill = continent) +  
    geom_density(alpha = 0.3)
```



```
COVID_data_latest %>%
  select(continent, gdp_per_capita) %>%
  drop_na() %>%
  ggplot() +
    aes(x = gdp_per_capita, fill = continent) +
    geom_density() +
    facet_wrap(~continent)
```



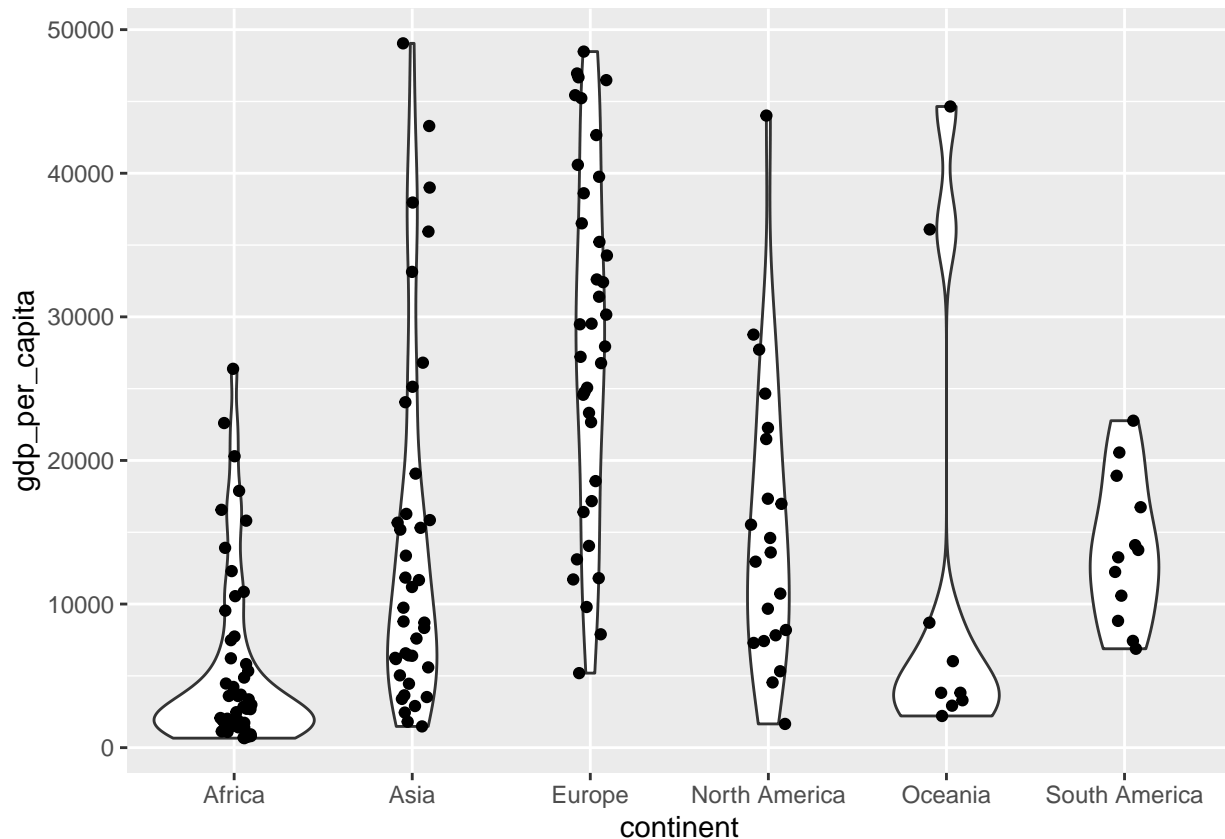
```
COVID_data_latest %>%
  select(continent, gdp_per_capita) %>%
  drop_na() %>%
  ggplot() +
    aes(x = continent, y = gdp_per_capita, fill = continent) +
    geom_violin() +
    geom_jitter(width = 0.1)
```



A fenti ábrán látszik, hogy Ázsiában a legtöbb országban viszonylag alacsony a GDP, viszont van néhány **kiurgo érték**, az átlagot felhúzza ebben a csoportban.

Ha szeretnénk **kizárni az elemzésünkben** az extrém értékeket, a **filter()** funkció bekezelevel a pipe-ba megépíthetjük a fenti ábrákat és táblázatokat úgy, hogy csak a 50,000-nél alacsonyabb GDP-jű országok kerüljenek az ábrára.

```
COVID_data_latest %>%
  select(continent, gdp_per_capita) %>%
  drop_na() %>%
  filter(gdp_per_capita < 50000) %>%
  ggplot() +
    aes(x = continent, y = gdp_per_capita) +
    geom_violin() +
    geom_jitter(width = 0.1)
```

```
COVID_data_latest %>%
  select(continent, gdp_per_capita) %>%
  drop_na() %>%
  filter(gdp_per_capita < 50000) %>%
  group_by(continent) %>%
  summarize(mean = mean(gdp_per_capita),
            sd = sd(gdp_per_capita))
```

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

```
## # A tibble: 6 x 3
##   continent      mean      sd
##   <fct>         <dbl> <dbl>
## 1 Africa         5444.  6183.
## 2 Asia          14591. 12710.
## 3 Europe         28661. 12390.
## 4 North America 15359. 10092.
## 5 Oceania        12392. 16121.
## 6 South America 13841.  5110.
```

Ha szeretnénk látni hogy a kisebb vagy nagyobb új esetszámmal jellemezhető országok (`new_cases_per_million_kat`) hogyan különböznek a GDP tekintetében kontinensenként akkor már **három változó** kapcsolatot kell ábrázolnunk. Ehhez a `facet_grid()` funkciót lehet használni, vagy különböző esztétikai elemeket (`aes()`) lehet a különböző változókhoz rendelni.

Gyakorlas

Hasznald a fent tanult modszereket, hogy megvizsgald a `total_cases_per_million` es a `gdp_per_capita_kat` változók közötti összefüggést.

- hasznald a fenti geomokat, es készíts legalább két különbozó 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 bármelyikeben hiányzó adat (NA) van. Ha ezt nem tesszük meg, a `cor()` függvény NA eredményt adhatna ha valamelyik változóban NA-val találkozunk.

```
COVID_data_latest %>%
  select(new_cases_per_million, gdp_per_capita) %>%
  drop_na() %>%
  cor()
```

```
##               new_cases_per_million gdp_per_capita
## new_cases_per_million             1.0000000      0.3296953
## gdp_per_capita                   0.3296953      1.0000000
```

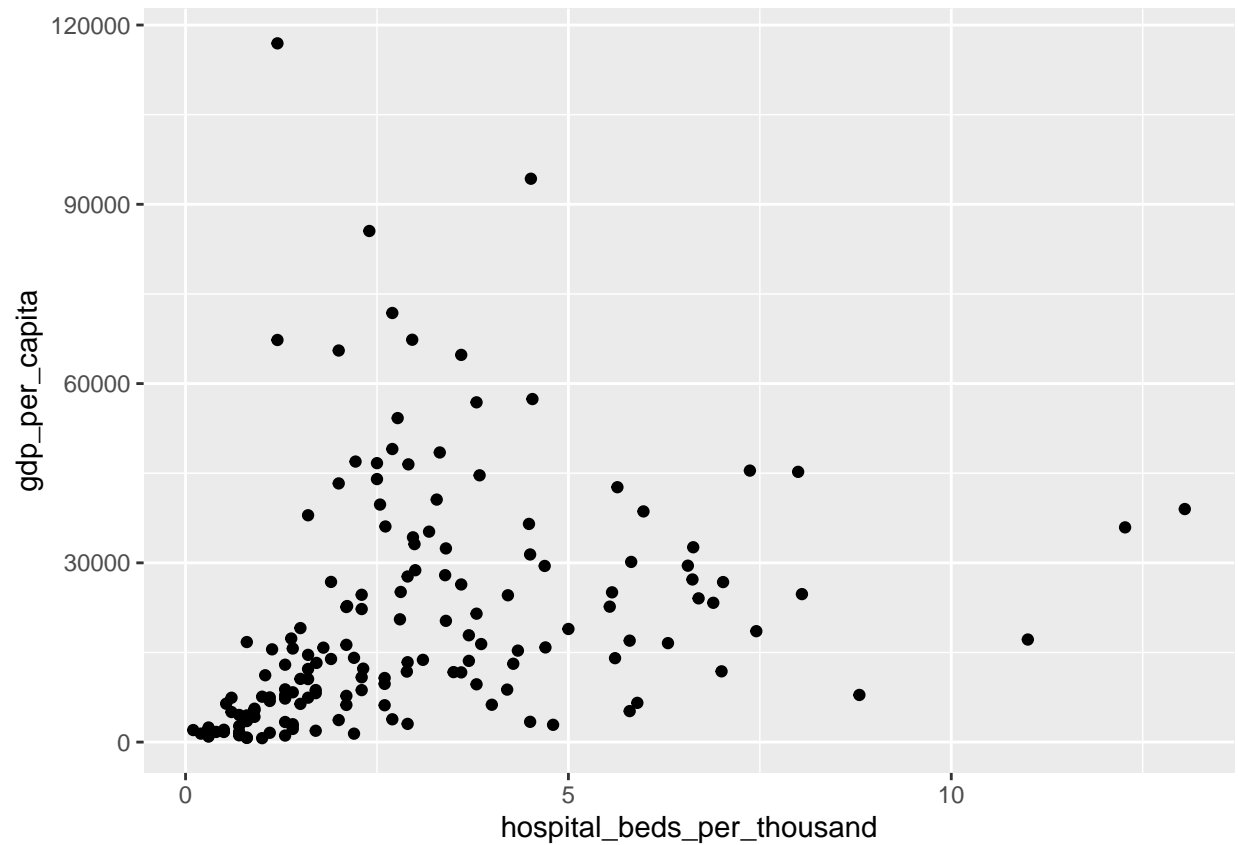
```
COVID_data_latest %>%
  select(new_cases_per_million, gdp_per_capita, hospital_beds_per_thousand) %>%
  drop_na() %>%
  cor()
```

```
##               new_cases_per_million gdp_per_capita
## new_cases_per_million             1.0000000      0.3082168
## gdp_per_capita                   0.3082168      1.0000000
## hospital_beds_per_thousand       0.2026815      0.2995055
##               hospital_beds_per_thousand
## new_cases_per_million             0.2026815
## gdp_per_capita                   0.2995055
## hospital_beds_per_thousand       1.0000000
```

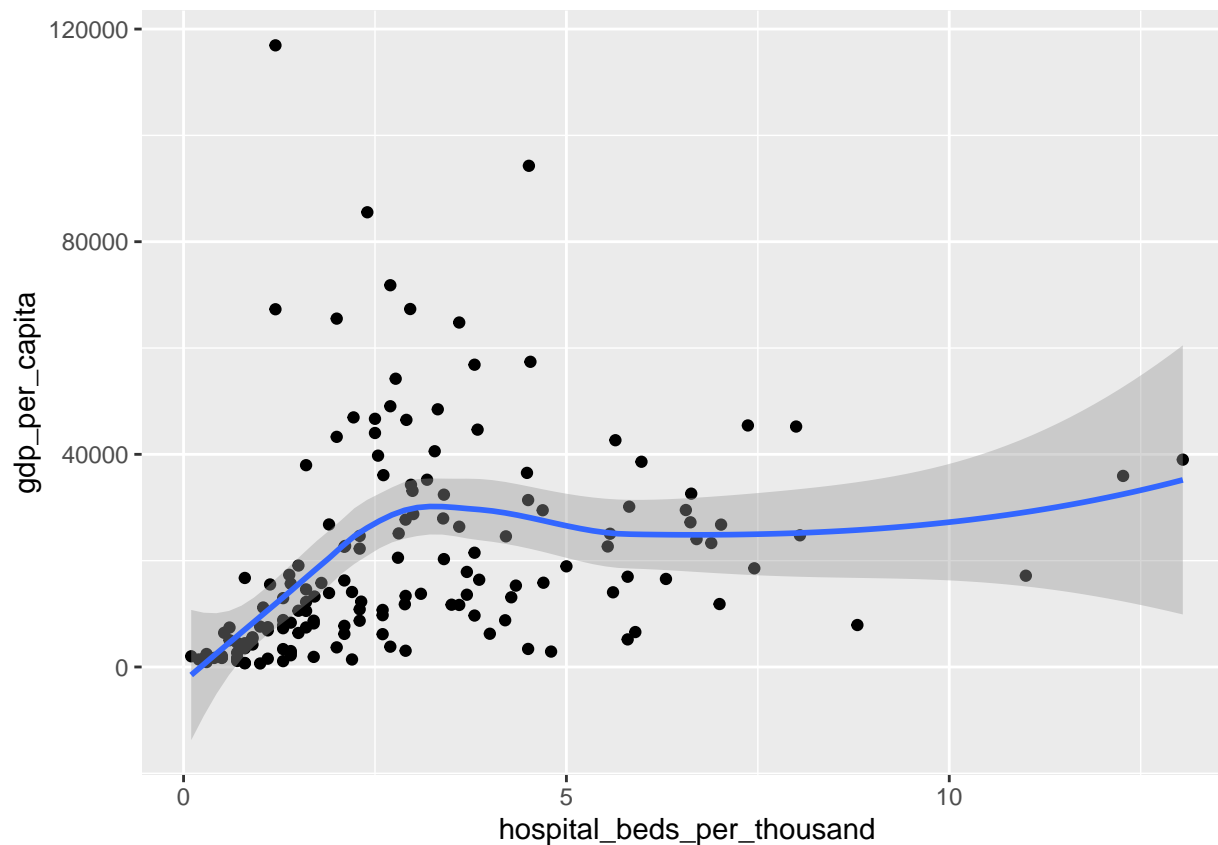
A numerikus változók közötti kapcsolatot általában pont diagrammal szoktuk ábrázolni (`geom_point()`)

A `geom_smooth()` layer hozzáadásával kaphatunk a pontok között meghúzódo trendrol egy képet. A kek vonal az ugyevezett trendvonal, a szurke sav a konfidencia intervallum. Ezekrol kesobb meg reszletesebben beszelunk majd

```
COVID_data_latest %>%
  select(hospital_beds_per_thousand, gdp_per_capita) %>%
  drop_na() %>%
  ggplot() +
  aes(x = hospital_beds_per_thousand, y = gdp_per_capita) +
  geom_point()
```



```
COVID_data_latest %>%  
  select(hospital_beds_per_thousand, gdp_per_capita) %>%  
  drop_na() %>%  
  ggplot() +  
    aes(x = hospital_beds_per_thousand, y = gdp_per_capita) +  
    geom_point() +  
    geom_smooth()
```



Gyakorlas

Milyen erős a kapcsolat a `aged_70_older` és a `gdp_per_capita` között?

- határozd meg a korrelációs együtthatót a változók között
- ábrázold a változók kapcsolatát

Több folytonos változó kapcsolata megjeleníthető például úgy, hogy az egyik változót egy színskálahoz rendeljük az alábbi módon.

```
COVID_data_latest %>%
  select(hospital_beds_per_thousand, gdp_per_capita, aged_70_older) %>%
  drop_na() %>%
  ggplot() +
    aes(x = hospital_beds_per_thousand, y = gdp_per_capita, col = aged_70_older) +
    geom_point() +
    scale_colour_gradientn(colours=c("green", "black"))
```

