

Statistik og dataanalyse

Tællende aktivitet 1

Peter Heilbo Ratgen - perat17@student.sdu.dk

11. oktober 2020

1 Databehandling

Datafilen indeholder data om landes CO2-udledning per indbygger i metriske ton per indbygger. Der er to kolonner i datafilen, "Country" og "CO2". Først gemmes .xlsx-filen som en .csv fil. Denne data er dog ikke komplet, det indeholder også de lande eller territorier der ikke har opgivet data. Det giver ikke mening at betragte disse lande i statistikken, derfor fjernes disse fra datasættet. Dette gøres ved at fjerne alle datapunkter der har "." som værdi. Også datapunkter der ikke er lande skal også fjernes fra datasættet, derfor fjernes alle datapunkter hvor "income" indgår i navnet med kommandolinjeværktøjet `sed`. Her er der fjernet 10 datapunkter. Derudover fjernes datapunktet "world" også, siden det noget der kan regnes frem til. Til sidst i databehandlingen skal , ændres til ., for at R importerer data som den korrekte type.

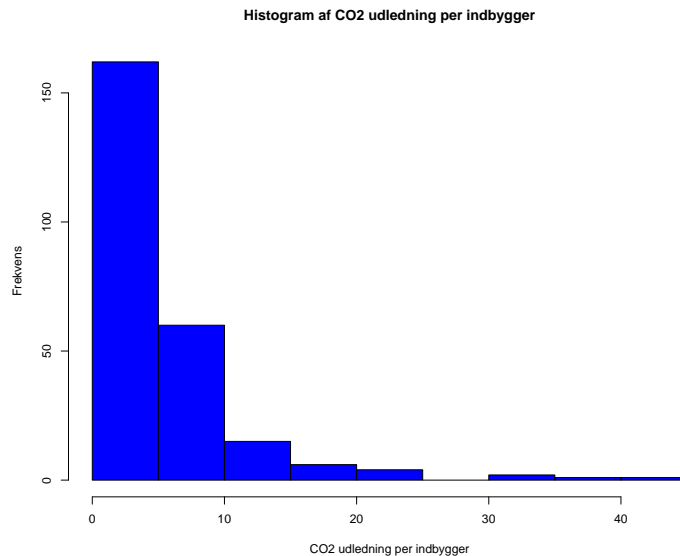
Nu kan data indporteres til R med:

```
1 co2 <- read.csv("CO2.csv", sep = ";", row.names = NULL)
```

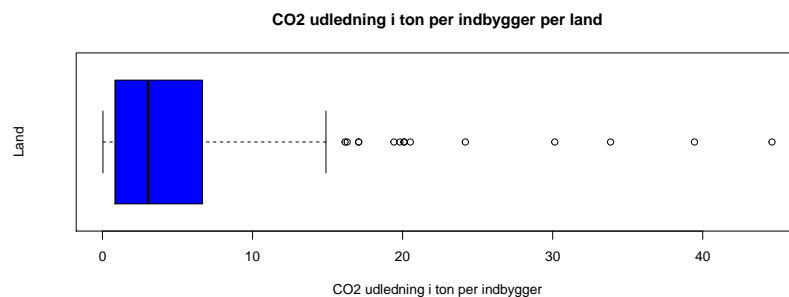
Her bruges `sep = ";"` til at indikere at værdierne i .csv formatet ikke er separeret med et komma, men med et semikolon.

2 Grafiske illustrationer

2.1 Histogram



2.2 Boxplot



Boxplottet fortæller at langt de fleste lande (75%, eller dem inden for det 3. kvartil) ikke udleder over 10 ton CO₂ per indbygger. Det fortæller også at der er nogle lande der udleder langt mere CO₂ per indbygger, disse er de lande der ligger uden for 1.5 gange den interkvartile afstand.

Danmark ligger indenfor det 3. kvartil med sine 6.51 ton CO₂ udledt per indbygger. Det vil også sige at Danmark ligger over medianen.

3 Deskriptiv analyse

Gennemsnit	4.9489
Median	3.0259
Typeinterval	0-5
Standardafvigelse	6.1881
Standard error	0.0246
Varians	38.2931
Minimum	0.0303
Maximum	44.6179
Q_1	0.8280
Q_3	6.6646

4 Normalfordeling

Som det ses på ovenstående histogram er data højreskæv. Normalfordelingen har en symmetrisk klokkeform. En anden indikator er at gennemsnittet og medianen ikke er tæt på at være den samme værdi.

Derfor er data ikke normalfordelt.

5 R-kode

```
1 co2 <- read.csv("C02.csv", sep = ";", row.names = NULL)
2
3 h <- hist(co2$C02)
4 pdf("co2plot.pdf", width = 10, height = 8)
5 plot(h, xlab = "C02 udledning per indbygger", ylab = "
  Frekvens"
6       , main = "Histogram af C02 udledning per indbygger",
7         col = "blue")
8 dev.off()
9 pdf("co2boxplot.pdf", width = 10, height = 4)
10 boxplot(co2$C02, main = "C02 udledning i ton per indbygger
  per land"
11          , col = "blue", horizontal = TRUE
12          , xlab = "C02 udledning i ton per indbygger"
13          , ylab = "Land", boxwex = 1.5)
14 dev.off()
15
16 #lande der ligger uden for 1.5 IRQ (potentielle outliers)
17 for (i in 1:length(co2$C02)) {
18   if (co2$C02[i] > 6.6646 + (6.6646 - 0.828) * 1.5) {
19     print(co2$Country[i])
20     print(co2$C02[i])
21   }
22 }
23
24 # Deskriptiv analyse
25 #standardafvigelse
26 sd(co2$C02)
27 summary(co2$C02)
28
29 #standard error
30 sd(co2$C02) / length(co2$C02)
31 var(co2$C02)
```