

# Datorlaboration 5

Johan Alenlöv, Josef Wilzén och Måns Magnusson

22 februari 2022

---

## Instruktioner

- Denna laboration ska göras i grupper om **två och två**. Det är viktigt för gruppindelningen att inte ändra grupper.
  - En av ska vara **navigatör** och den andra **programmerar**. Navigatörens ansvar är att ha ett helhetsperspektiv över koden. Byt position var 30:e minut. **Båda** ska vara engagerade i koden.
  - Det är tillåtet att diskutera med andra grupper, men att plagiera eller skriva kod åt varandra är **inte tillåtet**. Det är alltså **inte** tillåtet att titta på andra gruppers lösningar på inlämningsuppgifterna.
  - Använd gärna Teams för att ställa frågor. Det finns olika kanaler:
    - **Questions**: Skriv era frågor här. Svar kommer att ges öppet direkt i kanalen. Publicera inte kod till inlämningsuppgifter här (andra kan då se det). Det går bra att skriva frågor om inlämningsuppgifter här så länge ni inte inkluderar kod med lösningar till dessa uppgifter. Det går bra att publicera kod till övningsuppgifter här.
    - **Raise\_your\_hand**: Skriv här om ni vill ha hjälp men inte ställa er fråga öppet. Skriv något i stil med “Jag vill ha hjälp”. Då kommer en lärare att kontakta er när de har tid (i chatten på Teams). Vill flera ha hjälp så bildar de olika kommentarerna en kö, och hjälp kommer att ges i ordning efter kön. En “tumme upp” på kommentaren innebär att läraren har börjat hjälpa den aktuella studenten. Ett “hjärta” på kommentaren innebär att läraren har hjälpt klart studenten.
  - Använd inte å, ä eller ö i variabel- eller funktionsnamn.
  - Utgå från laborationsfilen, som går att ladda ned [här](#), när du gör inlämningsuppgifterna.
  - Spara denna som `labb[no]_grupp[no].R`, t.ex. `labb5_grupp01.R` om det är laboration 5 och ni tillhör grupp 1. Ta inte med hakparenteser eller stora bokstäver i filnamnet.  
**Obs!** Denna fil ska laddas upp på LISAM och ska **inte** innehålla något annat än de aktuella funktionerna, namn-, ID- och grupp-variabler och ev. kommentarer. Alltså **inga** andra variabler, funktionsanrop för att testa inlämningsuppgifterna eller anrop till markmyassignment-funktioner.
  - Om ni ska lämna i kompletteringar på del 2, döp då dessa till `labb5_grupp01._komp1.R` om det är första kompletteringstillfället. Se kurskanslenssidan för mer information om kompletteringar.
  - Laborationen består av två delar:
    - Övningsuppgifter
    - Inlämningsuppgifter
  - I laborationen finns det extrauppgifter markerade med \*. Dessa kan hoppas över.
  - Deadline för laboration framgår på [LISAM](#)
  - **Tips!** Använd “fusklapparna” som finns [här](#). Dessa kommer ni också få ha med på tentan.
-

# Innehåll

<b>I</b>	<b>Övningsuppgifter</b>	<b>4</b>
<b>1</b>	<b>Introduktion till grafik</b>	<b>5</b>
1.1	Visualisera en variabel . . . . .	5
1.1.1	Cirkeldiagram . . . . .	5
1.1.2	Barcharts . . . . .	6
1.1.3	Histogram . . . . .	8
1.2	Visualisering i flera variabler . . . . .	9
1.2.1	Sambandsdiagram . . . . .	9
1.2.2	Linjediagram . . . . .	10
1.2.3	Boxplot . . . . .	11
1.3	Grafiska inställningar och tillägg . . . . .	13
1.4	Spara figurer . . . . .	14
1.5	* Extraproblem: Skapa en egen graf . . . . .	14
<b>2</b>	<b>Slumptal och simulering</b>	<b>16</b>
2.1	Täthetsfunktioner mm . . . . .	16
2.2	<code>sample()</code> och <code>set.seed()</code> . . . . .	17
2.3	Exempel: <code>sum_of_dice()</code> . . . . .	18
2.4	Stora talens lag och Monte Carlo . . . . .	19
<b>3</b>	<b>Introduktion till R markdown och knitr</b>	<b>21</b>
3.1	Grunderna i markdown . . . . .	21
3.1.1	Grundläggande markdown . . . . .	22
3.1.2	Ekvationer . . . . .	23
3.2	Integrera R-kod med knitr . . . . .	23
3.2.1	Inline-kod . . . . .	23
3.2.2	R-block - “chunks” . . . . .	24
3.2.3	Grafik . . . . .	24
3.2.4	Tabeller . . . . .	25
<b>4</b>	<b>Input och output (I/O): Data på webben</b>	<b>26</b>
4.1	<code>downloader</code> . . . . .	26
4.2	Github . . . . .	27
4.3	Google docs . . . . .	27
<b>5</b>	<b>pxweb</b>	<b>29</b>
5.1	Navigera i SCB:s API . . . . .	29
5.1.1	Ladda ned data från SCB direkt med kod . . . . .	30
5.2	Namnstatistik 2021 . . . . .	31
5.3	Partistorlek i valet till kommunfullmäktige 2018 . . . . .	31
5.4	* Extraproblem: Partistorlek i valet till kommunfullmäktige 2014 . . . . .	31
5.5	* Extraproblem: Andra api:er . . . . .	31

<b>II</b>	<b>Inlämningsuppgifter</b>	<b>33</b>
<b>6</b>	<b>Inlämningsuppgifter</b>	<b>35</b>
6.1	estimate_pi() . . . . .	35
6.2	sum_of_random_dice() . . . . .	36
6.3	Miniprojektet del 1 . . . . .	39

# Del I

# Övningsuppgifter

## Kapitel 1

# Introduktion till grafik

I R finns en hel del funktionalitet för att arbeta med grafik. I det grundläggande R finns det som brukar kallas base graphs som är den grundläggande grafikfunktionaliteten. Utöver detta är paketet `ggplot2` mycket populär för visualisering. För en introduktion till grafisk visualisering av data är [1] ett standardverk.

I base-paketet fungerar grafiken på så sätt att vi lägger till lager för lager i en visualisering av data. Tänk dig att du ritar med en penna. Vi kommer i denna del använda oss av en del av de dataset som installeras tillsammans med R. Läs in dessa dataset på följande sätt. Undrar du vad materialen innehåller kan du använda `?iris`, `?mtcars`, `?Nile` för att få information om de olika variablerna.

```
data(iris)
data(mtcars)
data(Nile)
Nile <- as.data.frame(Nile)
Nile$years <- 1871:1970
```

Det finns många möjliga inställningar som går att göra. Dessa sammanfattas i dokumentationen för `par`. Sök efter `?par` i hjälpen för att få detaljerad information.

## 1.1 Visualisera en variabel

Vi inleder med att försöka visualisera data i en variabel.

### 1.1.1 Cirkeldiagram

Cirkeldiagram är ofta populärt, men bör generellt sett undvikas (se ex. [1]). Det är en typ av diagram som inte lämpas sig bra för människors tänkande och många har svårt att jämföra tårtbitar visuellt, [här](#) finns de vanligaste argumenten. För att skapa cirkeldiagram använder vi funktionen `pie()` på följande sätt.

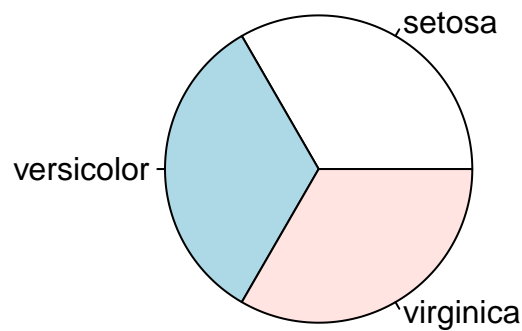
1. Som ett första steg måste vi beräkna frekvenser för den kategoriska variabel vi vill visualisera med `table()`. Testa `?table()`, `class(freqs)`, `str(freqs)`

```
freqs <- table(iris$Species)
freqs

      setosa versicolor  virginica 
       50         50         50
```

2. Baserat på denna frekvensfördelning är det därefter möjligt att skapa ett cirkeldiagram med `pie()`:

```
pie(freqs)
```



3. Vill vi ändra på labels gör vi det genom att ange en ny textvektor som labels, en rubrik anger vi med main:

```
pie(freqs, labels=c("Del 1", "Del 2", "Del 3"))  
pie(freqs, labels=c("Hej", "Hejsan", "Hej hej"), main="Cirkeldiagram")
```

4. På ett liknande sätt kan vi sedan ange vilka färger vi vill ange med argumentet col. Samtliga färger som går att använda i R finns [här](#).

```
pie(freqs, col=c("rosybrown1", "yellowgreen", "khaki2"))
```

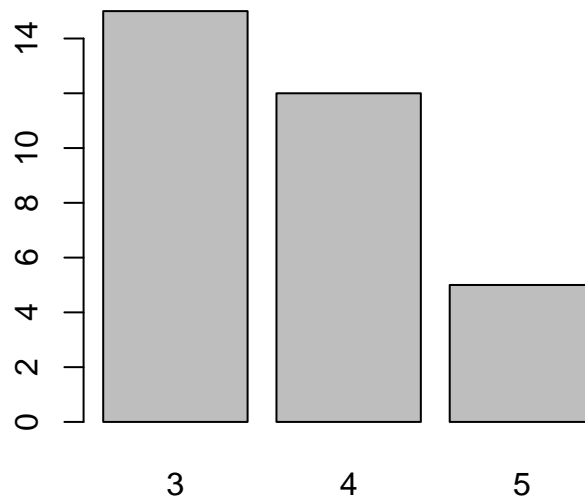
5. Vi kan självklart styra ännu mer i utformningen av cirkeldiagrammen. För mer hjälp för cirkeldiagram använd `?pie`. **Men som sagt undvik cirkeldiagram i största möjliga mån.**

### 1.1.2 Barcharts

Barcharts är enklare att tolka på ett korrekt sätt och är en av de mest grundläggande graftypeperna.

1. För att skapa en barchart behöver vi (på samma sätt som för cirkeldiagrammen) utgå från frekvenser när vi skapar vårt diagram. Vi börjar med det allra enklaste diagrammet:

```
freqCars <- table(mtcars$gear)  
barplot(freqCars)
```



2. Som framgår ovan får vi ett mycket enkelt diagram. Diagrammet använder sig av radnamn (`rownames()`) för `freqCars`. Prova att kolla hur `freqCars` ser ut och dess radnamn.
3. Att lägga till rubriker och titel gör vi på samma sätt som för cirkeldiagrammen ovan, med argumenten `main`, `xlab` och `ylab`.

```
barplot(freqCars, main="Cars", xlab="Gears", ylab="Counts")
```

4. Vi kan också välja att ha horisontella staplar med argumentet `horiz=TRUE`.

```
barplot(freqCars, main="Cars", horiz=TRUE)
```

5. Som för cirkeldiagrammet kan vi också byta färg om vi vill med `col`.

```
barplot(freqCars, main="Cars", col="red")
```

6. Vi kan också ha flera variabler i samma stapeldiagram, så kallade grupperade eller "stackade" stapeldiagram. Vi börjar med grupperade stapeldiagram.

```
carTable <- table(mtcars$vs, mtcars$gear)
barplot(carTable, main="Car Distribution by Gears and VS",
xlab="Number of Gears",
col=c("darkblue", "red"), legend= rownames(carTable))
```

7. På ett liknande sätt kan vi skapa "stackade" stapeldiagram:



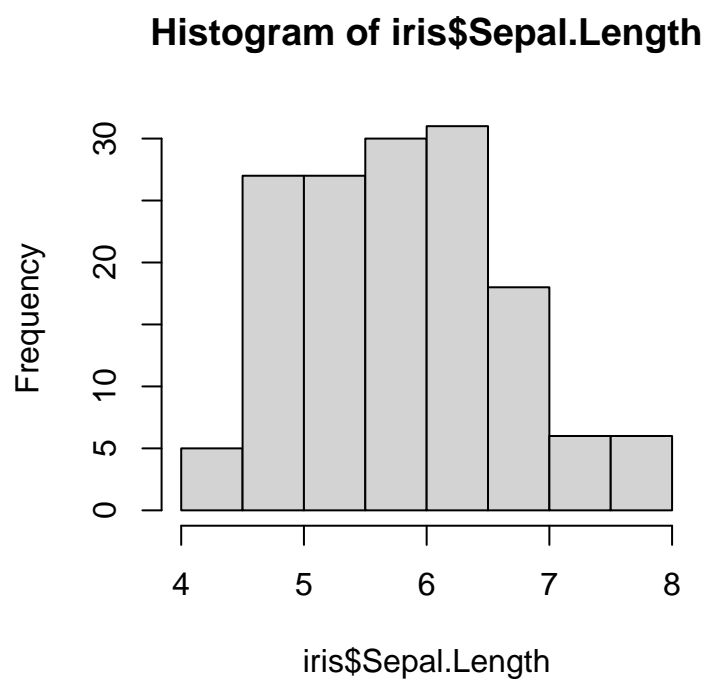
```
barplot(carTable, main="Car Distribution by Gears and VS",
xlab="Number of Gears", col=c("darkblue","red"), beside=TRUE)
```

### 1.1.3 Histogram

Histogram kan vi använda för att visualisera en kontinuerlig variabel.

1. För att skapa ett enkelt histogram använder vi funktionen `hist()`.

```
hist(iris$Sepal.Length)
```



2. Att ändra rubriker och färger görs på samma sätt som i övriga diagram.

```
hist(iris$Sepal.Length, col="blue", main="Min titel", xlab="X-titel", ylab="Y-titel")
```

3. När det gäller histogram kan det vara så att i vissa fall vill vi ha fler eller färre staplar. Detta styrs med `breaks`. Prova koden nedan, testa att ändra värde på `breaks` till 10, 20, 40 och 70.

```
hist(iris$Sepal.Length, breaks=40, col="red")
```

4. Histogram är diskreta, men vi kan enkelt i R skapa en uppskattning av den underliggande fördelningen visuellt.

```
dens <- density(iris$Sepal.Length)
plot(dens)
```

5. Vi kan styra om vi vill ha en absolut eller en relativ skala på histogrammet. Testa koden nedan. Hur skiljer sig y-axeln i de båda graferna? Vilket är defaultvärde för argumentet `freq`? (tips: kolla i dokumentationen)

```
hist(iris$Sepal.Length, breaks=40, col="red",freq = TRUE)
hist(iris$Sepal.Length, breaks=40, col="red",freq = FALSE)
```

6. Det går att kombinera ett histogram med en täthetskurva. Testa koden nedan, vad händer?

```
hist(iris$Sepal.Length, breaks=15, col="red",freq = FALSE)
lines(dens)
```

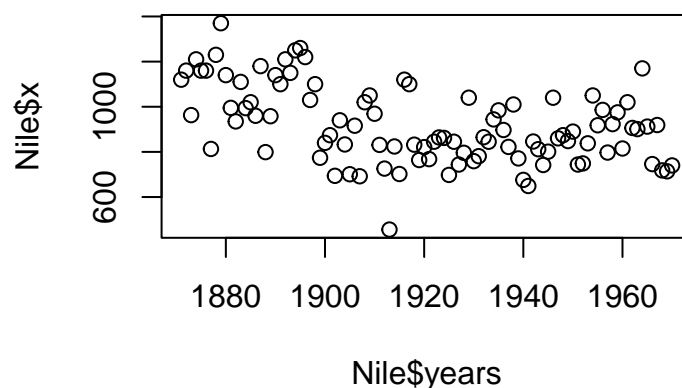
## 1.2 Visualisering i flera variabler

### 1.2.1 Sambandsdiagram

En av de vanligaste sätten att visualisera tvådimensionella data är med scatter plots (punktdiagram). Vi ska nu pröva att visualisera den historiska utvecklingen av Nilens vattennivåer.

1. Att skapa en vanlig scatterplots görs med funktionen `plot()`. Det är en generisk funktion<sup>1</sup> och i många fall använder vi `plot` för att visualisera olika typer av objekt. Testa `methods(plot)` för att se vilka metoder som finns för `plot()`. Grundutförandet ger dock en scatterplot på följande sätt:

```
plot(x = Nile$years, y = Nile$x)
```



2. Som tidigare kan vi också lägga till/förändra rubriker enkelt om vi vill.

```
plot(Nile$years, Nile$x, main="Water in the Nile", xlab="Years", ylab="Level")
```

3. Vill vi ändra färgen på våra punkter använder vi som vanligt parametern `col`.

---

<sup>1</sup>Mer om detta kursvecka 6.

```
plot(Nile$years, Nile$x, col="blue")
```

4. Vill vi att olika punkter ska ha olika färger anger vi bara en vektor med färgnamn.

```
colVector<- rep("blue",length(Nile$x))
colVector[Nile$years>1900] <- "red"
plot(Nile$years, Nile$x, col=colVector)
```

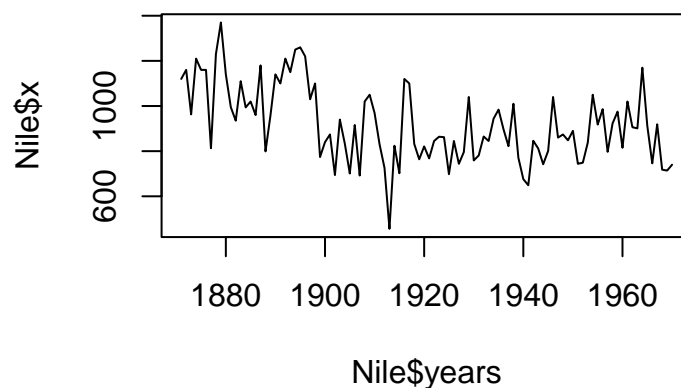
5. Studera hur `colVector` ser ut ovan så du förstår hur vektorn används för att styra färgen på punkterna. Prova att använd ytterligare än färg till punkterna efter 1945.
6. Vi kan också ändra hur punkterna ser ut och använda andra symboler. Det finns totalt 25 olika symboler i baspaketet. För att använda en punkttyp används argumentet `pch`. Med koden nedan kan vi snabbt se alla olika typer av punkter som finns i baspaketet. Se även dokumentation för funktionen `points()`

```
plot(1:25,rep(1,25), type="p", pch = 1:25,cex = 2.5)
# cex ändrar storleken
```

## 1.2.2 Linjediagram

1. Nilens vattennivåer är en tidsserie snarare än ett "vanligt" samband. För tidsserier vill vi ofta ha linjegrafer istället för enskilda punkter. För att ändra till en linjefgraf anger vi bara `type="l"`.

```
plot(Nile$years, Nile$x, type="l")
```



2. Precis som när det gäller punkter kan vi använda olika linjetyper med argumentet `lty`. Prova linjetyp 2 t.o.m. 10. Tips: använd en for-loop!

```
plot(Nile$years, Nile$x, type="o", lty=2)
```

3. Vill vi både ha linjer och punkter kan vi använda `type="o"`. Läs i dokumentationen för `plot()` under argumentet `type=`.

```
plot(Nile$years, Nile$x, type="o", lty=2, pch=3)
```

4. En sista typ av linjegrav är en trappstegsgraf:

```
plot(Nile$years, Nile$x, type="s")
```

5. Prova att lägg till argumentet `lwd=3` i plotten ovan. Vad innebär detta? Testa att köra `?par` Läs under rubrikerna `lwd` och `lty`.
6. När vi ska göra grafter över tidserier så är det ofta smidigt att låta värdena på x-axeln vara ett datumobjekt (mer om det kursvecka 6). Då kan vi automatiskt få en bra tidsskala på x-axeln när vi använder `plot()`. Nedan följer några exempel på hur ni kan göra:

```
data(Nile)
Nile <- as.data.frame(Nile)
Nile$years <- 1871:1970
temp<-paste0(Nile$years,"-01-01")
date_var<-as.Date(temp)
str(date_var)
plot(x = date_var,y = Nile$x,t="l",xlab="date")

# plotta två månader med data:
temp2<-paste0("2021-01-",1:31)
temp3<-paste0("2021-02-",1:28)
date_var2<-as.Date(c(temp2,temp3))
print(date_var2)

# skapa lite data:
x<-1:59
y<-10*log(x)-0.01*x^2+0.0001*x^3
plot(x = date_var2,y = y,t="l",xlab="date")

# mer exempel
temp4<-paste0(rep(2000,each=12),"-",1:12,"-01")
date_var3<-as.Date(temp4)
date_var3 y2<-rep(1:6,2)
plot(x = date_var3,y = y2,t="l",xlab="date")

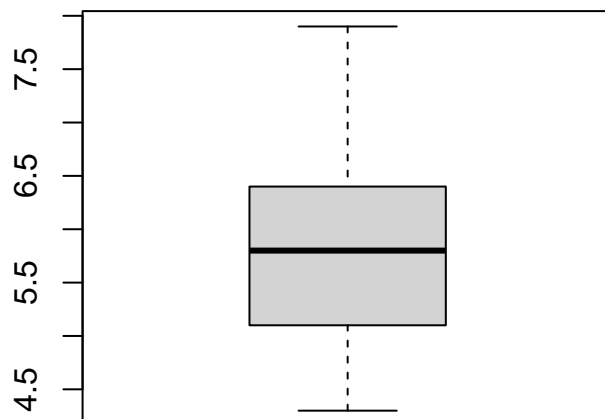
temp5<-paste0(rep(2000:2005,each=12),"-",1:12,"-01")
date_var4<-as.Date(temp5)
y3<-rep(1:24,3)
plot(x = date_var4,y = y3,t="l",xlab="date")
```

### 1.2.3 Boxplot

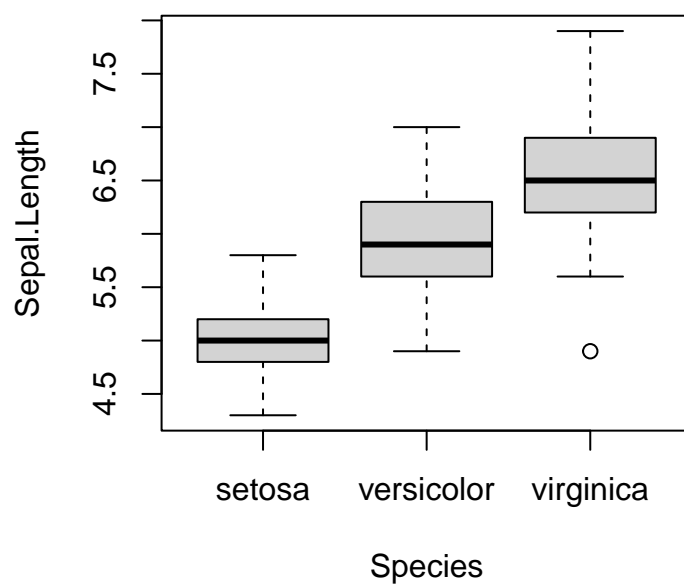
Vill vi jämföra olika fördelningar efter en kategorisk variabel gör vi det med fördel med en boxplot

1. Nedan finns kod för att producera en boxplot.

```
boxplot(iris$Sepal.Length)
```



```
boxplot(Sepal.Length~Species, data=iris)
```



2. `Sepal.Length~Species` är ett exempel på formel-objekt. Formel-objekt används i olika sammanhang i R, men ett vanligt exempel är när linjär regressions ska skattas med funktionen `lm()`. Testa att köra `?formula`. Kör sedan koden nedan:

```
Sepal.Length~Species
y<-Sepal.Length~Species
class(y)
str(y)
```

3. Precis som i tidigare diagram är det enkelt att lägga till färger.

```
boxplot(Sepal.Length~Species, data=iris, col=c("blue", "green", "red"))
```

4. Eller att lägga till rubriker.

```
boxplot(Sepal.Length~Species, data=iris, main="Blommor!")
```

## 1.3 Grafiska inställningar och tillägg

Ovan har vi sett en hel del av de figurer som går att producera. Nedan kommer lite mer inställningar och tillägg vi kan göra när vi arbetar med grafik i R.

1. För att kontrollera vilka värden som ska vara med på x- och y-axeln i en plot används argumenten `xlim=` och `ylim=`. Kör koden nedan. Ändra värdena på `xlim=` och `ylim=` och se vad som händer med plotten.

```
plot(Nile$years, Nile$x, type="s")
plot(Nile$years, Nile$x, type="s", xlim=c(1900,1945), ylim=c(600, 1200))
```

2. Vill vi lägga på punkter i grafen använder vi `points()`. Prova att lägga till punkterna i plotten:

```
points(Nile$years, Nile$x, pch=20)
```

3. Vi kan också lägga till godtyckliga linjer (t.ex. som referenser) med `abline()`.

```
abline(v=1920, lty=4)
abline(h=900, lty=10)
```

4. Vad händer om du ändrar värdet i "v=" och "h=" till numeriska vektorer?
5. `abline()` kan också användas för att rita ut räta linjer med räta linjens ekvation (om vi ex. anpassat en regressionsmodell):

$$f(x) = a + bx$$

```
abline(a=6130, b=-2.7)
```

6. Vill vi ha flera grafer i en använder vi `par(mfrow=c(3,2))`. Men det vi säger till R med detta kommando är att vi vill ha tre rader och två kolumner med figurer. Dessa struktur kommer att fyllas radvis. Prova att köra denna kod och skapa därefter 6 figurer (vilka som helst). Ändra till en 2×2 figur och skapa på samma sätt fyra figurer. `par(mfrow=c(1,1))` återställer till en plot-struktur med en figur. Argumentet `mfc01=c(3,2)` fyller kolumnvis.

## 1.4 Spara figurer

I många fall vill spara specifika grafer i olika format. I R finns ett antal olika format som kan användas. De vanligaste är TIFF, BNP, JPEG, PNG och PDF. I alla fall används den aktuella funktionen för formatet `tiff()` för TIFF, `pdf()` för PDF o.s.v.

1. De olika grafikfunktionerna använder olika argument, men gemensamt är att vi först anger vilket format vi vill använda, sedan skapar vi vår figur och därefter stänger vi av "utskriften" till denna fil. Kör koden nedan för ett exempel:

```
jpeg(filename = "minJPEG.jpeg", width = 480, height = 480)
plot(Nile$years, Nile$x, type="l")
dev.off()
```

2. Figurer som skrivs ut på detta sätt hamnar i "working directory". Leta reda på mappen som är ditt working directory och öppna jpeg-filen i ett bildvisningsprogram.
3. Upprepa nu uppgift 1, men spara filen i en annan mapp än i ditt working directory. Detta görs genom att ändra argumentet `filename` till en sökväg (path) som slutar på det faktiska filnamnet (i det här fallet "minJPEG.jpeg"). Sökvägen anges på samma sätt som när filer ska läsas från eller sparas på hårddisken.
4. För pdf:er finns det ett snabbare sätt att snabbt skriva ut en figur som pdf:

```
plot(Nile$years, Nile$x, type="l")
dev.copy2pdf(file="MinNilenPlot.pdf")
```

5. Prova nu att själv skriva ut en av dina figurer ovan i PNG- och TIFF-format. Se till att spara dem i en annan mapp än ditt working directory.
6. I Rstudio går det även att spara figurer genom att klicka på panelen "Plots", och sen klicka på "Export" → "Save as image". Då kan ni välja vilken mapp ni ska spara i storlek på figuren etc. Testa denna metod genom att spara en valfri plot någonstans på din dator.

## 1.5 \* Extraproblem: Skapa en egen graf

1. Ladda in datasetet geyser med `data(faithful)`.
2. Testa att göra två histogram, dels över `waiting` och `eruptions`.
  - (a) Ändra antalet `breaks` i histogrammen. Hur ser det ut om du har väldigt många eller väldigt få?
  - (b) Prova att ändra färg med `col="blue"`. Testa att ändra färgen till `col=c(1,2,3)`
  - (c) Ändra nu rubriken och axeltexterna med `xlab=` och `ylab=`.
  - (d) Spara ned de två histogrammen i en figur (förslagsvis över och under varandra) som en pdf.
3. Gör en scatterplot av `waiting` mot `eruptions`. Huvudrubrik ska vara "Old faithful".
4. Det verkar finnas två tydliga kluster.
  - (a) Ge de olika klustren olika färger.
  - (b) Ge de olika klustren olika punktsymboler.
5. Lägg till en bildtext (eng: legend) till figuren (kolla på `?legend`) enligt nedan. Testa att lägga "legend" på någon annan plats i plotten.

```
legend("topleft", pch = c(1, 2), col = c("red", "blue"), legend = c("clu1", "clu2"))
```

6. Spara ned även denna graf i pdf-format.



## Kapitel 2

# Slumptal och simulering

### 2.1 Täthetsfunktioner mm

En central del inom statistik och analys handlar om simulering och att hantera olika sannolikhetsfördelningar<sup>1</sup>. Det finns ett stort antal sannolikhetsfördelningar vi kan simulera ifrån/beräkna. Dessa har en täthetsfunktion (pdf), en kumulativ fördelningsfunktion (cdf) och en invers kumulativ fördelningsfunktion (q som i quantile). De flesta fördelningar har fyra varianter:

Prefix	Beskrivning	Exempel
r	simulera från fördelningen	<code>rnorm()</code>
d	täthetsfunktionen (pdf)	<code>dnorm()</code>
p	kumulativ fördelningsfunktion (cdf)	<code>pnorm()</code>
q	inversa kumulativa fördelningsfunktionen	<code>qnorm()</code>

För att se vilka fördelningar som finns förinstallerade med R se `?distribution`.

Vill vi exempelvis simulera 100 tal från  $\mathcal{N}(10, 1^2)$  gör vi på följande sätt:

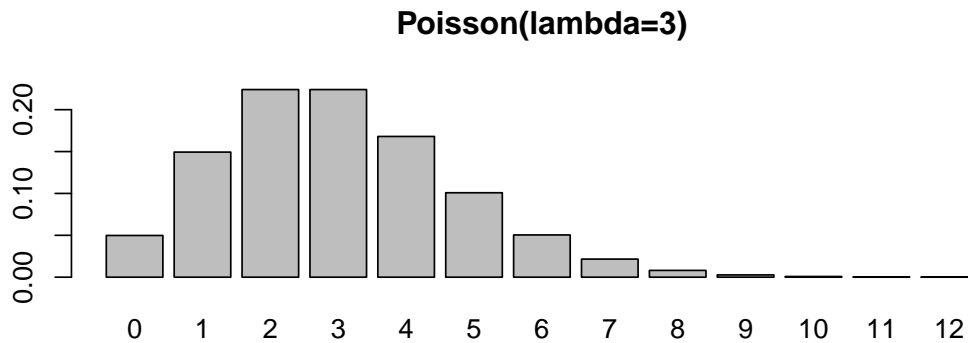
```
minNormal <- rnorm(n=100, mean=10, sd=1)
```

1. Skapa en vektor med 30 slumptal från den uniforma fördelningen  $\mathcal{U}(\min = 1, \max = 10)$ . [**Tips!** `runif()`]
2. Skapa en vektor med 200 slumptal från normalfördelningen med medelvärde 1 och varians 2, döp den till `minNorm`.
3. Skapa ett histogram över `minNorm` [**Tips!** `hist()`]
4. Skapa en vektor med slumptal av längd 50 från Poissonfördelningen med medelvärde 8, döp den till `minPoission`
5.  $X$  är poissonfördelad med medelvärde 3. Uppgiften är nu att skapa ett stapeldiagram (se koden nedan) över sannolikheterna att  $X$  antar värdena 0, 1, ..., 11, 12.

```
x <- dpois(x=0:12, lambda=3)
barplot(x, names.arg=0:12, main="Poisson(lambda=3)")
```

---

<sup>1</sup>Kolla även här eller här.



6. Låt  $X$  vara samma som ovan. Vad är sannolikheten att  $X$  antar följande värden [**Tips!** `ppois()`]:
  - (a) Sitt medelvärde?
  - (b) Mindre än eller lika med sitt medelvärde?
  - (c) Större än sitt medelvärde?
  - (d) Ett udda tal? Ni behöver inte ta hänsyn till tal som är större än 12.
  - (e) Talen 4 till 6?
7. Skapa en vektor med 50 slumpstal från t-fördelningen, med medelvärde 1 och 5 frihetsgrader (degrees of freedom), döp den till `minT`. **Tips:** För att ändra medelvärdet på en slumpvariabel kan en konstant adderas/subtraheras till variabeln.
8. Gör ett histogram över `minPoission` och över `minT`. Testa att ändra argumentet `breaks` i `hist()` till några olika värden.
9. Utgå från den normalfördelade variabeln  $X \sim N(\mu = 20, \sigma = 4)$  ( $\sigma$  är standardavvikelsen)
  - (a) Beräkna sannolikheten att variabeln antar värdet  $x \leq 25$ , Tips: `pnorm()`
  - (b) Beräkna sannolikheten att variabeln antar värdet  $x \geq 25$ , Tips: argumentet `lower.tail`.
  - (c) Säg att vi vill ta reda på vilket  $x$  som motsvarar en kumulativ sannolikhet på 0.75 (arean under täthetskurvan från vänster till  $x$ ). Detta kan göras med funktionen `qnorm()`. Kolla i dokumentation hur du ska använda `qnorm()` för att beräkna värdet på  $x$  som ger en kumulativ sannolikhet på 0.75.
10. Utgå från den normalfördelade variabeln  $X \sim N(\mu = 0, \sigma = 1)$ 
  - (a) Beräkna värdet på  $x$  som ger kumulativ sannolikhet på 0.025
  - (b) Beräkna värdet på  $x$  som ger kumulativ sannolikhet på 0.5
  - (c) Beräkna värdet på  $x$  som ger kumulativ sannolikhet på 0.975
  - (d) Beräkna värdet på  $x$  som ger kumulativ sannolikhet på 0.995

## 2.2 `sample()` och `set.seed()`

Vill vi dra slumpstal från ett antal element använder vi funktionen `sample()`. Med denna funktion kan vi dra ett stickprov (med eller utan återläggning) från en given vektor. Om vi exempelvis vill dra 5 slumpmässiga värden mellan 1 och 10 **med** återläggning gör vi det på följande sätt i R:

```
sample(x=1:10, size=5, replace=TRUE)
```

```
[1] 7 9 7 3 6
```

1. Dra ett stickprov ( $n = 5$ ) **utan** återläggning från sekvensen 10:20.
2. Upprepa uppgiften ovan men **med** återläggning.
3. Ett vanligt scenario är att vi vill dra ett slumpmässigt urval av rader från en data.frame (eller matris). Detta kan göras med hjälp av en indexvektor.

```
# slumpmässigt index:
rand_index<-sample(x = 1:nrow(mtcars),size = 5,replace = FALSE)
# dra ett slumpmässigt urval från data.frame:
mtcars[rand_index,]

# replace = TRUE, då kan size vara större än x
rand_index2<-sample(x = 1:nrow(mtcars),size = 100,replace = TRUE)
# dra ett slumpmässigt urval från data.frame
# där enskilda bilar finns med flera gånger
A<-mtcars[rand_index2,]
print(A)
dim(A)
dim(mtcars)
```

4. Vi kan på samma sätt dra slumpmässiga element från andra vektorer (exempelvis textvektorer). Tänk dig att du och några vänner ska organisera en fest. Skapa vektorn `namn`, som ska innehålla namnen för minst tre personer som textsträngar. Två personer behövs för att laga mat och två för att diska. Välj slumpmässigt vilka som ska göra de olika uppgifterna. Samma person ska kunna bli vald för båda uppgifterna.
5. Inte sällan vill vi att våra simuleringar och analyser ska vara reproducerbara, d.v.s. att vi ska få samma resultat varje gång vi gör en simulering. För detta används funktionen `set.seed()` och funktionen tar ett godtyckligt heltal för att initiera slumptalsgeneratoren i R. Upprepa uppgiften ovan två gånger, blir resultatet det samma? Upprepa två gånger till men kör först `set.seed(1234)` innan varje gång. Får du nu samma resultat?
6. Upprepa nu uppgiften ovan igen, men ändra argumentet `prob` så att de olika personer har olika sannolikhet att bli vald samt att du har sannolikheten 0 att bli vald. ;)

```
# ex på prob:
# kör flera gånger, vilka värden kommer med i urvalet?
sample(x=1:4, size=2, replace=FALSE,prob=c(0.8,0.1,0.05,0.05))

[1] 1 4

# med återläggning och många slumpdragningar:
y<-sample(x=1:4, size=2000, replace=TRUE,prob=c(0.8,0.1,0.05,0.05))
table(y)

y
 1    2    3    4
1563  217  106  114
```

## 2.3 Exempel: `sum_of_dice()`

1. Nu ska ni testa att simulera olika tärningskast. I uppgiften betyder D6 en vanlig 6-sidig tärning (med sidorna 1,2,3,4,5,6) där alla utfall har samma sannolikhet ( $1/6$ ).
  - (a) Skriv en funktion `my_dice()`, som genererar  $n$  stycken kast från en D6 och returnerar en vektor med kasten. [Tips! `sample()`]

- (b) Skriv en funktion `sum_of_dice()` som generar följande slumpstal:  $Y_k = \sum_{n=1}^N X_n$ , där  $X$  är ett tärningskast från funktionen `my_dice()`,  $N$  är ett heltal,  $k$  går från  $1, 2, 3, \dots, K$ . Funktionen ska ha  $N$  och  $K$  som argument, `sum_of_dice(N, K)`. Funktionen ska alltså kasta  $N$  stycken D6, summera dessa, sen upprepa det  $K$  stycken gånger.

```
sum_of_dice <- function(N,K){
  res <- integer(K)
  for (k in 1:K){
    res[k] <- sum(sample(1:6,size=N, replace=TRUE))
  }
  return(res)}

```

- (c) Testa nu `sum_of_dice(N, K)` med  $N = 3, 5$  och  $K = 100, 1000, 3000$ . Plotta resultaten i histogram och beräkna också medelvärde, standardavvikelse, min och max. Nedan ser ni några exempel på tester.

```
set.seed(3827)
sum_of_dice(N=3,K=10)

[1] 13  8  6 11 10 18  6  7 11 15

sum_of_dice(N=5,K=10)

[1] 20 24 16 24 21 19 17 11 16 20

```

- (d) Detta är ett klassiskt exempel på centrala gränsvärdessatsen.

## 2.4 Stora talens lag och Monte Carlo

En annan användning av slumpstal är den så kallade Monte Carlo metoden. I denna metod så använder vi slumpstal för att uppskatta väntevärden eller sannolikheter. Dessa väntevärden/sannolikheter skattas genom att vi gör en massa simuleringar och tar medelvärdet av dessa. En simulering kan t.ex. vara att dra ett slumpstal, det kan också vara lite mer komplicerat.

1. Om vi har ett slumpstal från  $\mathcal{U}(\min = 0, \max = 1)$  så kan vi räkna ut att väntevärdet är  $\frac{1}{2}$ . Testa om Monte Carlo metoden fungerar genom att köra följande.

```
set.seed(1234)
# Vi simulerar 10 slumpstal från fördelningen och räknar ut medelvärdet med mean()
u <- runif(10,min=0,max=1)
print(mean(u))

[1] 0.48923

# Nu simulerar vi istället 10 000 slumpstal från fördelningen och räknar ut medelvärdet med mean()
u <- runif(10000,min=0,max=1)
print(mean(u))

[1] 0.50022

```

2. Vi fortsätter med slumpstal från  $\mathcal{U}(\min = 0, \max = 1)$  så kan vi räkna ut att sannolikheten att få ett tal större än  $\frac{3}{4}$  är  $\frac{1}{4}$ . Vi kan räkna ut detta genom Monte Carlo metoden på följande sätt.

```

set.seed(1234)
# Vi simulerar 10 slumpstal
u <- runif(10,min=0,max=1)
# Vi skapar en logisk vektor med TRUE / FALSE om värdet är större eller mindre än 3/4
probs <- u > 3/4
# Vi uppskattar sannolikheten genom att räkna medelvärde på denna logiska vektor.
print(mean(probs))

[1] 0.1

# Vi upprepar allting med 10 000 slumpstal också
u <- runif(10000,min=0,max=1)
probs <- u > 3/4
print(mean(probs))

[1] 0.2465

```

3. Upprepa exemplena ovan och räkna ut medelvärdet för en normalfördelad med parametrarna  $\text{mean}=4$ ,  $\text{sd}=5$  samt räkna ut sannolikheten att variabeln är större än 9.

## Kapitel 3

# Introduktion till R markdown och knitr

R markdown och knitr är ett system för att skapa dynamiska rapporter med R, vilket innebär att vi väver ihop R-kod, data och text i ett och samma dokument., detta kallas "Literate programming". Med knitr och markdown kan vi skapa reproducerbara analyser med full spårbarhet hur beräkning, databearbetningar och grafik skapats. Vår slutprodukt kan bli rapporter i Word, PDF eller HTML.

R markdown består av två delar. Dels [markup-språket](#) markdown som används för att skapa text, ekvationer, bilder m.m. och dels knitr för att integrera dokumentet med R-kod. knitr kommer från ordet knit (sticka) och är en ordlek med innebörden att vi stickar ihop R med, i detta fall, markdown. knitr kan dock användas med andra ordbehandlare som L<sup>A</sup>T<sub>E</sub>X och L<sup>y</sup>X<sup>1</sup>.

Det som händer när vi **renderar** (= "kör") en Rmd fil är att knitr går igenom dokumentet, kör all R-kod och stoppar tillbaka svaret från R i markdownformat. När detta är gjort skapas ett word-, HTML- eller pdf-dokument från markdownfilen.

Det finns ett bra referensdokument [här](#) och en bra "fusklapp" [här](#).

Under senare tid har reproducerbarhet när det gäller statistiska analyser kommit att bli allt mer centralt. Reproducerbarhet innebär att ett experiment eller forskningsresultat ska kunna återupprepas - reproduceras - av andra forskare eller analytiker. Idag innebär ofta reproducerbarheten att det finns krav på att hela analyser, med både text, data och kod ska kunna reproduceras av andra forskare. Mer information om reproducerbar forskning finns [här](#).

### 3.1 Grunderna i markdown

1. För att skapa ett Rmd-dokument i R-Studio väljer vi New file → R markdown. Ange titel på dokumentet och HTML. Vi ska nu ha fått upp ett exempeldokument.
2. Vi borde fått upp ett dokument som borde se ut på följande sätt (första delen).

```
---
title: "My document"
author: "My name"
date: "1 januari 2021"
output: html_document
---

This is an R Markdown document. Markdown is a simple formatting
syntax for authoring HTML, PDF, and MS Word documents.
For more details on using R Markdown see <http://rmarkdown.rstudio.com>.
```

3. Prova att "knitta" dokumentet med knappen "knit HTML" i R-Studio, då skapas en HTML-fil baserat på R markdown-filen i samma mapp som .Rmd-filen.

---

<sup>1</sup>Alla dokument med laborationsinstruktioner i kursen har skapats med en kombination av L<sup>y</sup>X och knitr

4. Prova att “knitta” till PDF. För att göra detta måste en Tex-installation finnas på datorn. Det finns olika alternativ:

- (a) TinyTeX (rekommenderas), kan installeras med ett R-paketet tinytex, för mer info se här

```
install.packages("tinytex")
tinytex::install_tinytex()
# Avinstallera med:
# tinytex::uninstall_tinytex()
```

- (b) Andra alternativ är:

- i. Windows: MiKTeX eller TeXLive
- ii. Mac OS: MacTeX
- iii. Linux: texlive

5. Prova att “knitta” till word-dokument.

### 3.1.1 Grundläggande markdown

Lägg till det som framgår i kodblocken i ditt dokument och “knitta” ett HTML-dokument.

1. Rubriker skapas med #.

```
# Rubrik 1
## Rubrik 2
### Rubrik 3
```

2. Vill vi ha fet stil använder vi \*\* och vill vi ha kursiv stil använder vi \*.

```
I like both bold text and italic text.
```

3. Vill vi lägga till en länk använder vi hakparantes som anger länkens namn och därefter vanliga paranteser .

```
[Comics](http://xkcd.com/1239/) for the win!
```

4. Att lägga till en bild görs på ett liknande sätt som en länk, dock med ! innan hakparantesen.

```
![My pic](http://imgs.xkcd.com/comics/social_media.png)
```

5. Vill vi lägga till listor använder vi \* eller numrerar vår lista.

```
* lista (onummerad)
* knitr
+ R
+ text
+ bilder

1. lista (numrerad)
2. Rmd
+ knitr
+ md
```

6. Det går också enkelt att skapa tabeller.

Kolumnrubrik 1	Kolumnrubrik 2
Cell 1	Cell 2
Cell 3	Cell 4

7. Mer exempel på formatering finns i referensbladet för R markdown här.

### 3.1.2 Ekvationer

Ofta vill vi beskriva våra beräkningar med matematiska ekvationer. I markdown finns det möjlighet att skriva ekvationer med  $\text{\LaTeX}$  ekvationsystem. För exempel på hur det går att skriva  $\text{\LaTeX}$ -ekvationer, se följande [minihandbok](#).

För att skapa ekvationer använder vi  $\$$  och  $\$\$$ .

1. Vi börjar med att skapa följande ekvation i markdown.

$$a + b = 10$$

```
$$a+b=10$$
```

2. Vi kan skapa “inline”-ekvationer med  $\$$ . Denna ekvation:  $E = mc^2$ , skrivs i markdown som:

```
Denna ekvation: $E=mc^2$
```

3. Vi kan skapa vilka ekvationer vi vill med  $\text{\LaTeX}$  ekvationssystem. Fundera på koden nedan hur denna ekvation har byggts upp i  $\text{\LaTeX}$ .

$$\frac{\int_0^\infty x_a^2 dx}{10}$$

```
$$\frac{\int^{\infty}_{0}x^2_{a}dx}{10}$$
```

4. Skapa nu följande ekvationer i R markdown, ta hjälp av minihandboken vid behov.

- (a)  $y = \beta_0 + \beta_1 x$
- (b)  $\frac{(\alpha + \gamma)^2}{4}$
- (c)  $\sum_{i=1}^N x_i$

## 3.2 Integrera R-kod med knitr

Vi har nu gått igenom grunderna för att skapa ett markdowndokument. Den stora fördelen med markdown framgår dock först när vi kan integrera våra dokument med R-kod, grafik och tabeller.

Pröva att kopiera in koden nedan i ditt R-markdowndokument och “knitta” dokumenten till HTML eller pdf.

### 3.2.1 Inline-kod

1. Ibland kan vi vilja göra mindre, enklare beräkningar, direkt i ett dokument. För detta använder vi inlinekod. För att lägga till en enkel beräkning direkt i dokumentet används ``r [R-kod här] ``.



```
I know that 1 + 1 = `r 1 + 1`.
```

### 3.2.2 R-block - “chunks”

Med inline-kod kommer vi en liten bit mot ett dynamiskt dokument. Men vill vi ha mer komplicerade beräkningar, grafik eller tabeller måste vi skapa dessa i R-block, eller “chunks”. Dessa styrs med så knitr-alternativ. Samtliga knitr-alternativ för kodblock finns listade [här](#).

1. För att skapa ett R-block används följande kod.

```
```{r}
a <- 1 + 1
a
```
```

2. Om koden ovan körs kommer både R-koden att synas och resultatet som koden genererar. Vill vi inte visa R-koden anger vi bara `echo=FALSE` som knitr-alternativ.

```
```{r, echo=FALSE}
a <- 1 + 1
a
```
```

3. Vill vi istället dölja resultatet anger vi `results='hide'`.

```
```{r, results='hide'}
a <- 1 + 1
a
```
```

4. Vi kan också låta bli att köra koden med `eval=FALSE`.

```
```{r, eval=FALSE}
a <- 1 + 1
a
```
```

5. Det går självklart också att kombinera alternativ.

```
```{r, eval=FALSE, echo=FALSE}
a <- 1 + 1
a
```
```

### 3.2.3 Grafik

En av de stora fördelarna med knitr och Rmd är att vi i R kan skapa grafik som direkt skapas och sätts in i dokumentet.

1. För att lägga in ett diagram skapar vi diagrammet som vanligt i R. Under motorhuvens skapas en grafikfil som sätts in i dokumentet automatiskt.

```
```{r, echo=FALSE}
data(faithful)
hist(faithful$waiting)
```
```

2. Även här finns flera knitralternativ för att styra hur grafiken ska placeras i dokumentet. `fig.height` och `fig.width` styr storleken och `fig.align` styr om diagrammet ska vara höger-, vänster- eller centrerat till mitten.

```
```{r, echo=FALSE, fig.width=3, fig.height=3, fig.align='center'}
data(faithful)
hist(faithful$waiting)
```
```

### 3.2.4 Tabeller

Utöver grafik och text är tabeller vanligt i statistiska analyser och rapporter. Med funktionen `kable()` i knitr-paketet kan vi skapa tabeller direkt från R-objekt.

1. För att skapa tabeller direkt behöver vi dels använda funktionen `kable()`, men vi behöver också ange att funktionen direkt skapar en markdowntabell som ska ses som markdownkod. Därför behöver vi ange att resultatet ska vara `'asis'`.

```
```{r, echo=FALSE, results='asis'}
knitr::kable(head(faithful))
```
```

2. Vi kan, precis som med grafiken, styra tabellernas grundutseende en del.

```
```{r, echo=FALSE, results='asis'}
knitr::kable(head(faithful), digits = 2, align = c("l", "c"))
```
```

3. Det finns mer avancerade funktioner för att skapa tabeller automatiskt från ex. linjära regressionsmodeller med paketen `xtable` och `tables`.

## Kapitel 4

# Input och output (I/O): Data på webben

Allt mer data lagras på webben med olika former av molntjänster. Anledningen till att data lagras på webben kan vara flera:

- Reproducerbarhet/öppna data för andra forskare
- Samarbete kring datainsamling
- Stabilitet
- Versionshantera förändringar i data

Nedan finns exempel från vanliga molntjänster och hur man kan ladda ned data från dessa direkt till R. Ett av paketen för att hantera denna typ av datainläsning är **repmis**, som kan läsa in de flesta csv, xlsx och Rdata-filer. För att ladda ner filer från Dropbox till R kan paketet **rdrop2** användas, men det kommer inte vara en del av denna kurs, den intresserade kan läsa här.

### 4.1 downloader

Ibland kan det vara så att vi vill ladda ned filer från R, men inte läsa in dem. Vi kanske vill ladda ned ett antal filer och sedan läsa in dem en och en. För att ladda ned filer i R finns funktionen **download.file()**. Dock kan det ibland vara lite klurigt att få den att fungera för så kallade secure http (https) adresser. Av bekvämlighet har därför paketet **downloader** skapats som gör nedladdning av filer mycket enkelt och bekvämt oberoende av operativsystem.

1. För att ladda ned data anger vi dels sökvägen till den aktuella filen och sedan sökvägen dit filen ska laddas ned. I detta fallet laddas filen ned till min working directory.

```
library(downloader)
apple_remote <- "https://raw.githubusercontent.com/STIMALiU/KursRprgm2/main/Labs/DataFiles/Apple.txt"
apple_local <- paste0(getwd(), "/Apple.txt")
download(url = apple_remote, destfile = apple_local)
# den här filen hamnar i working directory.
apple_test <- read.table(file="Apple.txt", sep=";", header=TRUE)
head(apple_test, 3)
```

|   | Date       | Open   | High   | Low    | Close  | Volume   | Adj.Close |
|---|------------|--------|--------|--------|--------|----------|-----------|
| 1 | 2012-01-24 | 425.10 | 425.10 | 419.55 | 420.41 | 19226900 | 420.41    |
| 2 | 2012-01-23 | 422.67 | 428.45 | 422.30 | 427.41 | 10915800 | 427.41    |
| 3 | 2012-01-20 | 427.49 | 427.50 | 419.75 | 420.30 | 14758300 | 420.30    |

## 4.2 Github

github är framförallt en tjänst för att versionshantera programkod i molnet. Dock används det också mycket för att lagra enklare datamaterial. Särskilt material som förändras mycket där vi behöver kunna följa vilka förändringar som gjorts.

Kör nedanstående exempel för att läsa in filen polls från github. Filen innehåller (nästan) samtliga opinionsmätningar i Sverige sedan 1998. Mer information finns [här](#).

```
library(repmis)
data_url <- "https://github.com/MansMeg/SwedishPolls/raw/master/Data/Polls.csv"
polls <- repmis::source_data(data_url, sep = ",", dec = ".", header = TRUE)
```

*Downloading data from: https://github.com/MansMeg/SwedishPolls/raw/master/Data/Polls.csv*

*SHA-1 hash of the downloaded data file is:*  
*39a3efb02ab100ff0b642592d09b45ab9b21b56b*

```
head(polls[,1:12])
```

|   | PublYearMonth | Company  | M    | L   | C   | KD  | S    | V    | MP  | SD   | FI | Uncertain |
|---|---------------|----------|------|-----|-----|-----|------|------|-----|------|----|-----------|
| 1 | 2021-dec      | Ipsos    | 21.0 | 4.0 | 6.0 | 6.0 | 31.0 | 10.0 | 3.0 | 18.0 | NA | 14        |
| 2 | 2022-jan      | Sifo     | 20.9 | 2.9 | 7.0 | 4.4 | 31.1 | 9.6  | 3.0 | 19.2 | NA | NA        |
| 3 | 2022-jan      | Demoskop | 21.8 | 3.2 | 8.3 | 5.4 | 29.2 | 9.3  | 2.8 | 17.8 | NA | NA        |
| 4 | 2022-jan      | Novus    | 18.7 | 2.2 | 6.7 | 5.6 | 31.3 | 10.5 | 3.4 | 20.1 | NA | 8         |
| 5 | 2021-dec      | Sentio   | 20.9 | 2.9 | 6.8 | 5.9 | 28.2 | 8.5  | 4.3 | 20.0 | NA | NA        |
| 6 | 2021-dec      | Ipsos    | 20.0 | 3.0 | 7.0 | 5.0 | 29.0 | 10.0 | 4.0 | 20.0 | NA | 12        |

## 4.3 Google docs

Google docs eller google drive är en molntjänst för kalkylblad i molnet. Vi ska nu pröva att läsa in ett publicerat kalkylblad. [Här](#) finns dokumentet vi ska läsa in (det är samma data som faithful-datasetet i R). I sökvägen kan vi se textsträngen 1ZDJQXUiYg\_QPNY38SNyTHyOtc06kbwePmTHWkqxnXxA. Detta är det unika id:t för detta google-kalkylblad.

Vi kommer använda paketet googlesheets4 som klarar att göra mycket med google docs (som att spara dataset, läsa textfiler m.m.). Vi kommer dock endast använda det för att läsa in data i R.

För att läsa in vårt material gör vi det i tre steg.

1. Först läser vi in paketet.

```
library(googlesheets4)
```

2. I nästa steg laddar vi ned information om det aktuella kalkylbladet (hur många blad m.m.). I detta fall är det ett publikt datablad och vi behöver inte ange några användarnamn, däröfr kör vi `gs4_deauth()` först för att paketet inte ska fråga efter inloggning. För egna blad behöver vi logga in och skapa en koppling till vårt konto.

```
gs4_deauth()
google_worksheets <- read_sheet("1ZDJQXUiYg_QPNY38SNyTHyOtc06kbwePmTHWkqxnXxA")

v Reading from "Faithful".
v Range 'faith'.
```

3. Nu är bladet inläst och vi kan komma åt de olika kolumnerna som vanligt

```
head(google_worksheets)
```

```
# A tibble: 6 x 2  
  eruptions waiting  
    <dbl>    <dbl>  
1     3.6       79  
2     1.8       54  
3     3.33      74  
4     2.28      62  
5     4.53      85  
6     2.88      55
```

# Kapitel 5

## pxweb

Statistiska centralbyrån har utvecklat ett API för att ladda ned data direkt utan att först gå via deras webbplats. En koppling till deras API finns installerat som paketet **pxweb** i R. Detta paket fungerar för samtliga pxweb-apier. Men det största api:et är utan tvekan SCB:s.

1. Börja med att installera paketet **pxweb** och läs in det i R.

```
install.packages("pxweb")  
library(pxweb)
```

2. Det går att få en snabb introduktion till detta paket med funktionen **vignette()**.

```
vignette(topic="pxweb")
```

### 5.1 Navigera i SCB:s API

Ofta vet vi inte exakt vad för data vi vill ha utan vill navigera igenom SCB:s databaser på ett effektivt sätt. Detta gör vi med funktionen **pxweb\_interactive()**. Vill vi leta upp ett givet datamaterial som vi vill spara måste vi tillskriva datamaterialet till ett objekt.

Statistiska centralbyråns API består av två delar:

- Navigera mellan de olika datamaterialen.
  - Välja ut delar av datamaterialet vi vill ladda ned.
1. För att navigera i Statistiska centralbyråns API kan vi använda funktionen **pxweb\_interactive()** som listar förinstallerade api:er och ger möjlighet att välja vilket api vi vill hämta data från.

```
mitt_data <- pxweb_interactive()
```

2. För att navigera skrivs kommandon, siffror anger vilket menyval vi vill göra, bokstäver och tecken används för att backa/lista alla osv.
3. För att komma åt SCB välj först 1, sen 1 igen för att välja version 1, sen välj 2 för att ha på svenska och slutligen 1 för att komma åt databsen.
4. Vi ska nu leta upp andelen arbetslösa från den senaste arbetskraftsundersökningen och spara ned detta som ett datamaterial i R. Vi befinner oss nu i den översta menyn. Ange 2 för att gå vidare till menyn för statistik om Näringsverksamhet:
5. Prova att "gå tillbaka" till huvudmenyn med b. Gå sen in på **Arbetsmarknad**, använd a för att se alla alternativ om de är dolda.

6. Navigera dig fram till följande datamaterial som vi ska läsa in:  
[NAKUArblosaTAr] Arbetslösa 15-74 år (AKU) efter arbetslöshetstidens längd, kön och ålder. År 2005 - 2021
7. Nu är vi inne i den del av API:et som anger vilka delar av materialet du vill ladda ned. Varje variabel kommer nu dyka upp och vi får ange vilka data vi vill ha.  
Vill vi bara ha en kategori anger vi den siffran, vill vi ha allt material anger vi \* och vill vi ha delar anger vi det antingen som en sekvens med : eller avgränsat med , .
  - (a) För variabel ARBETSLÖSHETSTID ange **allt** (med \*) som den del av materialet du vill ha.
  - (b) För variabel KÖN ange **totalt** som den del du vill ha
  - (c) För variabel ÅLDER ange grupperna 15-24, 25-54 och 55-74 som den del du vill ha.
  - (d) För variabel CONTENTSCODE ang **Arbetslösa, 1000-tal** som den variabel du vill ha.
  - (e) Variabeln TID innehåller många kategorier. För att se alla kategorier, ange **a**.
  - (f) Välj nu ut tidsperioden 2011 till 2021.
8. Nu kommer det komma en massa frågor om hur du vill ladda ner datamaterialet. Först frågar den om du vill få koden för att ladda ner data direkt, välj y.
9. Sen frågar den om du vill ha resultatet som en JSON format eller R list, välj n.
10. R frågar nu om du vill ladda ned datamaterialet, ange y.
11. Nästa steg efterfrågas om du vill ha materialet i originalformat eller som en färdigformaterad data.frame i R., välj y.
12. Som ett sista steg efterfrågas nu om du vill ha citeringen för datamaterialet, välj y.
13. Nu ska materialet laddas ned. Du borde också få ut följande kod:

```
# PXWEB query
pxweb_query_list <-
  list("Arbetsloshetstid"=c("TOT", "1V", "2V", "3-4V", "5-26V", "27V-", "uppgsagn"),
       "Kon"=c("1+2"),      "Alder"=c("15-24", "25-54", "55-74"),
       "ContentsCode"=c("AM0401G0"),
       "Tid"=c("2011", "2012", "2013", "2014", "2015", "2016", "2017", "2018", "2019", "2020", "2021"))
# Download data
px_data <-
  pxweb_get(url = "http://api.scb.se/OV0104/v1/doris/sv/ssd/AM/AM0401/AM0401L/NAKUArblosaTAr",
            query = pxweb_query_list)
# Convert to data.frame
px_data_frame <- as.data.frame(px_data, column.name.type = "text", variable.value.type = "text")
```

14. Titta på det material du laddat ned.

### 5.1.1 Ladda ned data från SCB direkt med kod

Ofta vill vi ladda ned/komma åt data från SCB som en löpande del i en analysprocess där vi använder de senaste data från SCB. Då vill vi använda pxweb, inte interaktivt, utan som R-kod.

1. Spara ned koden du fick ovan och kör den för att ladda ned datamaterialet direkt. (Nu sparas den som px\_data\_frame).
2. Vi vill nu ladda ned data för alla tidpunkter. Ändra variabeln tid till '\*' och ladda ned datat på nytt.

## 5.2 Namnstatistik 2021

Använd pxweb för att ta reda på vilka namn som var vanliga att ge till nyfödda barn 2021. Målet är att hitta de 5 vanligaste flicknamnen och de 5 vanligaste pojknamnen som gavs till nyfödda barn 2021. Detta datamaterial finns under "[BE] Befolkning".

## 5.3 Partistorlek i valet till kommunfullmäktige 2018

Målet är nu att undersöka hur stora olika partier blev i kommunfullmäktige i valet 2018. Ni ska använda datamaterialet

"[ME0104T1] Kommunfullmäktigval - valresultat efter region och parti mm. Antal och andelar. Valår 1973 - 2018"

Detta datamaterial finns under "[ME] Demokrati". Målet är att hitta vilken kommun där varje parti blev som störst i kommunfullmäktige jämfört med alla kommuner i Sverige. Om ni tittar på ett parti, tex Socialdemokraterna, då vill ni ta reda på vilken kommun där de fick högst andel röster (jämfört med alla kommuner) i valet till kommunfullmäktige 2018. Sen ska ni upprepa detta för Moderaterna, Sverigedemokraterna och alla andra partier som sitter i riksdagen. I slutändan så ska ni få fram 8 kommuner, en för varje parti. Ni vill alltså svara på frågorna: "Vilken var Socialdemokraternas starkaste kommun i valet till kommunfullmäktige?", "Vilken var Moderaternas starkaste kommun i valet till kommunfullmäktige?", osv.

Tips på lösningsförslag:

- Ladda ner det datamaterialet från SCB med pxweb. Välj andelar och år 2018. Spara som en data.frame.
- Välj ut alla rader som motsvarar ett parti och spara som en ny variabel. Undvik att ta med "Riket".
- Sortera med avseende på andelar och spara namnet på den kommun som hade högst andel. Tips: `order()`
- Upprepa sedan för alla de 8 riksdagspartierna.
- Spara sedan ert resultat i en data.frame med tre kolumner: en med parti, en med kommun och en med andel röster.

## 5.4 \* Extraproblem: Partistorlek i valet till kommunfullmäktige 2014

Upprepa uppgiften ovan, men för valet 2014. Jämför sedan resultatet mellan 2014 och 2018.

## 5.5 \* Extraproblem: Andra api:er

Vi har nu prövat Statistiska centralbyråns API, men allt fler offentliga myndigheter (och företag) lägger ut sina resultat i form av ett pxweb-api. Med `api_catalogue()` är det möjligt att se vilka andra api:er som nu finns inkluderade.

1. Prova att navigera i något annat API än Statistiska centralbyråns och ladda ned data därifrån.



# Litteraturförteckning

- [1] Edward R Tufte and PR Graves-Morris. *The visual display of quantitative information*, volume 2. Graphics press Cheshire, CT, 1983.

## Del II

# Inlämningsuppgifter

## Inlämning

Utgå från laborationsmallen, som går att ladda ned här, när du gör inlämningsuppgifterna. Spara denna som labb[no]\_[liuID].R , t.ex. labb1\_\_josad732.R om det är laboration 1. Ta inte med hakparenteser i filnamnet. Denna fil ska laddas upp på LISAM och ska **inte** innehålla något annat än de aktuella funktionerna, namn- och ID-variabler och ev. kommentarer. Alltså **inga** andra variabler, funktionsanrop för att testa inlämningsuppgifterna eller anrop till markmyassignment-funktioner.

## Tips!

Inlämningsuppgifterna innebär att konstruera funktioner. Ofta är det bra att bryta ned programmeringsuppgifter i färre små steg och testa att det fungerar i varje steg.

1. Lös uppgiften med vanlig kod direkt i R-Studio (precis som i datorlaborationen ovan) utan att skapa en funktion.
2. Testa att du får samma resultat som testexemplen.
3. Implementera koden du skrivit i 1. ovan som en funktion.
4. Testa att du får samma resultat som i testexemplen, nu med funktionen.

## Automatisk återkoppling med markmyassignment

Som ett komplement för att snabbt kunna få återkoppling på de olika arbetsuppgifterna finns paketet **markmyassignment**. Med detta är det möjligt att direkt få återkoppling på uppgifterna i laborationen, oavsett dator. Dock krävs internetanslutning.

Information om hur du installerar och använder **markmyassignment** för att få direkt återkoppling på dina laborationer finns att tillgå [här](#).

Samma information finns också i R och går att läsa genom att först installera **markmyassignment**.

```
install.packages("markmyassignment")
```

Om du ska installera ett paket i PC-pularna så behöver du ange följande:

```
install.packages("markmyassignment", lib="mapp i din hemkatalog")
```

Tänk på att i sökvägar till mappar/filer i R i Windowssystem så används "\\", tex "C:\\Users\\Josef".

Därefter går det att läsa information om hur du använder **markmyassignment** med följande kommando i R:

```
vignette("markmyassignment")
```

Det går även att komma åt vignetten [här](#). Till sist går det att komma åt hjälpfilerna och dokumentationen i **markmyassignment** på följande sätt:

```
help(package="markmyassignment")
```

Lycka till!

## Kapitel 6

# Inlämningsuppgifter

För att använda `markmyassignment` i denna laboration ange:

```
library(markmyassignment)
lab_path <-
  "https://raw.githubusercontent.com/STIMALiU/KursRprgm2/main/Labs/Tests/d5.yml"
suppressWarnings(set_assignment(lab_path))
```

*Assignment set:*

*D5: Statistisk programmering med R: Lab 5*

*The assignment contain the following (2) tasks:*

- *estimate\_pi*
- *sum\_of\_random\_dice*

### 6.1 estimate\_pi()

Talet  $\pi$  är en matematisk konstant som dyker upp lite här och var, framförallt när man jobbar med cirklar och sfärer men även inom sannolikhetslära och fysik är den viktig. Talet är ett så kallat irrationellt tal, det betyder att det har en oändlig decimalutveckling utan repetitioner och saknar ett exakt uttryck. Man vill ofta försöka approximera denna konstant och nu ska vi göra det med hjälp av slumpstal och Monte Carlo metoden.

Om vi tänker oss att vi väljer ett  $x$  och  $y$  värde slumpmässigt mellan  $-1$  och  $1$ . Om vi tar dessa värden tillsammans blir det en punkt i kvadraten med sidlängd  $2$  centrerat runt origo. Inuti denna kvadrat kan vi rita in enhetscirkeln, det är cirkeln med radie  $1$  centrerad i origo. Vi kan nu beräkna sannolikheten att vår slumpmässigt utvalda punkt  $(x,y)$  hamnar inuti denna cirkel, den sannolikheten är  $\frac{\pi}{4}$ , arean av cirkeln (där vi vill hamna) delat med arean av kvadraten (där vi kan hamna). För att uppskatta  $\pi$  kan vi då slumpa massa  $(x,y)$  punkter räkna ut hur många av de som är inom enhetscirkeln, kvoten mellan de inom enhetscirkeln och antalet punkter totalt blir vår uppskattning av  $\frac{\pi}{4}$ . Multiplicera kvoten med  $4$  för att få ett estimat av  $\pi$ . Ett tips, för att se om en punkt  $(x,y)$  är inom enhetscirkeln kan vi testa om  $x^2 + y^2 \leq 1$ .

Vi ska nu skriva en funktion `estimate_pi()` som tar argumenten:

- $N$  är ett positivt heltal som är antalet slumpade punkter  $(x,y)$ .
- `my_seed` värdet som styr slumpalsgeneratorn, default ska vara `NULL`.

Funktionen ska returnera en lista som innehåller två variabler:

- `est` ert estimat på  $\pi$ .
- `punkter` som är en `data.frame` med två kolumner, en är  $x$  och den andra  $y$ . Dessa kolumner ska innehålla alla era slumpade värden.

Ett förslag på implementation kan vara följande:

- Ändra seeden till: `set.seed(my_seed)`

- Slumpa x och y värden med `runif()`. Tänk på att `min=-1`, `max=1`.
- Stoppa in era slumpade värden i en `data.frame`.
- Skapa en logisk vektor med `TRUE` eller `FALSE` om punkterna är inom enhetscirkeln eller inte.
- Beräkna ert estimat genom att räkna medelvärdet på den logiska vektorn.
- Skapa och returnera en lista med rätt namn. Använd t.ex. `list(est = ... , punkter = ... )`

Här kommer exempel på hur den kan fungera:

```
estimate_pi(N = 5, my_seed = 1234)

$est
[1] 3.2

$punkter
      x      y
1 0.60392 0.74118
2 0.24132 0.55451
3 0.80760 0.54396
4 -0.94695 0.92660
5 -0.40294 0.18943

output <- estimate_pi(N = 100, my_seed = 1234)
output$est

[1] 3.32

output <- estimate_pi(N = 10000, my_seed = 1234)
output$est

[1] 3.1624

output <- estimate_pi(N = 10000, my_seed = 445)
output$est

[1] 3.134
```

## 6.2 sum\_of\_random\_dice()

Funktionen `sum_of_dice()` i sektionen Slumptal och Statistik summerar värdet från ett fixt antal tärningar, nu ska ni skriva en funktion som kan summera ett slumpmässigt antal tärningar. Antalet tärningar ska vara Poissonfördelat med en parameter  $\lambda$  på följande sätt.

$$Y \sim \text{Po}(\lambda)$$

$$X = \sum_i^Y Z_i$$

där  $Z_i$  är utfallet från en sexsidig tärning. Skapa en funktion `sum_of_random_dice()`. Argumenten ska vara:

- `K`: antal dragningar från slumpfördelningen
- `lambda`: ett positivt kontinuerligt tal (parameter i poissonfördelningen)
- `my_seed`: ett slumpfrö som styr slumptalsgenereringen, default ska vara `NULL`.

Funktionen `sum_of_random_dice()` ska returnera en `data.frame` med summan av ögonen på de slumpmässigt antalet tärningar samt antalet kastade tärningar.

Ett förslag på hur detta kan implementeras finns här:

- Ändra seeden till: `set.seed(my_seed)`.
- Sätt upp en tom `data.frame` som namnges `result`, som ska ha K rader och 2 kolumner. Första kolumnen ska ha namnet `value` och den andra `dice`.
- Gör följande for-loop över vektorn `1:K`
  - Dra ett slumpstal från en poissonfördelning med parameter `lambda`. Spara slumptalet i `current_number`, vilket är antalet tärningar. Spara `current_number` i kolumnen `dice` på aktuell rad i `result`.
  - Anropa `sum_of_dice()` med argumenten `K=1` och `N=current_number`, spara resultatet i kolumnen `value` på aktuell rad i `result`. Observera att om `current_number=0` så ska summan bli 0.
- Returnera `result`.

**Obs!** Notera att funktionen `sum_of_dice()` ska vara inkluderad i filen som ni lämnar in.

**Obs!** Om du inte implementerar funktionen på detta sätt kan den fortfarande fungera korrekt, men eftersom slumptalen används i olika ordning kan det vara så att du inte får samma resultat som i exemplen nedan.

Testa om testfallen nedan fungerar:

```
sum_of_random_dice(K=5,lambda=3,my_seed=42)

  value dice
1    13    5
2     8    4
3    18    6
4    26    6
5    10    4

sum_of_random_dice(K=5,lambda=8,my_seed=4711)

  value dice
1    47   13
2    21    5
3    30    8
4    28    9
5    20    8

# No seed give different results
sum_of_random_dice(K=5,lambda=3)

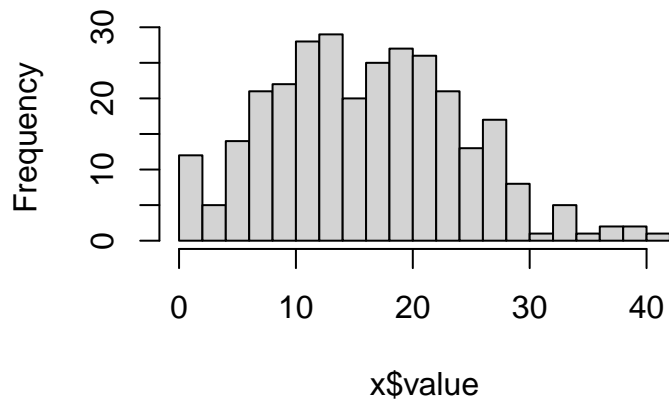
  value dice
1     3    1
2    13    4
3    15    4
4    14    4
5     1    1

sum_of_random_dice(K=5,lambda=3)

  value dice
1    12    3
2    14    3
3     5    2
4    16    3
5    11    3

x <- sum_of_random_dice(K=300,lambda=5,my_seed=42)
hist(x$value, 20)
```

### Histogram of x\$value



```
mean(x$value)

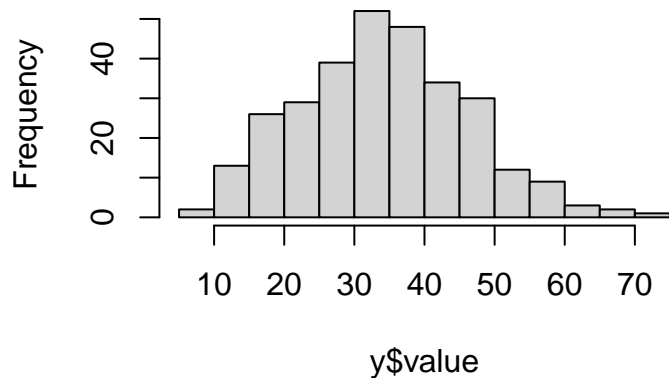
[1] 16.61

sd(x$value)

[1] 8.1899

y <- sum_of_random_dice(K=300, lambda=10, my_seed=4711)
hist(y$value, 20)
```

### Histogram of y\$value



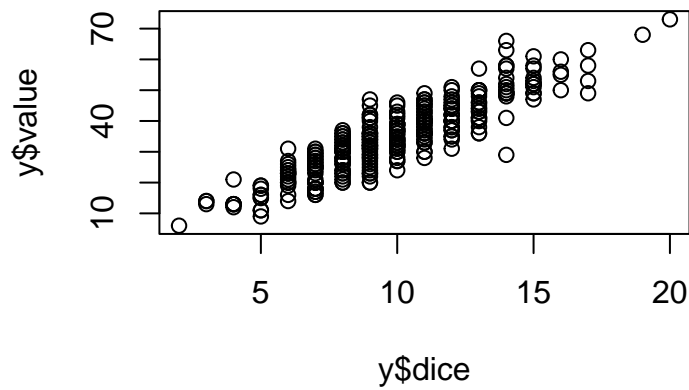
```
mean(y$value)

[1] 34.72

sd(y$value)

[1] 12.017

plot(y$dice, y$value)
```



### 6.3 Miniprojektet del 1

En del av denna laboration är att genomföra miniprojektet del 1. Se denna pdf för detaljer, finns även länk på kurshemsidan. Notera att det är en separat inlämning för miniprojektet del 1 på Lisam. I miniprojektet får (ska) ni använda `pxweb`, ni ska ha kod som fungerar och skapar reproducerbara rapporter. OBS: inga varningar ska synas i er färdiga rapport!