

Datorlaboration 6

Måns Magnusson

25 januari 2015

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.
 - Det är tillåtet att diskutera med andra grupper, men att plagiera eller skriva kod åt varandra är **inte tillåtet**.
 - Utgå från laborationsfilen som går att ladda ned [här](#)
 - Laborationen består av två delar:
 - Datorlaborationen
 - Inlämningsuppgifter
 - I laborationen finns det extrauppgifter markerade med *. Dessa kan hoppas över.
 - Deadline för labben framgår på [kurshemsidan](#)
-

Innehåll

I	Datorlaboration	3
1	Introduktion till grafik	4
1.1	Visualisera en variabel	4
1.1.1	Cirkeldiagram	4
1.1.2	Barcharts	5
1.1.3	Histogram	7
1.2	Visualisering i flera variabler	8
1.2.1	Sambandsdiagram	8
1.2.2	Linjediagram	9
1.2.3	Boxplot	9
1.3	Grafiska inställningar och tillägg	10
1.4	Spara figurer	11
1.5	* Extraproblem: Skapa en egen graf	11
2	Introduktion till R markdown och knitr	13
2.1	Grunderna i markdown	13
2.1.1	Grundläggande markdown	14
2.1.2	Ekvationer	15
2.2	Integrera R-kod med knitr	15
2.2.1	Inline-kod	15
2.2.2	R-block - "chunks"	15
2.2.3	Grafik	16
2.2.4	Tabeller	17
3	Input och output (I/O): Data på webben	18
3.1	Github	18
3.2	Google docs	19
3.3	Dropbox	19
4	pxweb	21
4.1	Navigera i SCB:s API	21
4.2	Ladda ned data från SCB direkt med kod	22
4.3	* Extraproblem: Andra api:er	23
II	Inlämningsuppgifter	25
5	Inlämningsuppgifter	27
5.1	fast_swe_pop()	27
5.2	gram_schmidt()	28
5.3	Miniprojektet del I	29

Del I

Datorlaboration

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

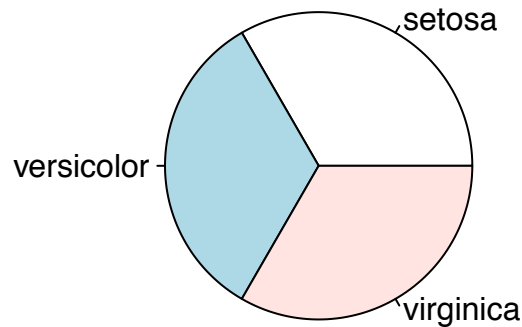
Cirkeldiagram är ofta populärt, men bör generellt sett undvikas (se ex. [1]). Det är en typ av diagram som inte lämpas sig bra för människors tänkande och många har svårt att jämföra tårtbitar visuellt, [här](#) finns de vanligaste argumenten. För att skapa cirkeldiagram använder vi funktionen `pie()`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

1. Som ett första steg måste vi beräkna frekvenser för den kategoriska variabel vi vill visualisera med `table()`:

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

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"))
```

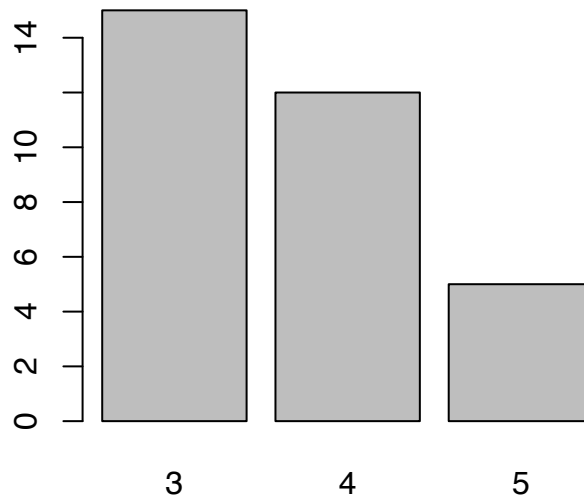
5. Vi kan självklart styra ännu mer i utformningen av cirkeldiagrammen. För mer hjälp för cirkeldiagram använd `?pie`. Men som sagt undvik cirkeldiagram i största möjliga mån.

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("red", "blue"))
```

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

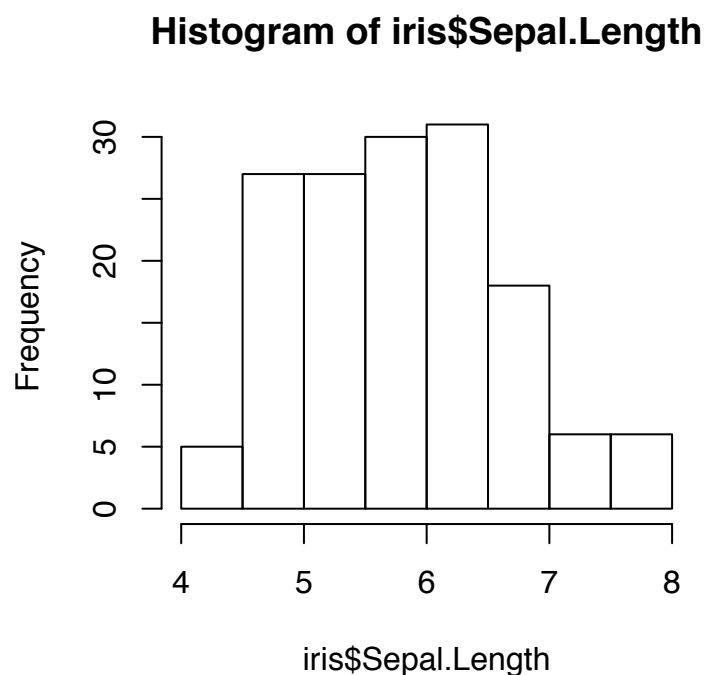
```
barplot(counts, main="Car Distribution by Gears and VS", xlab="Number of Gears", col=c('
```

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 följande sätt att skapa ett histogram:

```
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)
```

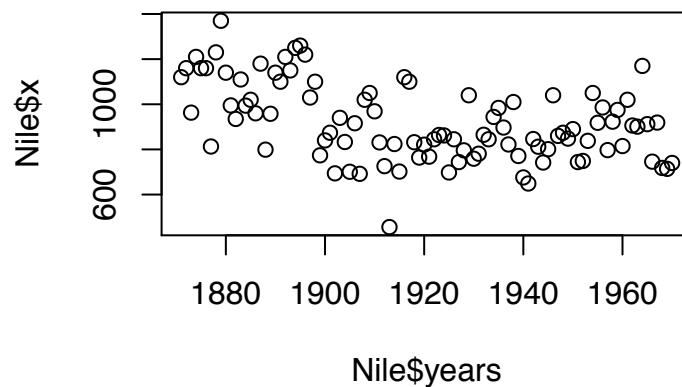

1.2 Visualisering i flera variabler

1.2.1 Sambandsdiagram

En av de vanligaste sätten att visualisera tvådimensionella data är med scatter plots. 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 och i många fall använder vi `plot` för att visualisera olika typer av objekt. Grundutförandet ger dock en scatterplot på följande sätt:

```
plot(Nile$years, Nile$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`.

```
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)
```

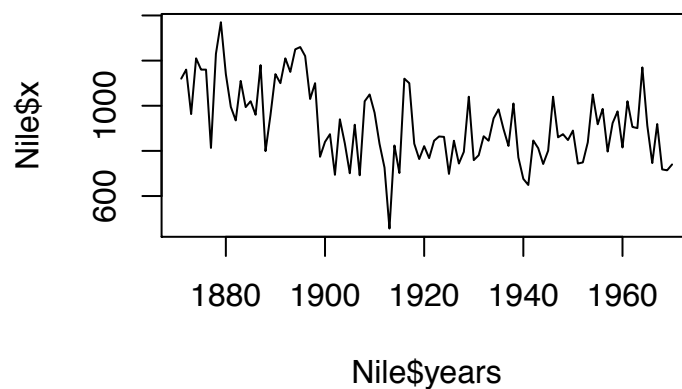
5. Studera hur `colVector` ser ut ovan så du förstår hur vektorn används för att styra färgen på punkterna. Prova att använd ytterligare än färg till punkterna efter 1945.
6. Vi kan också ändra hur punkterna ser ut och använda andra symboler. Det finns totalt 25 olika symboler i baspaketet. För att använda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

```
plot(1:25, 1:25, type="p", pch = 1:25)
```

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 linjegrav 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.

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

3. Vill vi både ha linjer och punkter kan vi använda `type='lo'`.

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

4. En sista typ av linjegrav är en trappstegsgraf:

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

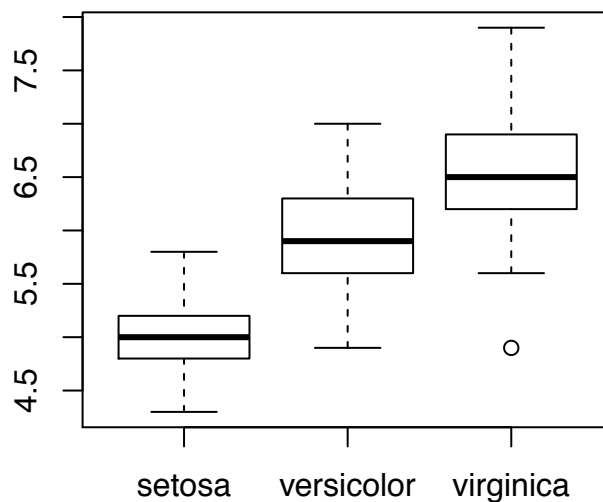
5. Prova att lägga till argumentet `lwd=3` i plotten ovan. Vad innebär detta?

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



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

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

3. 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. Ned 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 (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))`. Det kan vara krångligt att förstå exakt vad detta betyder. Men det vi säger till R med detta kommando är att vi vill ha tre rader och två kolumner med figurer. 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.

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”.
3. 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")
```

4. Prova nu att själv skriva ut en av dina figurer ovan i PNG- och TIFF-format.

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'`.
 - (c) Ändra nu rubriken och axeltexterna med `xlab=` och `ylab=`.
 - (d) Spara ned de två histogrammen i en figur (förlagsvis 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

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. 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 \LaTeX och \LaTeX .

Det som händer när vi **renderar** 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. 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 [\[här\]](#).

2.1 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 2015"
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. Prova att “knitta” till både PDF, HTML och word-dokument.
4. I R-Studio finns hjälp om R-markdown med knappen “?”. Prova att få fram hjälp på detta sätt.

2.1.1 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].

2.1.2 Ekvationer

Ofta vill vi beskriva våra beräkningar med matematiska ekvationer. I markdown finns det möjlighet att skriva ekvationer med \LaTeX ekvationsystem. För exempel på hur det går att skriva \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 \LaTeX ekvationssystem. Fundera på koden nedan hur denna ekvation har byggts upp i \LaTeX .

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

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

2.2 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.

2.2.1 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`.
```

2.2.2 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
```
```


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

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

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

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

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

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

- Det går självklart också att kombinera alternativ.

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

2.2.3 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.

- 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)
```
```

- 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)
```
```

2.2.4 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 3

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.

3.1 Github

github.com 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:
dc8c1f2931339df8d26756454bd95634da214109

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

| | PublYearMonth | Company | M | FP | C | KD | S | V | MP | SD | FI | Uncertain |
|---|---------------|----------|------|-----|-----|-----|------|-----|-----|------|-----|-----------|
| 1 | 2015-jan | Novus | 24.1 | 4.3 | 6.6 | 4.3 | 30.0 | 5.3 | 6.4 | 16.5 | 2.2 | 8.1 |
| 2 | 2015-jan | Inizio | 24.8 | 3.5 | 6.1 | 4.9 | 33.0 | 5.4 | 5.5 | 13.6 | 1.9 | 8.6 |
| 3 | 2015-jan | Sentio | 23.7 | 4.8 | 5.4 | 3.9 | 29.7 | 7.0 | 5.6 | 15.8 | 2.8 | NA |
| 4 | 2015-jan | Sifo | 23.1 | 5.0 | 6.7 | 3.4 | 33.3 | 5.6 | 6.7 | 13.4 | 1.7 | NA |
| 5 | 2015-jan | Demoskop | 24.4 | 4.8 | 5.3 | 4.1 | 32.3 | 5.3 | 6.6 | 14.4 | 2.1 | 10.5 |
| 6 | 2014-dec | YouGov | 23.2 | 3.7 | 5.7 | 5.5 | 30.6 | 6.0 | 5.7 | 16.9 | 2.0 | NA |

3.2 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 16uE3bb_7rHt_g4zOBXAVHM40N3YqRIXLPQwsE1yenAQ. Detta är det unika id:t för detta google-kalkylblad.

Vi kommer använda paketet RGoogleDocs 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 och skapar vi en koppling till det unika google doc id:t.

```
library(RGoogleDocs)

Loading required package: methods
Attaching package: 'RGoogleDocs'
The following object is masked from 'package:methods':
  getAccess

google_ws_object <- publicWorksheet("16uE3bb_7rHt_g4zOBXAVHM40N3YqRIXLPQwsE1yenAQ")
```

2. I nästa steg laddar vi ned information om det aktuella kalkylbladet (hur många blad m.m.). Under con anger vi en koppling till google doc. Är ett dokument öppet kan vi ange NULL. Annars behöver vi ange vårt google account om vi exempelvis vill komma åt dokument som inte är publika (ex. våra egna data från ett google forms).

```
google_worksheets <- getWorksheets(google_ws_object, con = NULL)
```

3. Det sista steget är att läsa in dokumentet i R. Det kan vi göra genom att ange vilket blad vi vill läsa in.

```
my_faithful <- sheetAsMatrix(google_worksheets[[1]], con = NULL, header = TRUE)
head(my_faithful)
```

| | eruptions | waiting |
|---|-----------|---------|
| 1 | 3.600 | 79 |
| 2 | 1.800 | 54 |
| 3 | 3.333 | 74 |
| 4 | 2.283 | 62 |
| 5 | 4.533 | 85 |
| 6 | 2.883 | 55 |

3.3 Dropbox

Vi kan också vilja läsa in data från dropbox på ett enkelt sätt. För att kunna komma åt data på dropbox behöver vi först dela den fil vi vill läsa in i R. Mer information om hur man delar filer med dropbox finns [\[här\]](#).

När vi delar en fil på dropbox för vi en länk. Nedan är ett exempel på en länk:

<https://www.dropbox.com/s/pbxm1lhmoax5zn6/google.csv?dl=0>

För att läsa in denna fil i R använder vi repmispaketet. Vill vi ha full kontroll på dropbox (skapa och läsa dolda filer m.m.) direkt från R finns paketet rDrop.

Länken från dropbox består av två delar, ett öppet dokumentid (pbxm1lhmoax5zn6 i exemplet ovan) och ett filnamn (google.csv). Med dessa två kan vi enkelt läsa in en fil i R på följande sätt:

```
library(repmis)
my_google <-
  repmis::source_DropboxData(file = "google.csv",
                             key = "pbxmlhmoax5zn6",
                             sep=";", dec = ",")
```

Downloading data from: <https://dl.dropboxusercontent.com/s/pbxmlhmoax5zn6/google.csv>
SHA-1 hash of the downloaded data file is:
45cc9b356f7121c149bdb405f2a70623d46add44

```
head(my_google)
```

| | Date | Open | High | Low | Close | Volume | Adj.Close |
|---|------------|--------|--------|--------|--------|----------|-----------|
| 1 | 2012-01-24 | 586.32 | 587.68 | 578.00 | 580.93 | 3055800 | 580.93 |
| 2 | 2012-01-23 | 586.00 | 588.66 | 583.16 | 585.52 | 3412900 | 585.52 |
| 3 | 2012-01-20 | 590.53 | 591.00 | 581.70 | 585.99 | 10576300 | 585.99 |
| 4 | 2012-01-19 | 640.99 | 640.99 | 631.46 | 639.57 | 6276500 | 639.57 |
| 5 | 2012-01-18 | 626.63 | 634.00 | 622.12 | 632.91 | 2761700 | 632.91 |
| 6 | 2012-01-17 | 631.98 | 631.98 | 625.68 | 628.58 | 1909300 | 628.58 |

Kapitel 4

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()**. Dock är denna labb mer utförlig än vignetten.

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

4.1 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 **interactive_pxweb()**. 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 är det tre delar vi behöver ange, vilket api (SCB i detta fall), version på api:et och api:ets språk. För att ta reda på dessa faktorer används funktionen **api_catalogue()** som då listar samtliga pxweb-apier som kommer förinstallerade.

```
api_catalogue()
```

2. För att anknyta till Statistiska centralbyråns api kan då följande kod användas:

```
mitt_data <- interactive_pxweb(api = 'api.scb.se', version = 'v1', lang = 'sv')
```

3. 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 1 för att gå vidare till menyn för statistik om arbetsmarknad:

```
...
21. [UF] Utbildning och forskning
Enter the data (number) you want to explore: ('esc' = Quit, 'b' = Back)
1: 1
```

4. Prova att “gå tillbaka” till huvudmenyn med **b**.
5. Navigera dig fram till följande datamaterial som vi ska läsa in:
[NAKUArblosaTK] Arbetslösa 15–74 år (AKU) efter arbetslöshetstidens längd, kön och ålder. Kvartal 2006K2 – 2014K2
6. Vi ska nu ladda ned detta datamaterial till vår R-session. R frågar nu om du vill ladda ned datamaterialet (eller om vi bara vill ha kod för att ladda ned materialet längre fram). Ange **y** (yes).
7. Nästa steg efterfrågas om du vill ha materialet i originalformat eller som en färdigformaterad data.frame i R. Ange **y** (yes).
8. Som ett sista steg efterfrågas nu om du vill ha koden för att ladda ned denna data direkt i R. Ange **y** (yes).
9. 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 **TABELLINNEHÅLL** ang **Arbetslösa**, **1000-tal** som den variabel du vill ha.
 - (e) Variabeln **KVARTAL** innehåller många kategorier. Om det är fler än 10 kategorier så döljs de mittersta kategorierna. För att se alla kategorier, ange **a**.
 - (f) Välj nu ut tidsperioden 2010K1 till senaste kvartalet hos SCB.
10. Nu ska materialet laddas ned. Du borde också få ut följande kod:

```
myDataSetName <-
get_pxweb_data(
url = "http://api.scb.se/OV0104/v1/doris/sv/ssd/AM/AM0401/AM0401L/NAKUArblosaTK",
dims = list(Arbetsloshetstid = c('*'),
Kon = c('1+2'),
Alder = c('15-24', '25-54', '55-64'),
ContentsCode = c('AM0401F0'),
Tid = c('2010K1', '2010K2', '2010K3', '2010K4', '2011K1', '2011K2', '2011K3', '2011K4', '2012K1',
clean = TRUE)
```

11. Titta på det material du laddat ned.

4.2 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 `myDataSetName`).
2. Prova att ändra argumentet `clean` till **FALSE**. Studera hur materialet ser ut. Detta är den “råa” data som laddas ned från SCB.
3. Vi vill nu ladda ned data för alla tidpunkter. Ändra variabeln `tid` till `'*'` och ladda ned datat på nytt.

4.3 * 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

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.

För att installera `markmyassignment` krävs paketet `devtools` (som därför först måste installeras):

```
> install.packages("devtools")
> devtools::install_github("MansMeg/markmyassignment")
```

För att automatiskt återkoppla en laboration behöver du först ange vilken laboration det rör sig om på följande sätt:

```
> library(markmyassignment)
> set_assignment("[assignment path]")
```

där `[assignment path]` är en adress du får av läraren till varje laboration.

För att se vilka uppgifter som finns i laborationen kan du använda funktionen `show_tasks()` på följande sätt:

```
> show_tasks()
```

För att få återkoppling på en uppgift använder du funktionen `mark_my_assignment()`. För att rätta samtliga uppgifter i en laboration gör du på följande sätt:

```
> mark_my_assignment()
```

Tänk på att uppgifterna som ska kontrolleras måste finnas som funktioner i R:s globala miljö. Du kan också kontrollera en eller flera enskilda uppgifter på följande sätt:

```
> mark_my_assignment(tasks="foo")
> mark_my_assignment(tasks=c("foo", "bar"))
```

Det går också att rätta en hel R-fil med samtliga laborationer. Detta är bra att göra innan du lämnar in din laboration. För att rätta en hel fil gör du på följande sätt:

```
> mark_my_assignment(mark_file = "[my search path to file]")
```

där `[my search path to file]` är sökvägen till din fil.

Obs! När hela filer kontrolleras måste den globala miljön vara tom. Använd `rm(list=ls())` för att rensa den globala miljön.

Kapitel 5

Inlämningsuppgifter

För att använda `markmyassignment` i denna laboration ange:

```
library(markmyassignment)

Loading required package: yaml
Loading required package: testthat
Loading required package: httr

lab_path <-
  "https://raw.githubusercontent.com/MansMeg/KursRprgm/master/Labs/Tests/d6.yml"
set_assignment(lab_path)

Assignment set:
D6 : Statistisk programmering med R: Lab 6
```

5.1 fast_swe_pop()

Du ska nu skapa en funktion, utan argument, som kopplar upp sig mot SCB:s api med `pxweb`, laddar ned befolkningsstatistik och summerar upp befolkningen i Sverige för varje år.

Tips! `aggregate()`

```
library(pxweb)
pop <- fast_swe_pop()
```

```
head(pop)

  year population
1 1968   7931193
2 1969   8004270
3 1970   8081142
4 1971   8115165
5 1972   8129129
6 1973   8144428
```

```
tail(pop)

  year population
41 2008   9256347
42 2009   9340682
43 2010   9415570
44 2011   9482855
```

| | | |
|----|------|---------|
| 45 | 2012 | 9555893 |
| 46 | 2013 | 9644864 |

5.2 gram_schmidt()

Vi ska nu göra en funktion för gram-schmidt ortogonalisering. Vi kommer implementera denna funktion på precis det sätt ni lärt er att göra Gram-Schmidt ortogonalisering. När vi gör det i datorn kan dock avrundningsfel göra att ortogonaliseringen blir numeriskt instabil. Mer information om gram-schmidt ortogonalisering finns [\[här\]](#).

För att implementera `gram_schmidt()` föreslås följande strategi:

1. Börja med att skapa funktionen `inner_prod()` som tar argumenten `x` och `y` och beräknar det inre produktrummet på följande sätt

$$\langle x, y \rangle = \sum_{k=1}^n x_k y_k$$

```
inner_prod(x = 1:4, y = -1:2)

[1] 10

inner_prod(x = 10:15, y = -(15:10))

[1] -920
```

2. Skapa sedan funktionen `proj()` som tar argumenten `u` och `v` och returnerar projektionsoperatoren:

$$\text{proj}_{\mathbf{u}}(\mathbf{v}) = \frac{\langle \mathbf{v}, \mathbf{u} \rangle}{\langle \mathbf{u}, \mathbf{u} \rangle} \mathbf{u}$$

```
proj(1:4, -1:2)

[1] -1.6667  0.0000  1.6667  3.3333

proj(10:15, -(15:10))

[1] 14.4503 13.4869 12.5236 11.5602 10.5969  9.6335
```

3. Skapa funktionen `vec_norm()` som tar argumentet `v` och returnerar en normaliserad konstant.

$$\|\mathbf{x}\| = \sqrt{\sum_{k=1}^n x_k^2}.$$

```
vec_norm(1:4)

[1] 5.4772

vec_norm(-1:2)

[1] 2.4495
```

4. Med dessa funktioner är det möjligt att sätta ihop funktionen `gram_schmidt(X)`. Funktionen ska kunna ta en godtycklig matris och returnera följande ortonormerade matris. Funktionen ska avbryta och returnera felmeddelandet "Not an object of class matrix." om `X` inte är klassen `matrix`.

$$\mathbf{E} = (\mathbf{e}_1 \quad \mathbf{e}_2 \quad \cdots \quad \mathbf{e}_n)$$

Nedan är två exempel på hur funktionen ska fungera.

```
X <- matrix(c(3,1,2,2), ncol=2)
gram_schmidt(X)

      [,1]      [,2]
[1,] 0.94868 -0.31623
[2,] 0.31623  0.94868

X <- matrix(c(1,-1,1,1,0,1,1,1,2), ncol=3)
gram_schmidt(X)

      [,1]      [,2]      [,3]
[1,]  0.57735  0.40825 -0.70711
[2,] -0.57735  0.81650  0.00000
[3,]  0.57735  0.40825  0.70711
```

5.3 Miniprojektet del I

En del av denna laboration är att genomföra miniprojektet del I. Se kurshemsidan för detaljer.