



SCRIPT OCH IMPORTERA DATA

Niklas Edvall

Script

För att slippa skriva saker för hand hela tiden är steg ett när man jobbar med R att skriva ett script för sin analys. I Rstudio finns ikonerna för att skapa nya filer uppe till vänster och där väljer man *R script* för att skapa ett nytt sådant.

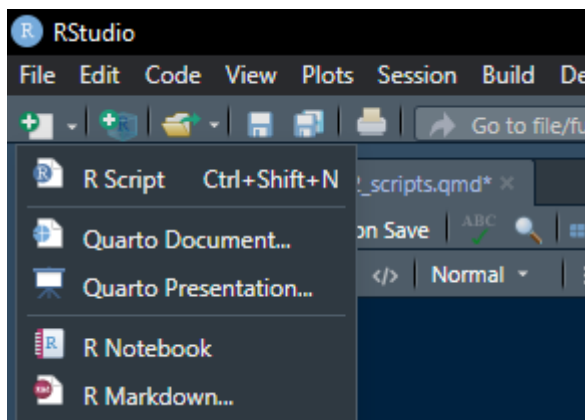


Figure 1: New File - R Script

I sitt script kan man nu skriva hela sin analys-pipeline från början till slut. Vill vi t.ex undersöka om det finns en statistisk signifikant skillnad mellan män och kvinnor att välja chokladglass före vanilj när de får frågan: *Skulle du hellre äta choklad- än vaniljglass?* Kan vi skriva följande analys-script.

Vi kodar kön som M eller F i variabeln `sex`, och om man svarade ja som Y eller nej som N i variabeln `ic.choco`

```
#Create variable for sex
sex <- c("M", "F", "M", "M", "F", "M", "F", "F", "M", "F", "F")

#Create variable for ice cream preference
ic.choco <- c("Y", "Y", "Y", "Y", "N", "N", "N", "N", "Y", "N", "N")

#Fishers exact test
fisher.test(sex, ic.choco)
```

Fisher's Exact Test for Count Data

```
data: sex and ic.choco
p-value = 0.08009
alternative hypothesis: true odds ratio is not equal to 1
95 percent confidence interval:
 0.6013099 1160.3870855
sample estimates:
odds ratio
 13.45266
```



Vi ser att testet resulterar i ett p-värde = 0.08, vilket betyder att vi inte kan påvisa statistiskt signifikant skillnad mellan grupperna på den klassiska 5%-nivån för statistisk signifikans.

Importera data och packages

Självklart är det galenskap att skriva in sin data manuellt med hjälp av funktionen `c()`

Istället vill vi importera data från en separat datafil som vi sparar. Det finns flera olika sätt att göra detta på, här använder vi ett paket som heter `readr` för att importera vår data som finns sparad i en csv-fil.

Paket? R har en massor av inbyggda funktioner men det finns oändligt många fler skräddarsydda funktioner att installera om man vill ha särskild funktionalitet. Dessa kommer i form av olika *packages* eller paket, och man kan ladda de man specifikt behöver för sitt aktuella script.

Första gången man vill använda ett paket måste man installera det med funktionen `install.packages()`, för att installera paketet `readr` anger vi `install.packages("readr")`.

När man vill använda ett paket i sitt script laddar man det med `library()`. Så, för att ladda vår data-fil med `readr` till en data frame vi kallar `dat` börjar vårt script med:

Fördelen med `readr` är att man också kan gå till *import dataset* och välja att importera data från en text-fil med `readr`. Då kan man även välja att t.ex exkludera vissa variabler eller ange olika format för variabler. Dialogrutan för att importera data skapar även en kod-snutt man klistra in i sitt script för att spara exakt parametrarna man använt för att importera data.

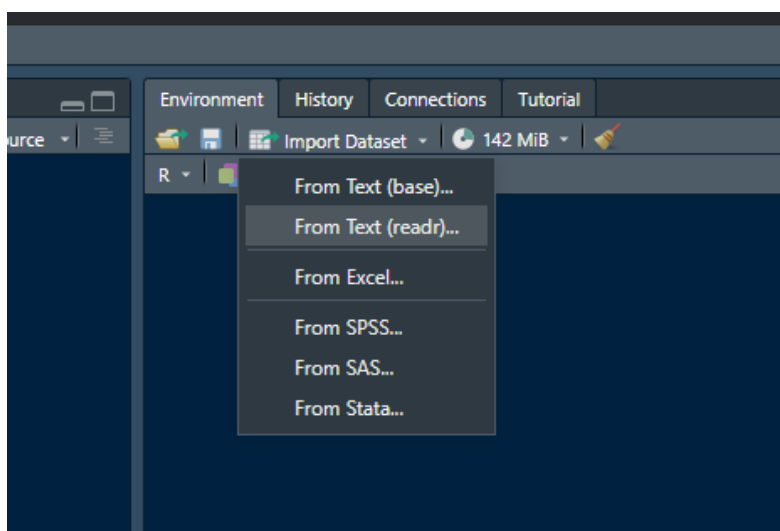


Figure 2: Import dataset

Vår exempel-data

Vår exempel-data innehåller ett unikt ID-nummer per deltagare, info om kön, ålder och självrapporterad hörselstatus. Hörtrösklar vid fyra frekvenser (0.5, 1, 2 och 4 kHz) per öra och de 25 frågor som återfinns på frågeformuläret Tinnitus Handicap Inventory (THI). Vi kommer använda datan för att se om vi kan besvara:

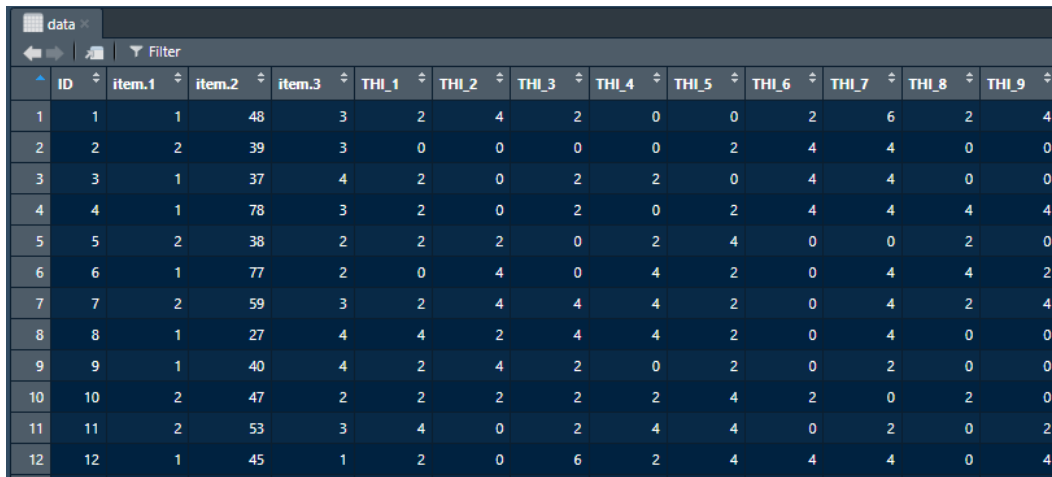
1. Är det skillnad mellan hur män och kvinnor besvarar THI?



2. Förändras hörseln med åldern?

Inspektera data

För att få en första överblick kan vi inspektera vår data. Enklast är att klicka på variabeln *dat* för vår data frame i det övre fönstret till höger i Rstudio.



	ID	item.1	item.2	item.3	THI_1	THI_2	THI_3	THI_4	THI_5	THI_6	THI_7	THI_8	THI_9
1	1	1	48	3	2	4	2	0	0	2	6	2	4
2	2	2	39	3	0	0	0	0	2	4	4	0	0
3	3	1	37	4	2	0	2	2	0	4	4	0	0
4	4	1	78	3	2	0	2	0	2	4	4	4	4
5	5	2	38	2	2	2	0	2	4	0	0	2	0
6	6	1	77	2	0	4	0	4	2	0	4	4	2
7	7	2	59	3	2	4	4	4	2	0	4	2	4
8	8	1	27	4	4	2	4	4	2	0	4	0	0
9	9	1	40	4	2	4	2	0	2	0	2	0	0
10	10	2	47	2	2	2	2	2	4	2	0	2	0
11	11	2	53	3	4	0	2	4	4	0	2	0	2
12	12	1	45	1	2	0	6	2	4	4	4	0	4

Figure 3: Vår data frame "dat"

Nedan ser vi med funktionen `dim()` att vår data har dimension 200x34, dvs 200 observationer (rader) för 37 olika variabler (kolumner). Funktionen `names()` returnerar namnen för alla kolumner i vår data. I funktionen `head()` specificerar vi att få tillbaka de första 3 raderna och 8 kolumnerna.

```
#Dimensions of dat frame  
dim(dat)
```

```
[1] 200 37
```

```
#Column names in dat frame  
names(dat)
```

```
[1] "ID"      "item.1" "item.2" "item.3" "THI_1"  "THI_2"  "THI_3"  "THI_4"  
[9] "THI_5"  "THI_6"  "THI_7"  "THI_8"  "THI_9"  "THI_10" "THI_11" "THI_12"  
[17] "THI_13" "THI_14" "THI_15" "THI_16" "THI_17" "THI_18" "THI_19" "THI_20"  
[25] "THI_21" "THI_22" "THI_23" "THI_24" "THI_25" "R500"   "R1000"  "R2000"  
[33] "R4000"  "L500"   "L1000"  "L2000"  "L4000"
```

```
#Look at first 3 rows and 8 columns of dat  
head(dat, c(3,8))
```

```
# A tibble: 3 x 8  
  ID item.1 item.2 item.3 THI_1 THI_2 THI_3 THI_4  
  <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>  
1     1     1    48     3     2     4     2     0  
2     2     2    39     3     0     0     0     0  
3     3     1    37     4     2     0     2     2
```



Notera att vi även ser vilken **typ** av variabel vi har i returen för `head()`, under kolumnens namn anges `<dbl>`, en förkortning för *double precision floating point* vilket betyder numerisk data. Vi behöver städa lite i vår data för att specificera vad som är numerisk data och vad som är kategorisk data.

Städa data

Kategoriska variabler

Kategoriska variabler i R kallas för faktor-variabler. Vi specificerar en faktor-variabel med funktionen `factor()`. Vill vi ange att variabeln *item.1* i vår data är en faktor är det dock ingen idé att bara ange:

```
#Make factor of item.1
factor(dat$item.1)
```

```
[1] 1 2 1 1 2 1 2 1 1 2 2 1 2 1 2 1 2 1 1 2 1 2 1 1 2 2 1 2 1 2 1 2 2 1 2 1 2
[38] 1 1 2 2 1 2 1 2 2 1 2 1 1 2 2 1 1 1 2 1 2 1 2 1 2 1 1 1 2 1 1 1 1 1 1 1
[75] 1 2 2 1 2 2 1 2 2 2 1 2 2 1 2 1 1 2 2 1 2 2 1 2 1 1 1 1 2 1 2 2 1 2 1 1
[112] 1 1 2 1 1 2 1 1 1 2 1 1 1 1 2 1 1 2 2 1 1 2 1 2 2 2 1 1 2 2 2 2 2 1 1 2 1
[149] 2 1 1 2 1 1 1 2 2 1 2 2 1 2 2 1 1 1 2 1 2 1 1 1 2 1 2 2 2 1 2 1 2 1 1 1 2
[186] 2 2 1 2 1 2 1 1 2 1 1 1 1 1 1 1
Levels: 1 2
```

Visserligen returneras *item.1* som en faktor, men vi måste spara den outputen till vår data frame genom att ange:

```
#Make factor and write to data frame
dat$item.1 <- factor(dat$item.1)
```

Vi råkar veta att *item.1* kodar för kön och kan göra detta tydligt genom att specificera ytterligare parametrar i funktionen `factor()` enligt nedan. Parametrarna `levels` (vilka kategorier variabeln innehåller) och `labels` (vad vi vill namnge dessa) åtskiljs med kommatecken.

Kom ihåg att ? `factor` alltid visar hjälpavsnittet för funktionen och ger exempel på hur den kan användas.

```
#Make factor, specify levels & labels, and write to data frame
dat$item.1 <- factor(dat$item.1,
                     levels = c(1,2),
                     labels = c("M", "F"))
```

Ändra namn på variabler i en data frame

Det vore även smidigt att ge kolumnen ett mer beskrivande namn. Minns tillbaka hur funktionen `names()` gav oss alla kolumn-namn ovan, vi kollar igen:

```
#Column names in data frame
names(dat)
```



```
[1] "ID"      "item.1" "item.2" "item.3" "THI_1"  "THI_2"  "THI_3"  "THI_4"
[9] "THI_5"   "THI_6"  "THI_7"  "THI_8"  "THI_9"  "THI_10" "THI_11" "THI_12"
[17] "THI_13"  "THI_14" "THI_15" "THI_16" "THI_17" "THI_18" "THI_19" "THI_20"
[25] "THI_21"  "THI_22" "THI_23" "THI_24" "THI_25" "R500"   "R1000"  "R2000"
[33] "R4000"   "L500"   "L1000"  "L2000"  "L4000"
```

Vi kan använda `names()` för döpa om variabler i en data frame. Vi ser att *item.1* är den andra variabeln i *data*, dvs har index 2. Men, det är strikt förbjudet att använda detta för att t.ex skriva: `names(dat)[2] <- "sex"`

Om kolumnen skulle få ett annat index, vilket ofta händer, kommer det sluta med att vi döper om någon annan kolumn och har förstört vår data totalt.

Istället hänvisar vi till ett dynamiskt index, `names(dat) == "item.1"`, som döper om alla kolumner med namnet *item.1* till *sex*

```
#Rename variable for sex
names(dat)[names(dat) == "item.1"] <- "sex"
```

På samma sätt definierar vi sen variabeln *item.3* som faktor-variabel och ändrar namn på den och variabel *item.2*

Skapa nya variabler

Vårt dataset innehåller hörtärsklar vid fyra frekvenser (0.5, 1, 2 och 4 kHz) per öra som vi kan använda för att räkna ut tonmedelvärde (*Pure Tone Average*; PTA4). Det är lätt att skapa/ange en ny variabel i vår data frame med `$` från medelvärdet av de fyra andra variablerna. Radbyte spelar ingen roll i scriptet, så länge alla symboler är på plats. Här skrivs varje referens till en variabel på egen rad för att det ska vara lättläsligt.

```
#Create variable for PTA4 Right
dat$PTA4.R <- (dat$R500 +
              dat$R1000 +
              dat$R2000 +
              dat$R4000) / 4

#Create variable for PTA4 Left
dat$PTA4.L <- (dat$L500 +
              dat$L1000 +
              dat$L2000 +
              dat$L4000) / 4
```

Vi har 25 variabler som heter *THI_1*, *THI_2*, *THI_3*.. osv. till *THI_25*. Dessa representerar svar på ett frågeformulär om tinnitus kallat *Tinnitus Handicap Inventory*. Svaren är redan kodade som antingen 0, 2 eller 4 vilket motsvarar de poäng man får för svarsalternativen. Maxpoäng är alltså $25 * 4 = 100$, och ju högre poäng desto mer besvärad är man av tinnitus.

Det vore intressant att skapa en ny variabel med varje försökspersons totala poäng på THI. Vi namnger denna variabel som *THIscore*. Det finns flera sätt att göra detta på. Enklast vore att helt enkelt summera de 25 variablerna:

```
#Sum all THI-variables to new variable THIscore
dat$THIscore <- dat$THI_1 + dat$THI_2 + dat$THI_3 ...
```



Det är funktionsdugligt, men det blir mycket att skriva och blir både svårläst, oflexibelt och ostabilt om t.ex en variabel får ett nytt namn eller kods på något annat sätt. Det finns alltid flera sätt att åstadkomma samma sak på med R, och vissa är smidigare än andra.

I det här fallet kan vi skapa en ny variabel som vi kallar *THI.names* med hjälp av funktionen `paste()` som klistrar ihop (eller 'konkatenerar') olika saker med en avskiljare (eller 'separator') som vi specificerar. Här konkatenerar vi bokstäverna "THI" med talen 1 till 25 avskiljda med ett understreck "_"

På så sätt innehåller då vår nya variabel *THI_names* namnen på alla de kolumner vi är intresserade av att summera.

```
#Create variable of column names relevant to THI
THI.names <- paste("THI", 1:25, sep="_")
```

Vi kan sen använda funktionen `apply()` för att applicera funktionen `sum()` på alla rader (rader specificerar vi med `MARGIN = 1`) i vår data frame som har ett namn som finns i variabeln *THI.names* och skriva resultatet till en ny variabel som vi igen kallar *THI.score*

```
#Calculate total THI score per subject
dat$THIscore <- apply(dat[,THI.names], MARGIN = 1, FUN = sum)
```

Analys

Vi kan nu besvara våra frågeställningar från början av dokumentet.

1. Är det skillnad mellan hur män och kvinnor besvarar THI?

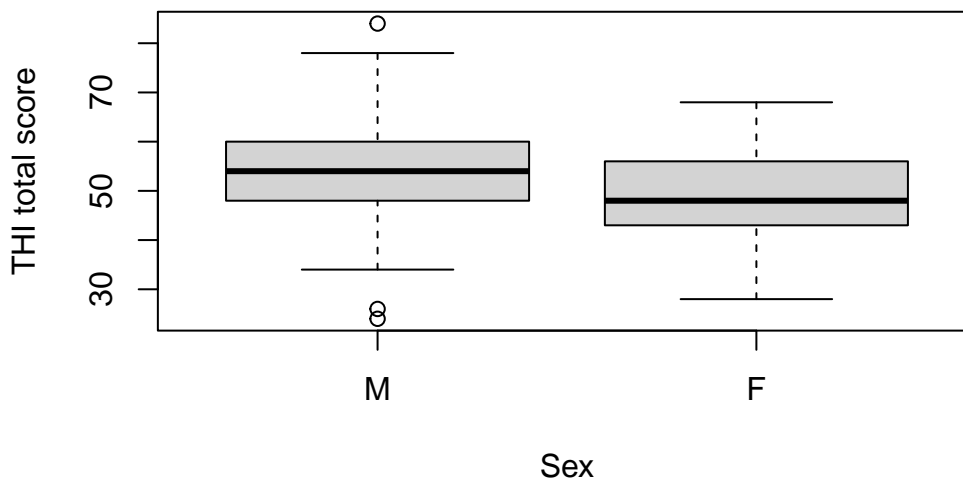
För att besvara denna fråga gör vi en box-plot och ett t-test mellan *THIscore* för män och kvinnor.

Vi börjar enkelt, men är du estetiskt lagd finns det massor av verktyg för att göra vackra figurer med R, se till exempel [galleriet för paketet ggplot2](#).

```
#Box plot of THI score for sex
plot(dat$sex, dat$THIscore,
     main = "THI score for Male (M) and Female (F)",
     xlab = "Sex",
     ylab = "THI total score")
```



THI score for Male (M) and Female (F)



```
#Create subset of women and men separately
F.THI <- dat[dat$sex == "F", "THIScore"]
M.THI <- dat[dat$sex == "M", "THIScore"]

#Perform t-test
t.test(F.THI, M.THI)
```

Welch Two Sample t-test

```
data: F.THI and M.THI
t = -3.9908, df = 193.44, p-value = 9.34e-05
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -7.951339 -2.691518
sample estimates:
mean of x mean of y
 49.25000  54.57143
```

Vår plot verkar indikera att män fått lite högre poäng på THI och vårt testresultat: **t = -3.99, p < 0.01** påvisar en statistiskt signifikant skillnad mellan grupperna.

Ovan skapade vi en ny variabel för män respektive kvinnors THIScore som vi sedan använde i funktionen `t.test()`. Vi hade också direkt kunnat skriva `t.test(dat[dat$sex == "F", "THIScore"], dat[dat$sex == "M", "THIScore"])` och fått samma resultat, men då är det svårt att utläsa vad som händer. Ett tredje alternativ är att använda paketet `arsenal` för att direkt ställa upp en deskriptiv tabell med funktionen `tableby()`.

Kom ihåg att om vi inte har paketet installerat kör vi först `install.packages("arsenal")` och sedan `library(arsenal)` för att aktivera det i sessionen.

```
#Read in necessary library
library(arsenal)
```



```
#Create table using tableby()
tab1 <- tableby(sex ~ THIScore, data = dat,
               digits = 1,
               total = FALSE)

#Print table
summary(tab1, text = TRUE)
```

	M (N=112)	F (N=88)	p value
THIScore			< 0.001
- Mean (SD)	54.6 (9.8)	49.2 (9.0)	
- Range	24.0 - 84.0	28.0 - 68.0	

På det här sättet får vi dessutom antal observationer, medelvärde och standardavvikelse utskrivet. Vi läser hjälpavsnittet genom att skriva `? tableby` i konsollen och ser att `digits = 1` kan användas för att specificera antalet decimaler, `total = FALSE` döljer kolumnen för total (män + kvinnor), och att p-värdet som rapporteras är: *equivalent to two-samples t-test*.

2. Förändras hörseln med åldern?

För att besvara den här frågan gör vi en enkel linjär regression av hur ökad ålder påverkar tonmedelvärdet PTA4 för höger och vänster öra separat. För linjär regression använder vi funktionen `lm()` och spar resultatet (vår regressionsmodell) som en ny variabel `linear.reg.L` när vi använder PTA för vänster öra och `linear.reg.R` när vi använder PTA för höger öra.

```
#Save linear regression model to variable
linear.reg.L <- lm(data = dat, PTA4.L ~ age)
linear.reg.R <- lm(data = dat, PTA4.R ~ age)

#Print summary of linear regression model
summary(linear.reg.L)
```

Call:

```
lm(formula = PTA4.L ~ age, data = dat)
```

Residuals:

Min	1Q	Median	3Q	Max
-18.3107	-6.1398	-0.1393	5.3527	20.5579

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	12.32806	1.69318	7.281	7.64e-12 ***
age	0.15301	0.03216	4.758	3.76e-06 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 7.952 on 198 degrees of freedom

Multiple R-squared: 0.1026, Adjusted R-squared: 0.09808

F-statistic: 22.64 on 1 and 198 DF, p-value: 3.758e-06



```
summary(linear.reg.R)
```

Call:

```
lm(formula = PTA4.R ~ age, data = dat)
```

Residuals:

Min	1Q	Median	3Q	Max
-20.1685	-4.8667	-0.4338	5.3292	20.4885

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	13.68233	1.64973	8.294	1.66e-14 ***
age	0.12957	0.03133	4.135	5.23e-05 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 7.748 on 198 degrees of freedom

Multiple R-squared: 0.07951, Adjusted R-squared: 0.07486

F-statistic: 17.1 on 1 and 198 DF, p-value: 5.232e-05

Genom att använda `summary()` för regressionsmodellen vi skapat skrivs en överblick av resultatet ut i konsollen. Vi ser att koefficienten `age`, för vänster öra, är 0.15 med ett p-värde långt under 0.001 (3.76e-06 är 0.00000376). Det betyder att det finns ett signifikant samband mellan ålder och tonmedelvärde där vår modell påvisar att PTA4 ökar med 0.15 dB per levnadsår. Förhållandet är signifikant även för höger öra (PTA4.R) med koefficienten 0.13.

Vi känner oss redo för att använda lite mer avancerade plot-alternativ. Vi specificerar att vi vill ha två plots bredvid varandra (med `par()`) och skapar en scatter-plot för vänster och höger öras PTA som en funktion av ålder. Funktionen `abline()` låter oss dessutom lägga till linjen från vår linjära regressionsmodell på respektive plot.

```
#Set plot space to two columns
par(mfrow=c(1,2))

#Scatter plot of PTA LEFT as function of age
plot(dat$age, dat$PTA4.L,
     pch = 20,
     main = "PTA Left ear vs age",
     xlab = "Age (years)",
     ylab = "PTA Left (dB HL)")

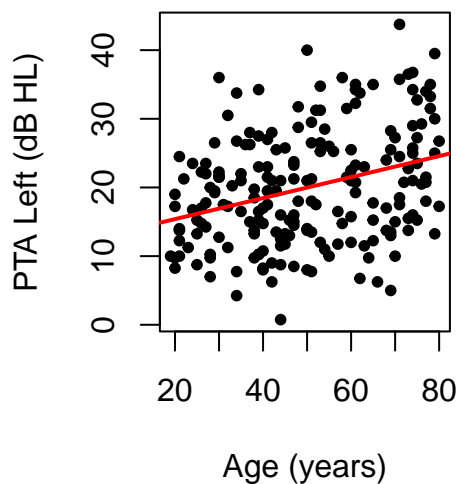
#Best fit linear regression LEFT
abline(linear.reg.L, lwd = 2, col = "red")

#Scatter plot of PTA RIGHT as function of age
plot(dat$age, dat$PTA4.R,
     pch = 20,
     main = "PTA Right ear vs age",
     xlab = "Age (years)",
     ylab = "PTA Right (dB HL)")
```

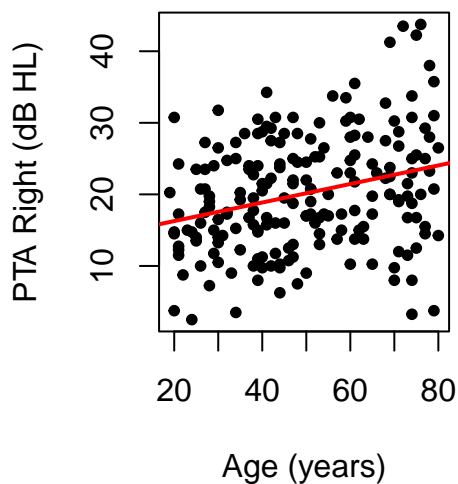


```
#Best fit linear regression RIGHT  
abline(linear.reg.R, lwd = 2, col = "red")
```

PTA Left ear vs age



PTA Right ear vs age



Sammanfattning

För att gå från rådata till figurer och statistisk analys som besvarar våra två frågeställningar behöver vi totalt runt 60 rader kod, inklusive kommentarer, i vårt script (nedan).

```
#Load libraries  
library(readr)  
library(arsenal)  
  
#Load data  
dat <- read_csv("data.csv")  
  
#Specify factor variables  
dat$item.1 <- factor(dat$item.1,  
                     levels = c(1,2),  
                     labels = c("M", "F"))  
  
dat$item.3 <- factor(dat$item.3,  
                     levels = c(1,2,3,4),  
                     labels = c("No", "Sometimes",  
                                "Often", "Always"))  
  
#Rename variables  
names(dat)[names(dat) == "item.1"] <- "sex"  
names(dat)[names(dat) == "item.2"] <- "age"  
names(dat)[names(dat) == "item.3"] <- "hearing"  
  
#Create variables for PTA  
dat$PTA4.R <- (dat$R500 + dat$R1000 + dat$R2000 + dat$R4000) / 4
```



```
dat$PTA4.L <- (dat$L500 + dat$L1000 + dat$L2000 + dat$L4000) / 4

#Create variable of column names relevant to THI
THI.names <- paste("THI", 1:25, sep="_")

#Calculate total THI score per subject
dat$THIScore <- apply(dat[,THI.names], MARGIN = 1, FUN = sum)

#Create table using tableby()
tab1 <- tableby(sex ~ THIScore, data = dat,
                digits = 1,
                total = FALSE)

#Print table
summary(tab1, text = TRUE)

#Box plot of THI score for sex
plot(dat$sex, dat$THIScore,
     main = "THI score for Male (M) and Female (F)",
     xlab = "Sex",
     ylab = "THI total score")

#Save linear regression model to variable
linear.reg.L <- lm(data = dat, PTA4.L ~ age)
linear.reg.R <- lm(data = dat, PTA4.R ~ age)

#Print summary of linear regression model
summary(linear.reg.L)
summary(linear.reg.R)

#Set plot space to two columns
par(mfrow=c(1,2))

#Scatter plot of PTA LEFT as function of age
plot(dat$age, dat$PTA4.L,
     pch = 20,
     main = "PTA Left ear vs age",
     xlab = "Age (years)",
     ylab = "PTA Left (dB HL)")
#Best fit linear regression LEFT
abline(linear.reg.L, lwd = 2, col = "red")

#Scatter plot of PTA RIGHT as function of age
plot(dat$age, dat$PTA4.R,
     pch = 20,
     main = "PTA Right ear vs age",
     xlab = "Age (years)",
     ylab = "PTA Right (dB HL)")
#Best fit linear regression RIGHT
abline(linear.reg.R, lwd = 2, col = "red")
```