

# Datalaboration 1 - 732G50

VT2021

## Introduktion

Denna datorövning behandlar visualisering av olika typer av variabler.

Denna laboration kommer att göras i R-studio.

Öppna ett nytt R-script **File** → **New file** → **R-script**.

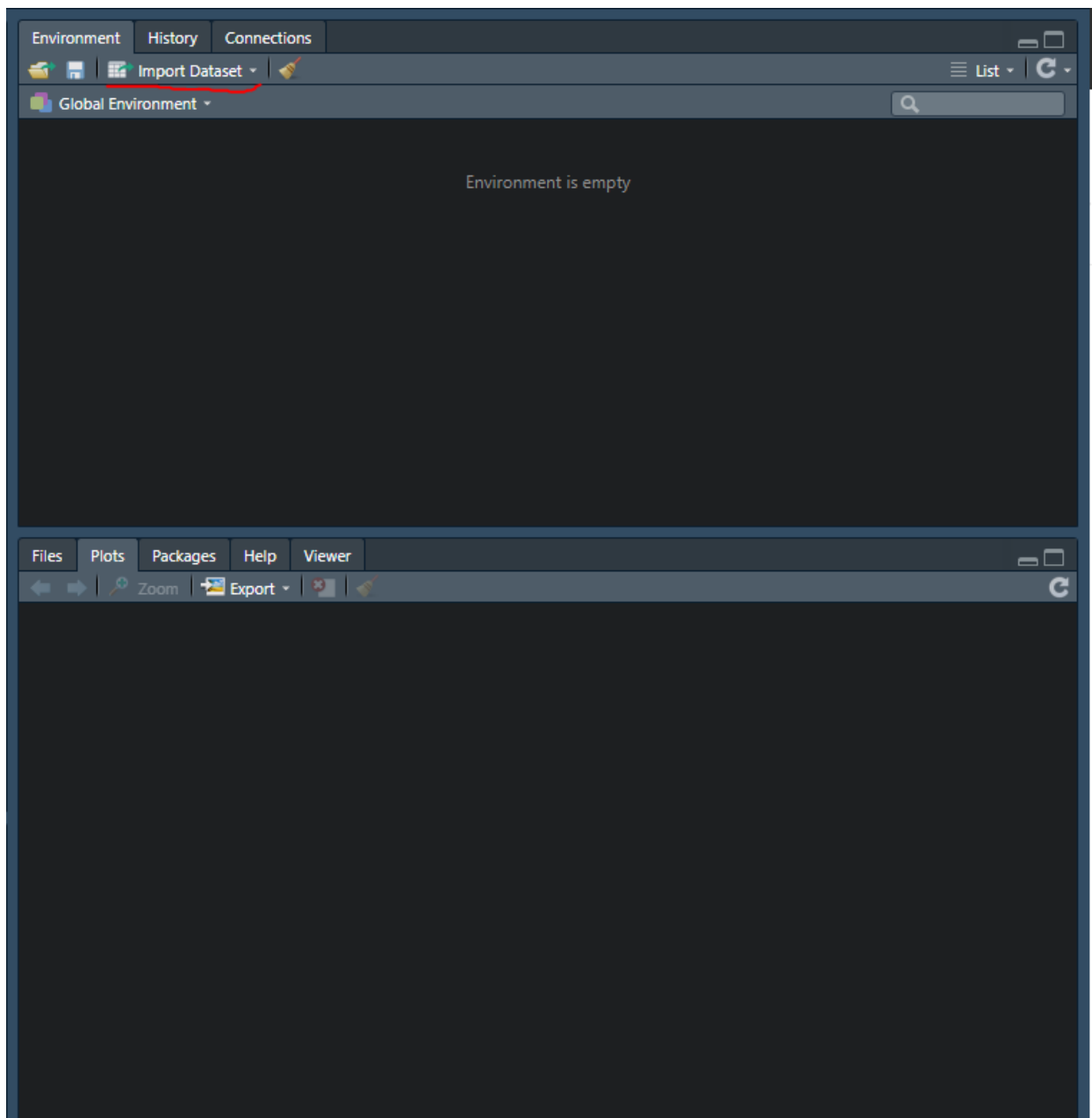
För att köra koden markerar du den koden du vill köra och trycker på **Run** uppe i högra hörnet. Du kan även trycka **Ctrl+Enter** för att köra den markerade koden.

## Ladda ned R-studio

- **R-studio**

## Uppgift 1

Börja med att ladda ner det datamaterial som finns på kursens LISAM-sida. Datamaterialet innehåller tre flikar med information. Den första fliken behandlar data om längd och vikt hos 100 män och 100 kvinnor. Den andra fliken behandlar information om civilstånd och bilmärken på Påhittade gatan. Den tredje fliken behandlar information om tentaresultat och ålder hos ett antal påhittade studenter.



För att importera datamaterialet till R-studio tryck på **Import Dataset** → **From Excel...** välj sedan sheet **Langd och vikt**. Välj även namnet `langd_vikt` i Name:

Import Excel Data

File/URL:

Data Preview:

Längd (double)	Vikt (double)	Kön (character)
185.4	78.9	Man
186.1	78.1	Man
172.0	65.4	Man
177.2	63.8	Man
174.1	58.3	Man
182.9	73.7	Man
173.9	61.5	Man
192.5	83.0	Man
182.6	66.1	Man
178.3	64.1	Man
182.5	68.9	Man
175.9	61.5	Man
189.6	89.2	Man
174.3	63.1	Man
184.9	65.9	Man
180.4	79.2	Man
182.3	67.8	Man
191.8	89.9	Man
177.4	66.5	Man
169.3	67.4	Man
181.3	70.2	Man

Previewing first 50 entries.

Import Options:

Name:  Max Rows:  ☒ First Row as Names

Sheet:  Skip:  ☒ Open Data Viewer

Range:  NA:

Code Preview:

```
library(readxl)
langd_vikt <- read_excel("c:/Users/Harje/Downloads/data1.xlsx",
  sheet = "Langd och vikt")
view(langd_vikt)
```

[Reading Excel files using readxl](#)

Du kommer sedan att se att du infört datamaterialet till höger i Global Environment.

## 1

Identifiera  $\bar{x}$ , s samt kvartilerna för längd samt vikt.

Dollar tecknet \$ skrivs efter din data.frame så väljer du sedan vilken variabel du är intresserad av att skriva ut, exempelvis Längd som nedan.

**Exemplet som visas har jag namnet på data.frame till langd\_vikt**

```
langd_vikt$Längd  
head(langd_vikt) # Skriver ut de första 5 observationerna.
```

Ni kommer använda följande funktioner som är inbyggda i R:

- `summary()` skriver ut kvartilerna min och max.
- `sd()` skriver ut standardavvikelsen.
- `mean()` skriver ut medelvärde.
- `median()` skriver ut medianen.

```
summary(langd_vikt$Längd) # summary() kvartilerna  
sd(langd_vikt$Längd) # sd() standardavvikelsen  
mean(langd_vikt$Längd) # mean() medelvärde  
median(langd_vikt$Längd) # median() medianen
```

Genom att byta ut **Längd** till **Vikt**

Gruppera datamaterialet efter kön.

```
data_kvinna <- langd_vikt[which(langd_vikt$Kön == "Kvinna"),]  
data_man <- langd_vikt[which(langd_vikt$Kön == "Man"),]
```

```
summary(data_kvinna$Längd) # summary() kvartilerna  
sd(data_kvinna$Längd) # sd() standardavvikelsen  
mean(data_kvinna$Längd) # mean() medelvärde  
median(data_kvinna$Längd) # median() medianen
```

Gör nu samma för männen genom att skriva `data_man$Längd` istället för `data_kvinna$Längd`. Är det någon skillnad mellan man och kvinna?

## 2

Nästa steg är att undersöka fördelningen av Längd och Vikt med histogram. Genom att använda funktionen `hist()` som är inbyggd i R.

`main = "text"` bestämmer titlen för ploten är.

`ylab = "text"` bestämmer vad som ska stå på y-axel

`xlab = "text"` bestämmer vad som ska stå på x-axel

`col = "red", col = 2, col = c(1:8)` sätter olika färger på ploten.

```
hist(langd_vikt$Längd, main= "Histogram över längd båda könen")
```

```
hist(data_kvinna$Längd, main= "Histogram över längd för Män")
```

```
hist(data_man$Längd, main= "Histogram över längd för Kvinna")
```

Vad kan vi säga om fördelningen av de olika variablerna? Skillnader mellan könen?

## Uppgift 2

Vi ska nu gå till bladet som heter Bilmärke och fortsätta med uppgifterna. Importera Bilarken från samma datamaterial från tidigare uppgift genom att välja sheet **Bilmarke**.

Exemplet som visas har jag namnet på `data.frame` till **Bilmarke**

1

Skapa en frekvenstabell över variablerna Bil och Civilstånd med hjälp av funktionen `table()`.

```
table(Bilmarke$Bil) # table() skriver ut en frekvenstabell av vald variabel
```

```
table(Bilmarke$Civilstånd, Bilmarke$Bil) # Du kan även skriva in fler variabler in i funktionen
table(Bilmarke$Bil, Bilmarke$Civilstånd) # Du kan välja den ordningen som passar dig bäst.
```

Vilket bilmärke är den vanligaste? Hur stor andel av bilarna är fordar? Hur många har civilstånd par?

2.

Att presentera tabeller är inte alltid det enklaste sättet för läsaren att få en bild av data. Det är oftast enklare (och snyggare) att skapa ett diagram.

Här kommer vi måsta kombinera två olika funktioner `table()` samt `barplot()` testa gärna att experimentera med olika färger

```
barplot(table(Bilmarke$Bil))
```

```
barplot(table(Bilmarke$Bil), col=1:8) # col="red", col=8:1
```

```
barplot(table(Bilmarke$Civilstånd), col=1:2) # Två grupper blir det två olika färger.
```

```
Bil_par <- Bilmarke[which(Bilmarke$Civilstånd == "Par"),]
Bil_singel <- Bilmarke[which(Bilmarke$Civilstånd != "Par"),]
```

```
barplot(table(Bil_par$Bil),main="Bilmärken för par")
barplot(table(Bil_singel$Bil),main="Bilmärken för singlar")
```

Om ni vill ha dem i samma bild kan ni sätta genomskinlighet på detta sätt.

```
blue <- rgb(0, 0, 1, alpha=0.2)
red <- rgb(1, 0, 0, alpha=0.2)
barplot(table(bilmarke$Civilstånd,bilmarke$Bil),main="Bilmärken för alla kön", col=c(blue,red))
legend("topleft", c("Par", "Singel"), col=c(blue,red), pch=19)
```

## Uppgift 3

Läs nu in det sista bladet som heter Elevålder som ni gjort tidigare med att välja sheet **Elevvalder**.

Exemplet som visas har jag namnet på `data.frame` till **Elevvalder**

### 1

Börja med att visualisera fördelningen av dessa två variabler. Här kan vi antingen välja ett histogram eller stolpdigram på Ålder eftersom man kan tänka sig att värdena kan anta decimalvärden.

Svara på följande frågor

- (a) Hur stor andel fick full pott, alltså 20 poäng, på tentan?
- (b) Vilka åldrar är mest frekventa?

```
barplot(table(Elevvalder$Ålder),xlab = "Ålder",main="Stolpdigram över åldrar", ylab="Antal")
```

```
barplot(table(Elevvalder$Resultat),xlab = "Resultat",main="Stolpdigram över resultat", ylab="Antal")
```

#### a)

Detta kan svaras på att summera hur många som fick resultatet 20.

`sum()` summerar en vektor. Om du skriver `Elevvalder$Resultat == 20` så kommer det upp två olika värden `TRUE` och `FALSE` detta värde kallas `Boolean`, sant eller falskt. Med funktionen `sum()` kan man summera alla sanna värden.

```
sum(Elevvalder$Resultat == 20)
```

#### b)

Du kan antingen skriva ut en tabell som vi gjorde i 2.1 eller kolla stolpdigrammet.

### 2

R har inga grundfunktioner för dotplot så i stället testar vi att skapa boxplot och försök att dra slutsatser för ålder och resultat. Har vi några `outliers` dvs några observationer som sticker ut?

```
boxplot(Elevvalder$Ålder, col="lightblue", main="Boxplot över ålder",  
        ylab="Ålder")  
boxplot(Elevvalder$Resultat, col="lightblue", main="Boxplot över resultat",  
        ylab="Resultat")
```

### 3

Ibland kan man vilja dela upp en kvantitativ variabel, t.ex. ålder, i olika klassindelningar. Detta kallas med ett annat ord diskretisering.

Vi ska nu gruppera åldrarna till 6 olika åldersgrupper.

0-19, 20-24, 25-29, 30-34, 35-39, 40-

I R-studio använder vi följande kommandon för att göra detta.

- $A == B$  betyder  $A = B$ , A är lika med B.
- $A < B$  betyder  $A < B$ , A är mindre än B.
- $A <= B$  betyder  $A \leq B$ , A är mindre eller lika med B.
- $A >= B$  betyder  $A \geq B$ , A är större eller lika med B.
- $A != B$  betyder  $A \neq B$ , A är inte lika med B.

I detta exempel räcker det med att använda  $<=$  och  $>=$  tillsammans med funktionen `which()`.

Börja med att skapa en tom variabel NA står för Not Applicable eller Not Available, kan även tolkas som Null.

```
Eleva1der$Alderklass <- NA  
head(Eleva1der) # Skriver automatiskt ut de första 5 observationerna i datamaterialet.
```

```
Eleva1der$Alderklass[which(Eleva1der$Ålder <= 19)] <- "0-19"  
# Alla under 20 kommer ha gruppsnamnet 0-19
```

```
Eleva1der$Alderklass[which(Eleva1der$Ålder >= 20 & Eleva1der$Ålder <= 24)] <- "20-24"  
# Alla med åldern mellan 20 och 24 kommer få gruppsnamnet 20-24  
Eleva1der$Alderklass[which(Eleva1der$Ålder >= 25 & Eleva1der$Ålder <= 29)] <- "25-29"  
# Alla med åldern mellan 25 och 29 kommer få gruppsnamnet 25-29  
Eleva1der$Alderklass[which(Eleva1der$Ålder >= 30 & Eleva1der$Ålder <= 34)] <- "30-34"  
# Alla med åldern mellan 30 och 34 kommer få gruppsnamnet 30-34
```

```
Eleva1der$Alderklass[which(Eleva1der$Ålder >= 35 & Eleva1der$Ålder <= 39)] <- "35-39"  
# Alla med åldern mellan 35 och 39 kommer få gruppsnamnet 35-39
```

```
Eleva1der$Alderklass[which(Eleva1der$Ålder >= 40)] <- "40-"  
# Alla med åldern 40 och uppåt kommer få gruppsnamnet 40-
```

```
barplot(table(Eleva1der$Alderklass),main="Stolpsdiagram åldersklasser",xlab="Ålder",ylab="Antal")
```

Saknas det en grupp i stolpsdiagrammet och i så fall varför?