

Introduktion til R III

12. maj 2023

Kristian G. Kjellmann (kgk@socsci.aau.dk)
&
Rolf L. Lund (rolfll@socsci.aau.dk)

Institut for Sociologi og Socialt Arbejde



AALBORG UNIVERSITY
DENMARK

Dagens program



1. Opsamling på datahåndtering
2. Eksport af statistiske modeller (med stargazer)
3. Opsætning af modeller i tekstpublikationer
4. Øvelser

I kan lave publicérbare modeller med R

I kan løse forskelligartede datahåndteringsudfordringer i det samme datasæt

Har I styr på jeres datahåndteringsbegreber?

Subsetting

Filtrering

Variabelændringer

Rekodning

Missingværdier

Der er to dele i at få resultaterne af en statistisk model frem i R:

1. Specificér modellen (fx med en funktion som `lm()` for lineære modeller)
2. Få koefficienter og resultater fra modellen (med brug af funktionen `summary()`)

Tre ting krævet for at specificere en model:

1. Et datasæt
2. Funktion for typen af model, man vil lave
3. Formel, der specificerer det sammenhæng, som man vil modellere

Al datahåndtering i datasæt skal ske *inden* man laver modellen.

Man specificerer en formel med R's formelsyntax, fx $y \sim x_1 + x_2 + x_3$.

Eksempel:

```
grsp_model <- lm(grspnum ~ eduyrs + wkhtot, data = ess18)
```

En model i R er blot endnu en type objekt.

For at se resultatet af modellen, skal man derfor spørge R korrekt om det.

Ved blot at kalde modellen får man begrænsede resultater:

```
grsp_model
```

Call:

```
lm(formula = grspnum ~ eduyrs + wkhtot, data = ess18)
```

Coefficients:

(Intercept)	eduyrs	wkhtot
-3668.7	891.1	888.4

Ved at bruge `summary()` funktionen gives de relevante resultater:

```
summary(grsp_model)
```

```
Call:
lm(formula = grspnum ~ eduyrs + wkhtot, data = ess18)

Residuals:
    Min       1Q   Median       3Q      Max
-64614  -16639  -10324   -1448  295649

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  -3668.7    20645.3   -0.178  0.8590
eduyrs         891.1     856.4    1.041  0.2984
wkhtot        888.4     458.2    1.939  0.0529 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 121600 on 707 degrees of freedom
(575 observations deleted due to missingness)
Multiple R-squared:  0.007698, Adjusted R-squared:  0.004891
F-statistic: 2.742 on 2 and 707 DF,  p-value: 0.06511
```


Som standard vil R behandle tekstvariable i statistiske modeller som *unordered factors*; altså nominalt skalerede variable.

R vil desuden tage kategorien, der kommer først i alfabetisk rækkefølge, som referencekategori.

Dette kan lede til uhensigtsmæssige resultater. Man bør derfor altid tage aktiv stilling til, hvordan den kategoriske variabel skal behandles, inden man laver modellen.

Overvej følgende:

Skal variablen behandles som *ordinal* eller *nominal*? (ordered eller unordered)

Skal variablen behandles som *interval skaleret*? (konvertér til numerisk)

Hvis variablen skal behandles som nominal, hvilken kategori skal så være *referencekategorien*? (kan fx ændres med `relevel` i `mutate` funktionen)

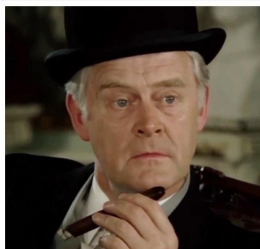
HUSK: Der er forskel på hvad variablen *er*, og hvordan vi behandler den i en model!

Eksport af statistiske modeller



Call:	
lm(formula = grspnum ~ eduyrs + wkhtot, data = ess18)	
Residuals:	
Min	1Q Median 3Q Max
-64614	-16639 -10324 -1448 2955649
Coefficients:	
Estimate Std. Error t value Pr(> t)	
(Intercept)	-3668.7 20645.3 -0.178 0.8590
eduyrs	891.1 856.4 1.041 0.2984
wkhtot	888.4 458.2 1.939 0.0529

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1	
Residual standard error: 121600 on 707 degrees of freedom	
(575 observations deleted due to missingness)	
Multiple R-squared: 0.007696, Adjusted R-squared: 0.004891	
F-statistic: 2.742 on 2 and 707 DF, p-value: 0.06511	



Effekt af års uddannelse og arbejdstid på løn

Antal års uddannelse	891,15
	(856,41)
Arbejdstimer om ugen	888,39
	(458,22)
Constant	-3.668,75
	(20.645,31)
Observations	
710	
R ²	
0,01	
Adjusted R ²	
0,005	
Residual Std. Error	
121.581,20 (df = 707)	
F Statistic	
2,74 (df = 2; 707)	

Note: *p<0,05; **p<0,01; ***p<0,001

Outputtet af en model i R konsollen er ikke kønt og egner sig dårligt til at fremstille i en publikation (rapport, artikel eller andet).

Der findes forskellige pakker til at lave pæne outputs af modeller fra R. En god pakke til dette er pakken `stargazer`.

`stargazer` tillader bl.a. at eksportere output fra en statistisk model til en HTML-fil. Indholdet af en HTML-fil kan kopieres direkte over i programmer som Microsoft Word.

`stargazer` har ufatteligt mange tilpasningsmuligheder.

For at danne en output-fil af en model med `stargazer` kræves som minimum:

- En model (det objekt hvor modellen er gemt - ikke outputtet af `summary()`)

- En filtype (`stargazer` kan danne forskellige filtyper, som kan styres med argumentet `type`. `type = "HTML"` danner et HTML-output)

- Et filnavn, som outputtet skal gemmes til (dette gøres med argumentet `out = <filnavn>`)

OBS! Output-filer fra `stargazer` gemmes i din arbejdssti (`getwd()`)

Eksempel på brug af stargazer

```
stargazer(  
  grsp_model,  
  type = "html",  
  out = "grsp_model2.html",  
  star.cutoffs = c(0.05, 0.01, 0.001),  
  decimal.mark = ",",  
  digit.separator = ".",  
  covariate.labels = c("Antal års uddannelse",  
                       "Arbejdstimer om ugen"),  
  dep.var.labels.include = FALSE,  
  dep.var.caption = "",  
  digits = 2,  
  title = "Effekt af års uddannelse og arbejdstid på løn"  
)
```

Eksport af statistiske modeller



Argument	Forklaring
<code>type</code>	Hvilken filtype skal output gemmes som? ("html" anbefales)
<code>out</code>	Hvad skal filen hedde? (Husk at ende med ".html", hvis <code>type = "html"</code> - vær desuden opmærksom på arbejdssti (tjek med <code>getwd()</code>)
<code>star.cutoffs</code>	Skæringsværdier for p-værdier til stjernnotation. Specificeres som vector med tre tal, svarende til værdi for hhv. *, **, *** (fx <code>c(0.05, 0.01, 0.001)</code>)
<code>decimal.mark</code>	Sæt hvilket tegn, der skal adskille decimaler
<code>digit.separator</code>	Sæt hvilket tegn, der skal adskille tusinde
<code>digits</code>	Bestem antal decimaler
<code>covariate.labels</code>	Ændrer mærkater for uafhængige variable. Skrives som en vector (<code>c()</code>) i samme rækkefølge som i modellen
<code>dep.var.labels.include</code>	Bestem hvorvidt mærkat for afhængig variabel skal inkluderes (logisk værdi, TRUE/FALSE)
<code>dep.var.caption</code>	Bestem overskrift for afhængig variabel (udelad ved at angive tom tekstværdi: "")
<code>title</code>	Giv output en overskrift