

02402 - Statistics Projekt 2: Bmi undersøgelse

Daniel F. Hauge - s2011687

- **a) Descriptiv analyse**
 - Fördelning
 - Nøgle tal
 - Plots
 - **b) Multipel lineær regressionsmodel**
 - Model parametre
 - Model kontrol
 - Alders koefficient konfidensinterval
 - Hypotesetest
 - **g) Backward selection**
 - Correlation
 - Confidensintervaller
 - MLE Summary
 - Slut model
 - **h) Prædiktioner**
 - Prediction
 - Confidence
 - Vurdering

Dette projekt består af en statistisk analyse af et dataset med mennesker. Projektet forsøger at kaste lys på BMI med statistik. Projektet henvender sig til læsere som er familier med projekt beskrivelsen fra statistik kursen 02402 fra DTU og befander sig komfortabelt i statistiske begreber og metoder.

Projektet er lavet som en R-notebook, og indeholder derfor oversider hvor R er brugt som redskab til udregning og plot-opsætninger. Boken med koden indekser at der er blevet kodet noget i R til at udregne, tegne eller gemme værdier til senere brug. R koden kan findes på følgende måder:

- Se den medfølgende "R" fil, eller via på:
- Se den medfølgende "note" fil, eller på:
- **g) Alders koefficient konfidensinterval**
- Projektet kan også findes i renderet form med kode afnet på: <https://htmlpreview.github.io/?https://github.com/DanielHauge/02402-statistics/blob/master/Project1/ProjectLab.html>

a) Descriptiv analyse

Data)

Datamaterialet indeholder 847 observationer med 5 egenskaber på mennesker for projektets problemstilling. Variable på hver observation er som følger:

Variable Navn	Måle type	Måle enhed	Forklaring
id	N/A	Heltal	Et unikt tal der kan bruges som identifikation for observationen.
age	Kvantitativt	År	Personens alder målt i år.
fastfood	Kvantitativt	dage pr. år	Antal dage per år personen har spist fastfood.
bmi	Kvantitativt	Bmi	Dette er en enhed som prøver at beskrive kroppens størrelse. Formel: $\text{bmi} = \frac{\text{vægt}}{\text{højde}^2}$. BMI står for "Body Mass Index".
logbmi	Kvantitativt	Log(Bmi)	Dette er bmi'en der er blevet log-transformeret.

Nedenfor ses første og sidste observation som et præsentierende eksempel:

id	bmi	age	fastfood	logbmi
<id>	<bmi>	<age>	<fastfood>	<logbmi>
1	21.2963	44	0.0	3.054533
847	21.2963	24	78.2	3.054533

Fordeling

Nøgle tal

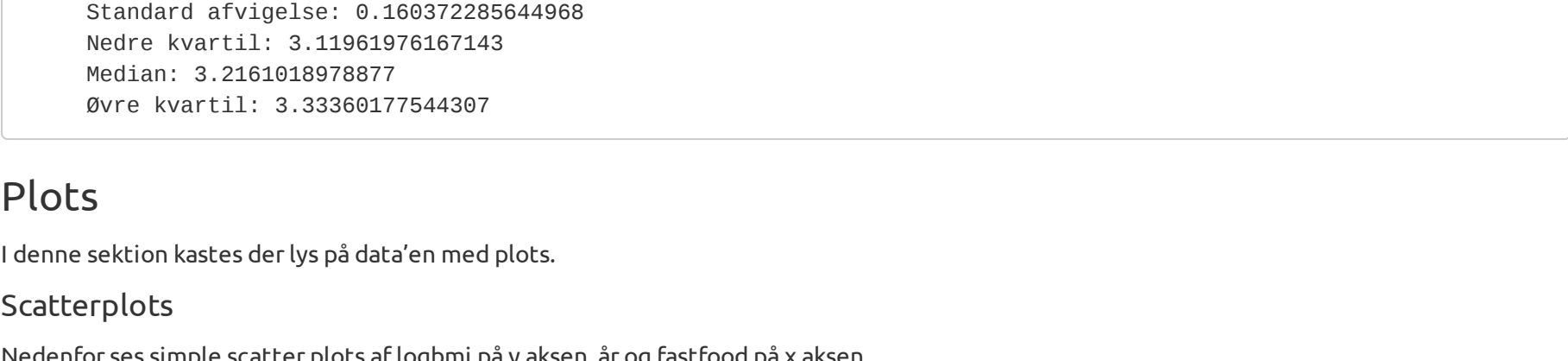
fastfood: Antal Observationer: 847 Minimum: 0.0 Maximum: 78.2 Variation: 3866.1831280813 Standard afvigelse: 32.4512381754146 Median: 6 Q1: 0 Q3: 24	age: Antal Observationer: 847 Minimum: 18 Maximum: 78 Variation: 311.28248872555 Standard afvigelse: 14.332798790713 Median: 32 Q1: 24 Q3: 37	logbmi: Antal Observationer: 847 Minimum: 2.2584847320382 Maximum: 3.2181818788777 Variation: 0.025182798929911 Standard afvigelse: 0.1802723864848 Median: 3.054533 Q1: 3.054533 Q3: 3.054533
---	--	---

Plots

I denne sektion kastes der lys på data'en med plots.

Scatterplots

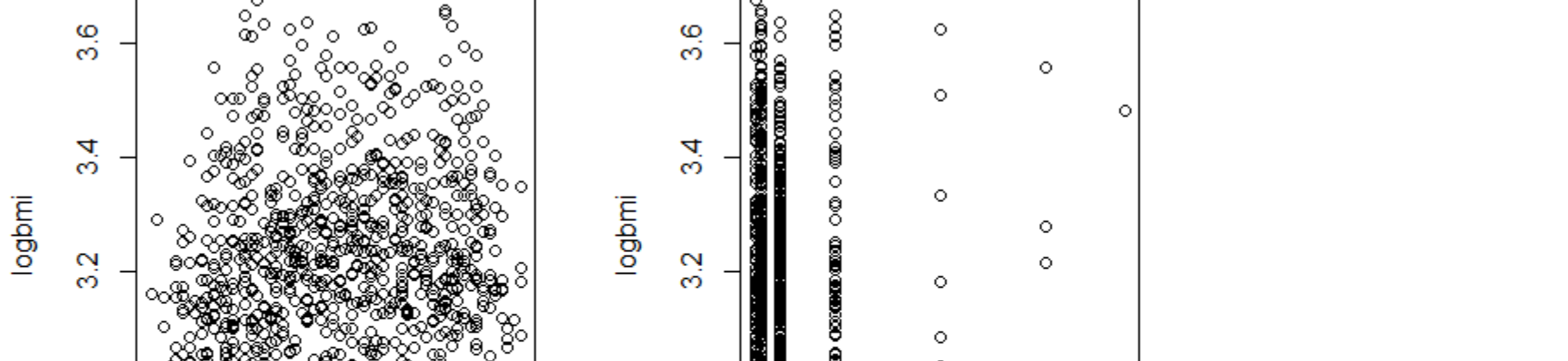
Nedenfor ses simple scatter plots af logbmi på y-aksen, år og fastfood på x-aksen.



En lille ting der bør nævnes er at fastfood har en kategoriserende natur, men er i kvantitativ form, derfor ses observationerne af fastfood til figne på linjen.

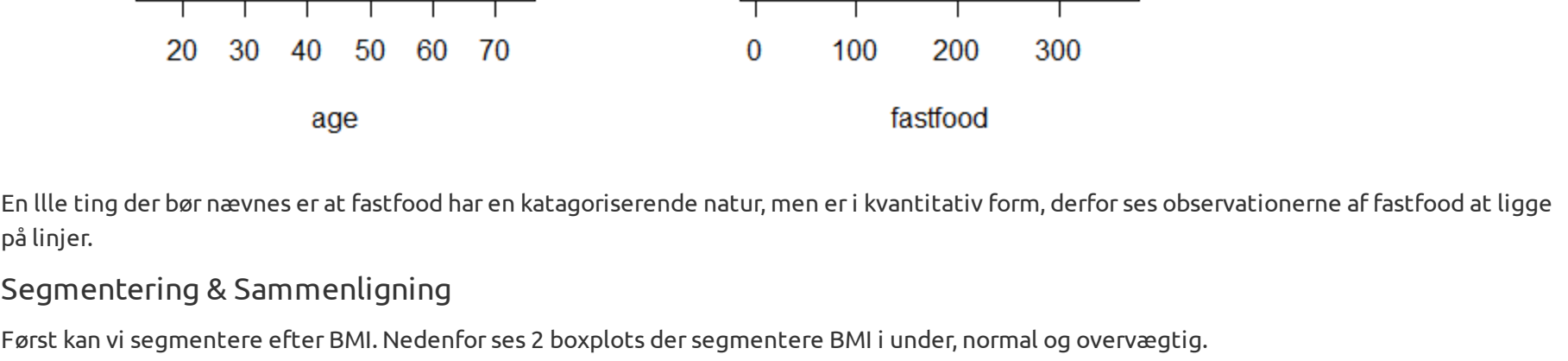
Segmentering & Sammenligning

Først kan vi segmentere efter BMI. Nedenfor ses 2 boxplots der segmenterer BMI i under, normal og overvejgt.



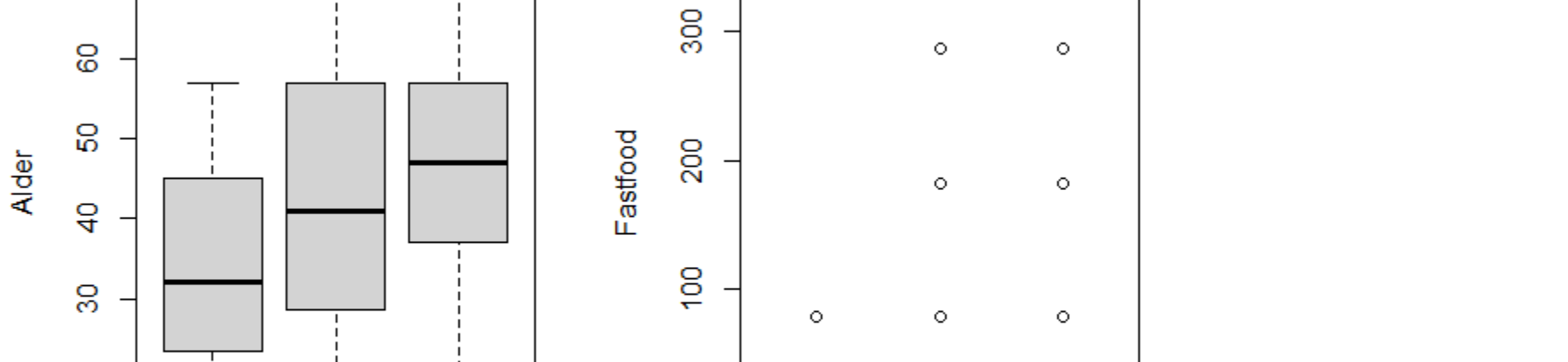
Man kan se at der er en tendens til at mennesker med højere BMI typisk er ældre, samt at de med højere BMI også findes flere ekstremer tilfælde af fastfood.

Nedenfor ses 2 boxplots hvor segmenteringen følger fastfood.



Med disse boxplots kan man se at det generelt er den yngre del af tilgruppen der spiser mest fastfood. Det ses også at med højere fastfood forbrug er spredningen af BMI endnu større og generelt højere BMI, hvilket kunne indikere at fastfood har en øgende indflydelse på bmi. Dog bør der bruges på flere faktorer for at konkludere, da det nok ikke kun er fastfood eller alder der har indflydelse på bmi.

Nedenfor ses histogrammer for bmi, fastfood og alder.



Histogrammene viser en højere normal fordeling logbmi. Histogrammet viser også at fastfood generelt ikke er en daglig ting, men mere til særlige dage. Dog er der sjældne tilfælde hvor fastfood er tæt på en dagligdag ting. Slikpræven lader også til at følge en mere uniform fordeling. Hvis vi ser bort fra helt unge og helt gamle mennesker, har vi nogenlunde lige mange mennesker i alders grupperne fra slut 20'erne til start 60'erne, så det er en ret normal fordeling.

b) Multipel lineær regressionsmodel

Det forudsættes at residualerne er normal fordelt. Denne forudsætning er vigtig for den multi lineære regression.

Multipel lineær regressionsmodel:

- $\log_{bmi} = \beta_0 + \beta_1 \times \text{age} + \beta_2 \times \text{fastfood} + \epsilon, \epsilon_i \sim N(0, \sigma^2)$

c) Model parametre

Bruger R til udregning af modellens parametre.

```
Call:
lm(formula = logbmi ~ age + fastfood, data = D_model)

Residuals:
    Min       1Q   Median       3Q      Max
-0.37643  -0.11384  -0.01488   0.09736   0.48839

Coefficients:
(Intercept)  3.1124288  0.0193517 160.835  < 2e-16 ***
age          0.0023744  0.0003889  0.186  3.58e-09 ***
fastfood     0.0005484  0.0001732  3.119  0.00188 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1573 on 837 degrees of freedom
Multiple R-squared:  0.84497, Adjusted R-squared:  0.84259
F-statistic: 19.66 on 2 and 837 Df, p-value: 4.53e-09
```

- $\beta_0 = 3.1124288$ er intercepten og angiver en start-kontekst for tilgængeligheden for de faktorer der ikke betragtes i modellen.
- $\beta_1 = 0.0023744$ er koefficienten der angiver indflydelsen størrelsen af første forklarings variable, (i dette tilfælde hvor stor en indflydelse alder har på modellen)
- $\beta_2 = 0.0005404$ er koefficienten der angiver indflydelsen størrelsen af anden forklarings variable, (i dette tilfælde hvor stor en indflydelse fastfood har på modellen).
- $r^2 = 0.84497$ er standardiseret R^2.

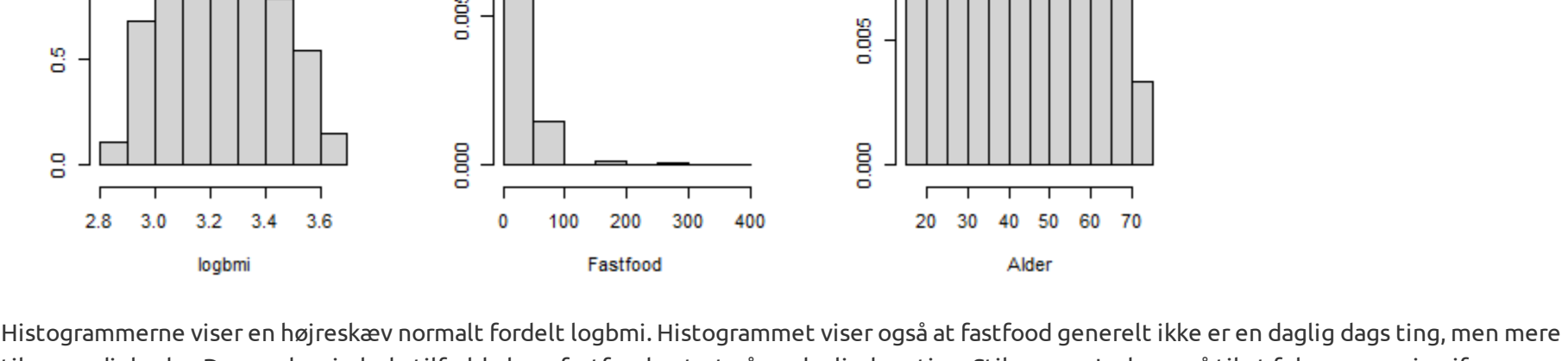
Denne fortolkning burde ikke bruges hvis forklarings variablerne er kollinerale. Men i dette tilfælde ser det ud til at være fint. Det kan blandt andet ses ved at finde correlation mellem forklaringsvariablerne (væs senere) og ved P-værdien i tidligere R-output.

d) Model kontrol

Det er vigtigt at residualerne er normal fordelt for at den multipel lineære regression passer. Nedenfor ses et plot af de fittede værdier imod dets residual. Hvis der ses bort fra de enkelte ekstremer residualer, er der ikke umiddelbart noget system over residualerne, samt at residualerne følger og plottet rimeligt tæt.



Udover ovenstående, kan vi også lade med et "Tins wally" experiment. Nedenfor ses 9 box plots som experimentet er lavet på. Experimentet forsøger at udpege en eventuel afvigning fra en normal fordeling. Der kan ikke umiddelbart ses nogen afvigning fra at residualerne er normal fordelt.



NB Vi kender ikke den helt rigtige model, men vi går nok ikke helt galt ved at antage modellens residualer er normal fordelt.

e) Alders koefficient konfidensinterval

Bruger formel:

$$\hat{\beta}_1 \pm t_{\alpha/2} \times \text{SE}(\hat{\beta}_1)$$

Værdier aflæst fra tidligere R-output:

- $\hat{\beta}_1 = 0.0023743902$
- $\text{SE}(\hat{\beta}_1) = 0.0003880714$
- $t_{\alpha/2} = 1.9629222725$
- $0.0023743902 \pm 1.9629222725 \times 0.0003880714$

Herved konfidensintervallet: [0.00110886, 0.00317834] Det er det samme resultat der fås med R-funktion 'confint'.

	2.5 %	97.5 %
(Intercept)	3.074463234	3.150432872
age	0.001618881	0.003178342
fastfood	0.000293259	0.000803957

f) Hypotesetest

Givet nullhypotesen:

$$H_0: \beta_1 = 0.001 \implies \beta_1 - 0.001 = 0$$

Den alternative hypotesen:

$$H_1: \beta_1 - 0.001 \neq 0$$

Givet signifikansniveauet $\alpha = 0.05$ kan vi teste hypotesen ved at udregne p-værdien.

Beregner test-størrelsen med formel:

$$t_{\text{stat}} = \frac{\hat{\beta}_1 - \beta_0}{\text{SE}(\hat{\beta}_1)}$$

Derefter kan test-størrelsen bruges til at slå op i fordelingen for at finde p-værdien med formel:

$$p = 2P(T > |t_{\text{stat}}|)$$

Finder test-størrelsen for indflydelsen størrelsen af alder ved værdien $0.001: t_{\text{stat}} = \frac{0.0023743902 - 0.001}{0.0003880714} = 3.5333194163$

Slår op i fordelingen med test-størrelsen, 837 frihedsgrader og får $p = 2P(T > 3.5333194163) = 0.0004328512$

P-værdien er meget lav, hvilket betyder der er stærk evidens imod hypotesen. Vi afser derfor hypotesen med et signifikans niveauet på 0.05, da p-værdien er mindre.

g) Backward selection

Med backward selection forsøger vi at reducere modellen ved at ignorere variable der ingen indflydelse på modellen har. Vi starter med modellen som beskrevet i opgave del b.

Nedenfor ses udregninger med R.

Correlation:

De følgende 3 correlation'er er hhv. bmi-fastfood, bmi-age, fastfood-age og er udregnet med R's cor-funktion.

[1]	0.0604365	0.1672487	-0.2856725
-----	-----------	-----------	------------

I blandt de 3 variable har ingen af dem en betydelig klar correlation, derfor kan vi ikke på det grundlag fjerne dem. Hvis der havde været en variabel der angav minutter siden fødsel, havde den haft en høj correlation med alder og kunne derfor ignoreres.

Confidens intervaller

De følgende confidens intervaller er for den multipel lineære regressions models parametre.

	2.5 %	97.5 %
(Intercept)	3.074463234	3.150432872
age	0.001618881	0.003178342
fastfood	0.000293259	0.000803957

På R's output ses det at alle variablerne er statistisk signifikant med de lave p-værdier, hvilket betyder at der er evidens imod ide'en om at variablerne alder og fastfood ikke har en betydning på logbmi. Af denne grund pluses det altså alder og fastfood har en statistisk effekt på logbmi, og har derfor ikke fjernes fra modellen. Med andre ord har alder og fastfood en unik indflydelse på logbmi.

Slut model

Herved den endelige model med dets parametre:

$$\log_{bmi} = -3.074463234 + 0.0023744 \times \text{age} + 0.0005404 \times \text{fastfood} + \epsilon, \epsilon_i \sim N(0, 0.1573^2)$$

h) Prædiktioner

I dette afsnit er R's predict funktion brugt til beregninger følgende.

id	age	fastfood	logbmi
841	3.258993	0.027973	3.548935
842	3.218075	0.051892	3.519949
843	3.232445	0.023231	3.541208
844	3.232445	0.023231	3.541208
845	3.229870	0.02687	3.538883
846	3.229841	0.02687	3.538883
847	3.218170	0.051899	3.521443

Confidence

$\hat{\beta}_1 \approx 0.0023743602$
 $\hat{\sigma}^2_{\beta_1} = 0.0003889714$
 $t_{1-\alpha/2} = 1.9628022725$
 $0.0023743602 \pm 1.9628022725 \times 0.0003889714$

Hermed konfidensinterval: [0.001610886, 0.003137834] Det er det samme resultat der fås med R's funktion 'confint'.

Vurdering

Givet test data er modellens prædiktion. Alle test punkternes P-værdi (værdi tilpasset med indflydelse fra alder og fastfood) passer i både prediction og confidence intervallerne. Confidence intervaller reflectere midelværdien hvor prædiktions intervaller reflectere en enkelt værdi. Hvis vi udregner den rigtige bmi ved at tage værdien til exponential funktionen af de første test intervaller får følgende:

Prediktion:	[38.689893473547, 34.674824722567]
-------------	------------------------------------

Confidence:	[25.7380604941637, 25.739311329542]
-------------	-------------------------------------

Det giver meget god mening at modellen prædiktører at en eventuel næste værdi observeret vil være mellem 18 og 34 BMI. Det er måske et lidt stort interval, men det giver meget god mening. Udover prædiktionen viser den også at middeværdien burde ligge mellem 25.16 og 25.73, hvilket også er meget plausibelt.

Givet residualerne følger der ud til at være normalt fordelt med en konstant variance som vist i opgave del d, kan modellen umiddelbart godt bruges til at prædiktører.