

# Matematisk Statistik: Vejledende besvarelse af reeksamen

Steffen Lauritzen og Niels Richard Hansen

20. august, 2020

## Spørgsmål 1.1

$X$  har nu samme fordeling som  $\exp(\xi + \sigma Z)$  hvor  $Z \sim N(0, 1)$  Vi har derfor

$$\mathbf{E}(X) = \int_{-\infty}^{\infty} e^{\xi + \sigma z} \varphi(z) dz = e^{\xi} \int_{-\infty}^{\infty} e^{\sigma z} \varphi(z) dz = e^{\xi} e^{\sigma^2/2} = \exp(\xi + \sigma^2/2)$$

som ønsket.

## Spørgsmål 1.2

Vi bemærker først, at  $\xi = \log \mu - \sigma^2/2$ . Vi har så for log-likelihood funktionen, idet vi ignorerer irrelevante konstantled ( $\sigma^2$  antages jo kendt)

$$\ell_n(\mu) = \sum_{i=1}^n \frac{(\log x_i - \log \mu + \sigma^2/2)^2}{2\sigma^2}$$

og videre

$$S_n(\mu) = \sum_{i=1}^n \frac{(\log \mu - \log x_i - \sigma^2/2)}{\mu \sigma^2}$$

samt

$$\begin{aligned} I_n(\mu) &= \sum_{i=1}^n \frac{\mu \sigma^2 / \mu - (\log \mu - \log x_i - \sigma^2/2) \sigma^2}{\mu^2 \sigma^4} \\ &= \frac{n}{\mu^2 \sigma^2} - \sum_{i=1}^n \frac{(\log \mu - \log x_i - \sigma^2/2)}{\mu^2 \sigma^2} = \frac{n}{\mu^2 \sigma^2} - \frac{1}{\mu} S_n(\mu). \end{aligned}$$

## Spørgsmål 1.3

Dette kan udledes vha deltametoden: Vi har at

$$\hat{\xi}_n = \frac{\sum_i \log x_i}{n}$$

og da MLE er ækvivariant har vi så

$$\hat{\mu}_n = \exp(\hat{\xi}_n + \sigma^2/2).$$

Idet vi ved, at  $\hat{\xi}_n \sim N(\xi, \sigma^2/n)$  anvender vi deltametoden for funktionen

$$f(t) = \exp(t + \sigma^2/2)$$

og får  $f'(t) = \exp(t + \sigma^2/2)$  og derfor er

$$\hat{\mu}_n \stackrel{\text{as}}{\sim} N\left(\mu \exp(2\xi + \sigma^2), \frac{\sigma^2}{n}\right) = N\left(\mu, \frac{\mu^2 \sigma^2}{n}\right).$$

Alternativt kan vi bruge resultatet fra det forrige spm. Da  $\mathbf{E}S_n(\mu) = 0$  er Fisher informationen

$$i_n(\mu) = \mathbf{E}(I_n(\mu)) = \frac{n}{\mu^2 \sigma^2} - \mathbf{E}(S_n(\mu)) = \frac{n}{\mu^2 \sigma^2} - 0 = \frac{n}{\mu^2 \sigma^2}$$

og derfor

$$\hat{\mu}_n \stackrel{\text{as}}{\sim} N\left(\mu, \frac{\mu^2 \sigma^2}{n}\right).$$

### Spørgsmål 1.4

Idet likelihood ratio teststørrelsen er ækvivariant, kan vi omformulere hypotesen til parameteren  $\xi$  som  $H_0 : \xi = \log 6 - 1/8$ . LR testet bliver derfor et simpelt Z-test og forkaster for numerisk store værdier af

$$Z = \sqrt{10} \frac{\sum_i \log x_i / 10 - \log 6 + 1/8}{\sigma} = 0.829$$

svarende til en  $p$ -værdi på 0.343, så observationerne understøtter hypotesen fint.

Man kan naturligvis også direkte beregne likelihood ratio størrelsen og bruge de asymptotiske resultater. Det giver samme konklusion men lidt andre talværdier.

### Spørgsmål 1.5

Idet  $m(\mu) = \mathbf{E}(X) = \mu$  er

$$\tilde{\mu}_n = \frac{\sum X_i}{n}.$$

### Spørgsmål 1.6

Vi skal finde  $\mathbf{V}(X)$  og har derfor brug for andet moment i fordelingen:

$$\mathbf{E}(X^2) = \int_{-\infty}^{\infty} e^{2\xi + 2\sigma z} \varphi(z) dz = e^{2\xi} e^{2\sigma^2} = \exp(2\xi + 2\sigma^2)$$

og videre

$$\mathbf{V}(X) = \mathbf{E}(X^2) - (\mathbf{E}(X))^2 = \exp(2\xi + 2\sigma^2) - \exp(2\xi + \sigma^2) = \mu^2(e^{\sigma^2} - 1).$$

Vi har tidligere beregnet Fisher information til  $i_n(\mu) = n/(\mu^2 \sigma^2)$ . Idet  $\tilde{\mu}_n$  per definition er en central estimator af  $\mu$  giver Cramer–Raos ulighed at

$$\mathbf{V}(\tilde{\mu}_n) = \frac{\mu^2}{n}(e^{\sigma^2} - 1) \geq \frac{\mu^2 \sigma^2}{n}.$$

Idet vi kan rækkeudvikle eksponentialfunktionen har vi

$$e^{\sigma^2} - 1 = \sigma^2 + \sum_{k=2}^{\infty} \frac{(\sigma^2)^k}{k!}$$

så forskellen er især stor for store værdier af  $\sigma^2$ , så MLE er klart at foretrække i det tilfælde, når  $n$  er stor, idet den nedre grænse er lig med MLEs asymptotiske varians.

### Spørgsmål 1.7

Simulation til sammenligning af estimatorer. Der arbejdes med tre forskellige værdier af  $\sigma$ .

```

M <- 5000
n <- 10
truexi <- 0
truesigma <- 0.1
truemean <- exp(truexi+truesigma^2/2)

```

Klargøring af arrays til resultaterne

```

muhat <- rep(0, M)
mutilde <- rep(0, M)
xihat <- rep(0,M)

```

Solve simulationerne

```

for (i in 1:M)
{
  simx <- rnorm(n) # simulerede standardnormalfordelte
  simy <- exp(truexi+truesigma*simx)

  xihat[i] = mean(log(simy))
  muhat[i] <- exp(xihat[i]+truesigma^2/2)      # mle

  mutilde[i]<-mean(simy)                      # alternativ estimator
}

```

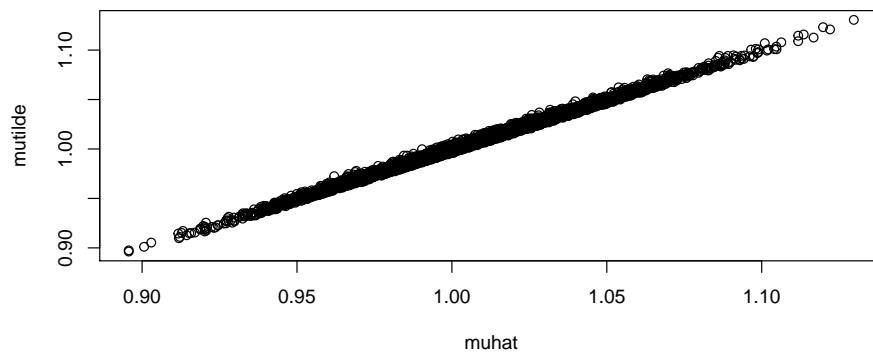
Scatterplot.

Der er naesten ingen forskel paa de to estimatorer, naar sigma er lille.

```

simData <- data.frame(muhat,mutilde)
plot(simData)

```



Resultaterne

```
truemean
```

```
[1] 1.005013
```

```
mean(mutilde)
```

```
[1] 1.005311
```

```
mean(muhat)
```

```
[1] 1.005832
```

```
sd(muhat)
```

```
[1] 0.03240237
```

```
sd(mutilde)
```

```
[1] 0.03245558
```

MLE har lidt mindre varians. Vi ser i stedet på mean square error: MLE har lige akkurat den mindste MSE

```
msehat <- (mean(muhat)-truemean)^2+sd(muhat)^2  
msehat
```

```
[1] 0.001050585
```

```
msetilde <- (mean(mutilde)-truemean)^2+sd(mutilde)^2  
msetilde
```

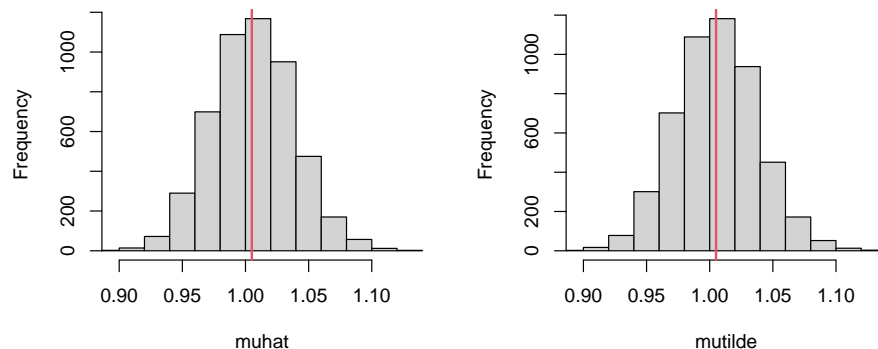
```
[1] 0.001053454
```

Ønsker fælles akser

```
myRange <- range(c(muhat, mutilde))
```

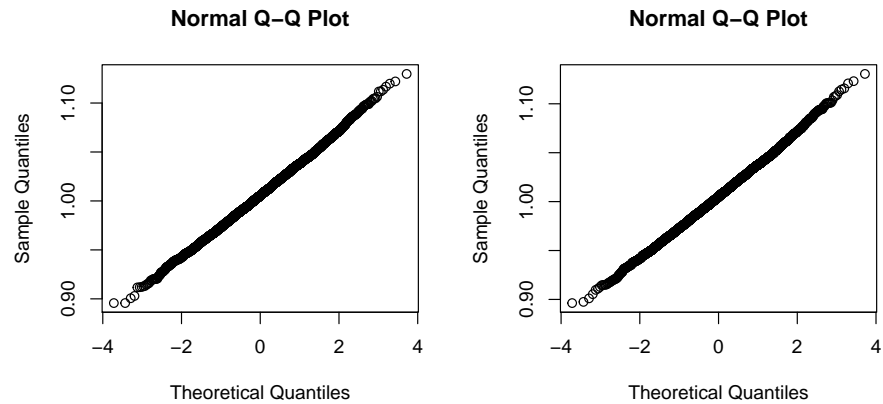
Histogrammer med sande værdi som lodret linie. De er næsten helt ens.

```
par(mfrow=c(1,2))  
hist(muhat, xlim=myRange, main="")  
abline(v=truemean, col=2, lwd=2)  
hist(mutilde, xlim=myRange, main="")  
abline(v=truemean, col=2, lwd=2)
```



Og de er begge ret fint normalfordelte

```
par(mfrow=c(1,2))  
qqnorm(muhat)  
qqnorm(mutilde)
```



Gentag for større værdi af truesigma.

```
truesigma <- 0.5
truemean <- exp(truexi+truesigma^2/2)
```

Klargøring af arrays til resultaterne

```
muhat <- rep(0, M)
mutilde <- rep(0, M)
xihat <- rep(0, M)
```

Solve simulationerne

```
for (i in 1:M)
{
  simx <- rnorm(n) # simulerede standardnormalfordelte
  simy <- exp(truexi+truesigma*simx)

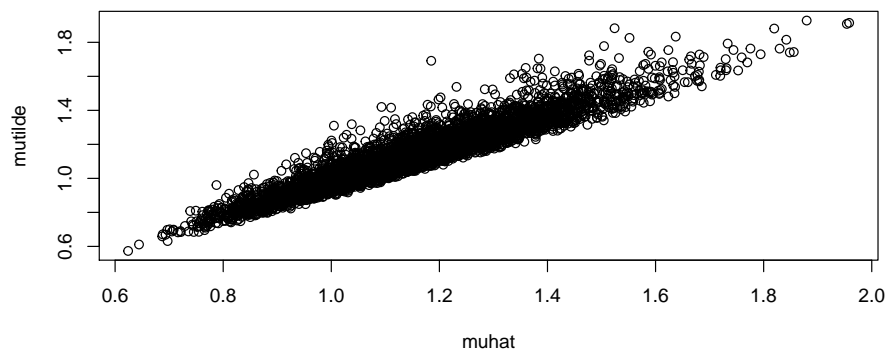
  xihat[i] = mean(log(simy))
  muhat[i] <- exp(xihat[i]+truesigma^2/2)      # mle

  mutilde[i]<-mean(simy)                      # alternativ estimator
}
```

Scatterplot.

Nu er forskellen lidt større.

```
simData <- data.frame(muhat,mutilde)
plot(simData)
```



Resultaterne

```
truemean
```

```
[1] 1.133148
```

```
mean(mutilde)
```

```
[1] 1.132086
```

```
mean(muhat)
```

```
[1] 1.147048
```

```
sd(muhat)
```

```
[1] 0.1824333
```

```
sd(mutilde)
```

```
[1] 0.1897178
```

MLE har lidt mindre varians, men overvurderer middelværdien. Vi ser i stedet på mean square error: MLE har lige akkurat den mindste MSE

```
msehat <- (mean(muhat)-truemean)^2+sd(muhat)^2  
msehat
```

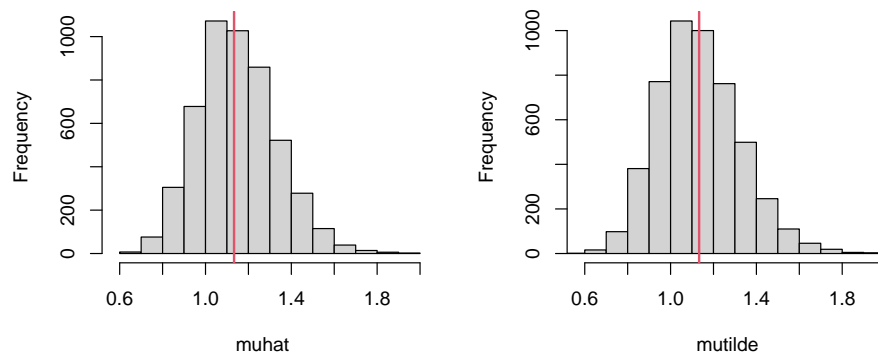
```
[1] 0.0334751
```

```
msetilde <- (mean(mutilde)-truemean)^2+sd(mutilde)^2  
msetilde
```

```
[1] 0.03599396
```

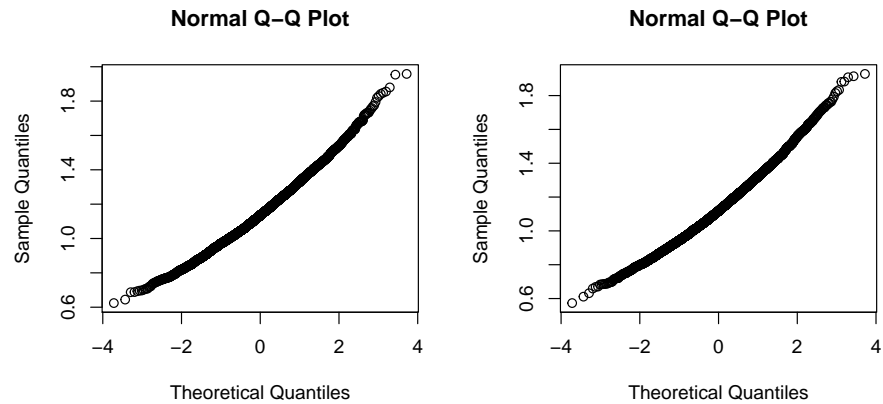
Histogrammer med sande værdi som lodret linie. De er ikke helt ens længere.

```
myRange <- range(c(muhat, mutilde))  
par(mfrow=c(1,2))  
hist(muhat, xlim=myRange, main="")  
abline(v=truemean, col=2, lwd=2)  
hist(mutilde, xlim=myRange, main="")  
abline(v=truemean, col=2, lwd=2)
```



Og de er ikke længere pænt normalfordelte.

```
par(mfrow=c(1,2))  
qqnorm(muhat)  
qqnorm(mutilde)
```



Gentag for meget større værdi af truesigma.

```
truesigma <- 1
truemean <- exp(truexi+truesigma^2/2)
```

Klargøring af arrays til resultaterne

```
muhat <- rep(0, M)
mutilde <- rep(0, M)
xihat <- rep(0, M)
```

Solve simulationerne

```
for (i in 1:M)
{
  simx <- rnorm(n) # simulerede standardnormalfordelte
  simy <- exp(truexi+truesigma*simx)

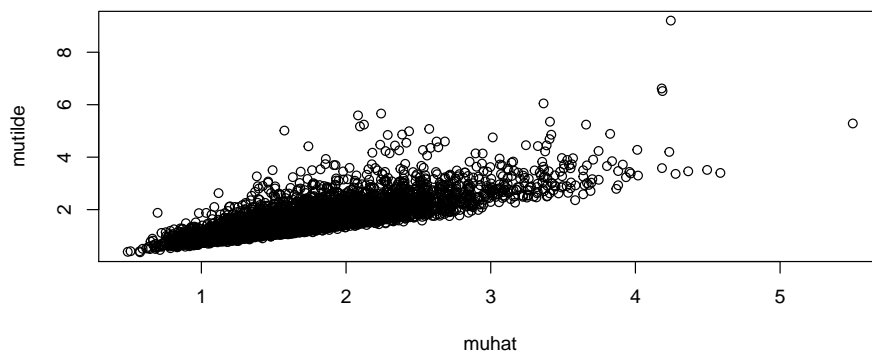
  xihat[i] = mean(log(simy))
  muhat[i] <- exp(xihat[i]+truesigma^2/2) # mle

  mutilde[i]<-mean(simy) # alternativ estimator
}
```

Scatterplot.

Nu er de slet ikke ens.

```
simData <- data.frame(muhat,mutilde)
plot(simData)
```



Resultaterne

```
truemean
```

```
[1] 1.648721
```

```
mean(mutilde)
```

```
[1] 1.647151
```

```
mean(muhat)
```

```
[1] 1.737378
```

```
sd(muhat)
```

```
[1] 0.5697215
```

```
sd(mutilde)
```

```
[1] 0.6792565
```

MLE har nu klart mindre varians, men overvurderer stadig middelværdien. Vi ser i stedet på mean square error: MLE har klart den mindste MSE.

```
msehat <- (mean(muhat)-truemean)^2+sd(muhat)^2  
msehat
```

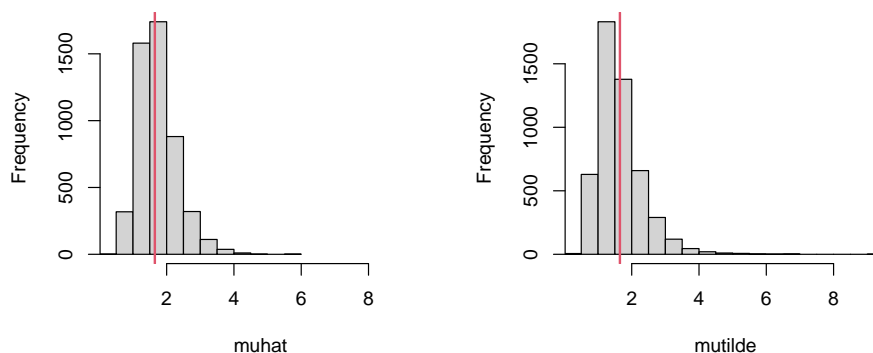
```
[1] 0.3324427
```

```
msetilde <- (mean(mutilde)-truemean)^2+sd(mutilde)^2  
msetilde
```

```
[1] 0.4613919
```

Histogrammer med sande værdi som lodret linie. De er nu meget forskellige og  $\tilde{\mu}$  er sommetider meget for stor.

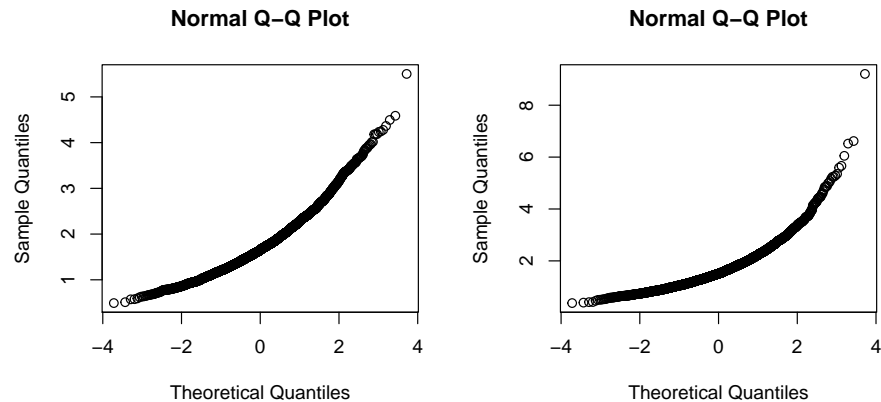
```
myRange <- range(c(muhat, mutilde))  
par(mfrow=c(1,2))  
hist(muhat, xlim=myRange, main="")  
abline(v=truemean, col=2, lwd=2)  
hist(mutilde, xlim=myRange, main="")  
abline(v=truemean, col=2, lwd=2)
```



Og de er nu meget langt fra at være normalfordelte, Så  $n = 10$  er ikke nok til at den asymptotiske normalfordeling er en god approximation.

```
par(mfrow=c(1,2))  
qqnorm(muhat)  
qqnorm(mutilde)
```





## Spørgsmål 2.1

Vi indlæser data til opgaven.

```
restaurant <- read_csv("restaurant.txt",
  col_types = cols(D = col_factor(
    levels = c("m", "ti", "o", "to", "f", "l", "s"))))
```

Først tabuleres  $T \times B$ .

```
table(restaurant$T, restaurant$B)
```

```
      br fr mi
fa  8 20  8
ta  8 20  8
```

Vi ser, at alle seks kombinationer af  $T$  og  $B$  forekommer, så  $\dim(L_{T \times B}) = 6$ . Bemærk at alle indgange er positive og rækkerne er identiske, så  $T \wedge B = 1$  og faktorerne opfylder balanceligningen og er således geometrisk ortogonale.

Dernæst tabuleres  $B \times D$ .

```
table(restaurant$B, restaurant$D)
```

```
      m ti o to f l s
br 0  0 0  0 0 8 8
fr 8  8 8  8 8 0 0
mi 0  0 0  0 0 8 8
```

Heraf fremgår det, at designrafen har to sammenhængskomponenter. En svarende til frokost på hverdage, og en svarende til serveringer lørdage og søndage (altså i weekenden). Så minimum har to niveauer, og kan opfattes som en indikator for om det er en weekendservering. Bemærk iøvrigt også at indenfor hver sammenhængskomponent er designet fuldstændigt balanceret, så også  $B$  og  $D$  er geometrisk ortogonale.

Der er tre ikke-trivielle minima tilbage at undersøge. Det er  $T \times B \wedge D$ ,  $D \wedge T$  og  $T \wedge W$ . Først ser vi at  $1 = B \wedge T \geq W \wedge T \geq 1$ , så  $W \wedge T = 1$ . Tabellen

```
table(interaction(restaurant$T, restaurant$B), restaurant$D)
```

```
      m ti o to f l s
fa.br 0  0 0  0 0 4 4
ta.br 0  0 0  0 0 4 4
fa.fr 4  4 4  4 4 0 0
```

```

ta.fr 4 4 4 4 0 0
fa.mi 0 0 0 0 0 4
ta.mi 0 0 0 0 0 4

```

viser at  $T \times B \wedge D = W$  og tabellen

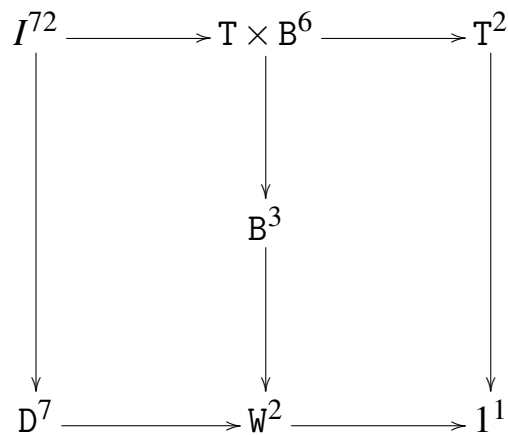
```
table(restaurant$T, restaurant$D)
```

```

m t i o t o f l s
fa 4 4 4 4 4 8 8
ta 4 4 4 4 4 8 8

```

viser at  $D \wedge T = 1$ . Fakturstrukturdiagrammet (med identitetsfaktoren  $I$  tilføjet) er som følger.



Diagrammet er ovenfor annoteret med dimensioner vi kender på nuværende tidspunkt.

## Spørgsmål 2.2

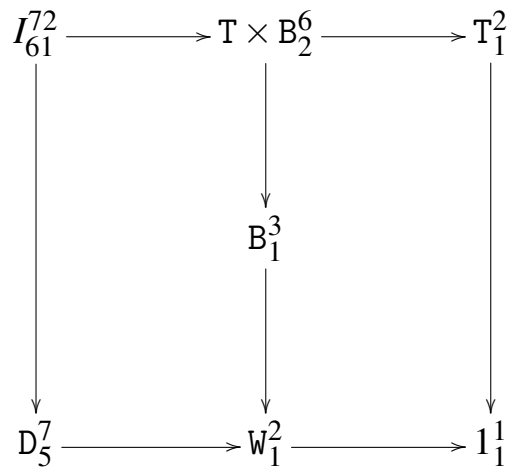
Tabellerne ovenfor viser at balanceligningen (sætning 14.8) er opfyldt for  $T \times B$  og  $D$ ,  $D \times B$ ,  $D \times T$ , og  $B \times T$ , og disse faktorer er derfor geometrisk ortogonale. De øvrige faktorer i designet på nær  $W$  og  $T$  opfylder en ordningsrelation og er således geometrisk ortogonale ifølge lemma 14.11. Da  $B \geq W \geq B \wedge T = 1$  følger det af lemma 14.12 at også  $W$  og  $T$  er geometrisk ortogonale.

Alternativt kunne man indføre  $W$  i data som en faktor og lave tabellen derfra, men det er lidt mere bøvlet, eller finde  $T \times W$ -tabellen fra  $T \times B$ -tabellen

	weekend	hverdag
fa	16	20
ta	16	20

og indse at den ligeledes opfylder balanceligningen.

Konklusionen er, at designet er geometrisk ortogonalt, og da det ligeledes er  $\wedge$ -stabilt kan vi bruge sætning 14.21 til at beregne  $V_G$ -dimensioner. Diagrammet nedenfor er annoteret med disse dimensioner



Så

$$\dim(L_D + L_{T \times B}) = 1 + 1 + 5 + 1 + 1 + 2 = 11$$

(eller alternativt  $\dim(L_D + L_{T \times B}) = \dim(L_D) + \dim(L_{T \times B}) - \dim(L_W) = 7 + 6 - 2 = 11$ ) og

$$\dim(L_D + L_T + L_B) = 1 + 1 + 5 + 1 + 1 = 9.$$

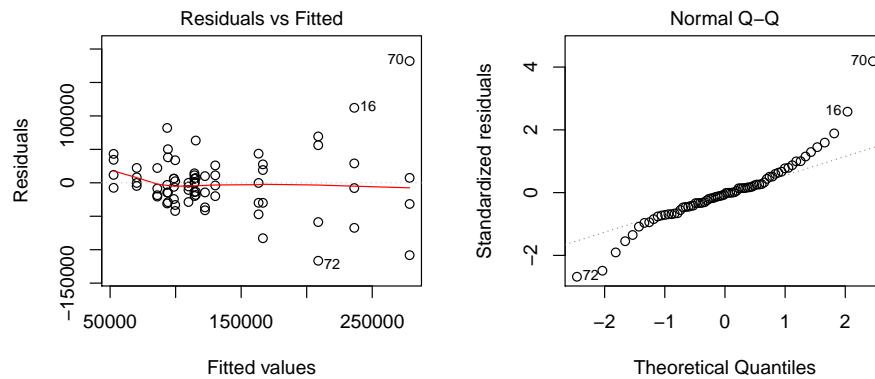
### Spørgsmål 2.3

Vi fitter de to modeller i R.

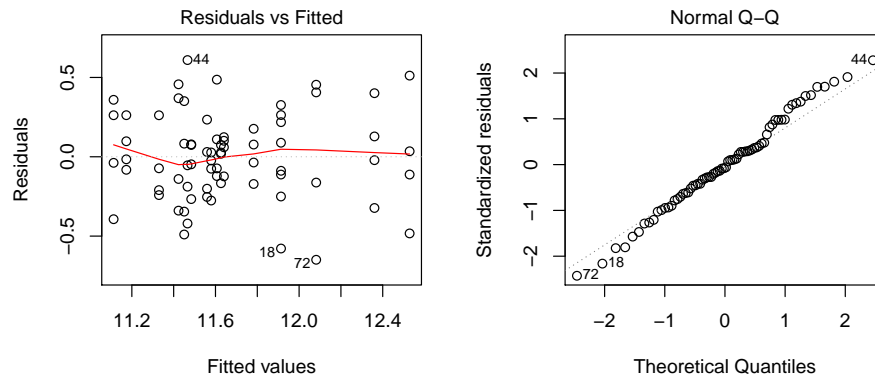
```
rest_lm <- lm(X ~ D + B * T, data = restaurant)
log_rest_lm <- lm(log(X) ~ D + B * T, data = restaurant)
```

Dernæst ser vi på residualplot og qqplot for residualerne for de to modeller.

```
plot(rest_lm, 1:2)
```



```
plot(log_rest_lm, 1:2)
```



Modellen for  $X$  fitter ikke data særligt godt. Der er en tydelig tragtform af residualplottet som viser, at variansen ikke er konstant. QQplottet viser også at residualerne ikke er normalfordelt, men har tungere haler end normalfordelingen.

Modellen for  $\log(X)$  fitter data meget bedre. Der er ingen åbenlyse systematiske afvigelser i residualplottet, og QQplottet viser at residualerne ser pænt normalfordelte ud.

## Spørgsmål 2.4

Det er klart at  $L_D + L_B \subseteq L_D + L_T + L_B$  da vi her blot lægger  $L_T$  til  $L_D + L_B$ . Endvidere er  $L_T + L_B \subseteq L_{T \times B}$ , så når vi lægger  $L_D$  til på begge sider fås

$$L_D + L_T + L_B \subseteq L_D + L_{T \times B}.$$

Vi benytter nedenfor modellen for  $\log(X)$  som udgangspunkt, da den fitter data bedst. Vi kan vælge at teste en effekt af reklameplatformen ( $T$ ) på omsætningen på flere måder. Her vælger vi først at teste den additive model mod modellen med vekselvirkningen, og dernæst at teste den additive effekt af  $T$ .

```
log_rest_lm %>% anova()
```

Analysis of Variance Table

```
Response: log(X)
      Df Sum Sq Mean Sq F value    Pr(>F)
D       6  4.7309   0.7885   9.2982 3.121e-07 ***
B       1  4.0118   4.0118  47.3091 3.751e-09 ***
T       1  0.7398   0.7398   8.7237 0.004456 **
B:T     2  0.1175   0.0587   0.6927 0.504121
Residuals 61  5.1727   0.0848
```

---

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

Vi aflæser her en  $p$ -værdi på omkring 0.5 for  $F$ -testen for vekselvirkningen mellem  $T$  og  $B$ , så vi kan ikke afvise den additive hypotese. Det efterfølgende test for at der ikke er en additiv effekt af  $T$  har en  $p$ -værdi omkring 0.0045, som er relativt lille, og vi afviser hypotesen om, at der ikke er en effekt af reklameformen på omsætningen. Vi konkluderer altså på basis af analysen at der er en additiv effekt af reklameplatformen på omsætningen på en log-skala.

Bemærk at anova udregner sekventielle test på en lidt anden måde, end hvis vi udførte dem et ad gangen, men konklusionen er den samme.

```
lm(log(X) ~ D + B + T, data = restaurant) %>% anova()
```

Analysis of Variance Table

```
Response: log(X)
      Df Sum Sq Mean Sq F value    Pr(>F)
```

```

D          6 4.7309  0.7885  9.3898 2.368e-07 ***
B          1 4.0118  4.0118 47.7751 2.845e-09 ***
T          1 0.7398  0.7398  8.8097 0.004233 **
Residuals 63 5.2902  0.0840
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Endelig kunne vi også have testet modellen  $D+B$  direkte

```
anova(lm(log(X) ~ D + B, data = restaurant), log_rest_lm)
```

Analysis of Variance Table

```

Model 1: log(X) ~ D + B
Model 2: log(X) ~ D + B * T
  Res.Df    RSS Df Sum of Sq    F Pr(>F)
1      64 6.0300
2      61 5.1727  3   0.85724 3.3697 0.0241 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

der giver en  $p$ -værdi på omkring 0.024, hvilket er en lille  $p$ -værdi, men den efterlader os ikke med en klar konklusion. Endvidere belyser dette ene test ikke hvorvidt der er en vekselvirkning eller ej.

## Spørgsmål 2.5

Vi tager udgangspunkt i den additive model fra spørgsmålet ovenfor af log-omsætningen, hvor effekten af reklameformen på omsætningen udtrykkes ved parameteren tilhørende faktoren  $T$ .

```
summary(lm(log(X) ~ D + B + T, data = restaurant)) %>% coef()
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	11.15019519	0.10799406	103.2482268	5.318629e-72
Dti	0.24979254	0.14488924	1.7240241	8.960772e-02
Do	0.27693412	0.14488924	1.9113506	6.051101e-02
Dto	0.45288933	0.14488924	3.1257624	2.682046e-03
Df	0.31059686	0.14488924	2.1436848	3.592506e-02
Dl	0.48484573	0.13553147	3.5773664	6.738595e-04
Ds	0.03787857	0.13553147	0.2794817	7.807907e-01
Bmi	0.70814493	0.10245216	6.9119570	2.844753e-09
Tta	0.20272641	0.06830144	2.9681132	4.232502e-03

Vi ser fra ovenstående summary at forskellen på de to reklameformer er estimeret til 0.2027 (parameteren hørende til Tta), som er positiv. Estimatet udtrykker at ta platformen på giver en forøgelse af omsætningen estimeret med en faktor  $\exp(0.2027) = 1.225$  i forhold til fa platformen. Vi kan udregne et (standard) 95% konfidensinterval på log-skalaen, men det er mere informativt at transformere det til den oprindelige skala.

```
confint(lm(log(X) ~ D + B + T, data = restaurant), "Tta")
```

```

      2.5 %      97.5 %
Tta 0.06623687 0.3392159

```

```
exp(confint(lm(log(X) ~ D + B + T, data = restaurant), "Tta"))
```

```

      2.5 %      97.5 %
Tta 1.06848 1.403846

```

Intervallat udtrykker altså at ta platformen forøger salget med en faktor mellem 1.068 og 1.404 i forhold til fa platformen.

### Spørgsmål 3.1

Vi indlæser først data.

```
bodyfat <- read_csv("bodyfat.txt")
n <- nrow(bodyfat)
n ## Der er 240 observationer
```

```
[1] 240
```

Fra sætning 5.5 (se også s. 29 i de supplerende noter) har vi, at MLE for  $\mu$  er gennemsnittet

```
colMeans(bodyfat)
```

```
      Weight      BodyFat
178.04917    19.12167
```

og MLE for  $\Sigma$  er  $\frac{1}{n}S$ . Den simpleste måde at beregne MLE på er nok ved brug af cov og en reskalering (da den funktion beregner  $\frac{1}{n-1}S$ .)

```
Sigma_hat <- (n - 1) / n * cov(bodyfat)
Sigma_hat
```

```
      Weight      BodyFat
Weight  695.1475 132.66356
BodyFat 132.6636  67.30078
```

Alternativt kan MLE beregnes på den her måde

```
as.matrix(bodyfat) %>%
  scale(scale = FALSE) %>%
  crossprod() %>%
  "/"(n)
```

```
      Weight      BodyFat
Weight  695.1475 132.66356
BodyFat 132.6636  67.30078
```

Hypotesen  $H_0$  er et specialtilfælde af  $H$  givet ved (16) i de supplerende noter (og identisk med hypotesen  $H$  givet ved (18) i eksempel 6.5), så korollar 6.3 giver at MLE af  $\mu$  fortsat er gennemsnittet som udregnet ovenfor, og MLE af  $\Sigma$  er diagonalmatricen

```
Sigma_hat %>%
  diag() %>%
  diag()
```

```
      [,1]      [,2]
[1,] 695.1475 0.00000
[2,]  0.0000 67.30078
```

Vi beregner først korrelationen.

```
rho <- Sigma_hat[1, 2] / sqrt(Sigma_hat[1, 1] * Sigma_hat[2, 2])
rho
```

```
[1] 0.6133426
```

```
cor(bodyfat$Weight, bodyfat$BodyFat) ## Alternativ beregning af korrelation
```

```
[1] 0.6133426
```

Vi tester hypotesen med et kvotienttest, og eksempel 6.5 giver at  $-2\log Q$  er

```
logQ <- - n * log (1 - rho^2)
logQ
```

```
[1] 113.2579
```

Sætning 6.4 giver at under  $H_0$  er  $-2\log Q$  asymptotisk  $\chi^2_1$ -fordelt, så  $p$ -værdien er

```
pchisq(logQ, df = 1, lower.tail = FALSE)
```

```
[1] 1.894542e-26
```

som er ekstremt lille, og vi afviser derfor hypotesen  $H_0$  om uafhængighed.

## Spørgsmål 3.2

Vi udregner først  $\Sigma^{-1}$  under hypotesen  $H_1$ . Idet

$$\det(\Sigma) = \sigma_1^2 \sigma_2^2 (1 - 1/4) = \frac{3}{4} \sigma_1^2 \sigma_2^2$$

har vi at

$$\Sigma^{-1} = \frac{4}{3\sigma_1^2 \sigma_2^2} \begin{pmatrix} \sigma_2^2 & -\frac{1}{2}\sigma_1 \sigma_2 \\ -\frac{1}{2}\sigma_1 \sigma_2 & \sigma_1^2 \end{pmatrix} = \frac{2}{3} \begin{pmatrix} \frac{2}{\sigma_1^2} & -\frac{1}{\sigma_1 \sigma_2} \\ -\frac{1}{\sigma_1 \sigma_2} & \frac{2}{\sigma_2^2} \end{pmatrix}.$$

Som opgaven er formuleret, er det selvfølgelig også godt nok at verificere direkte, at  $\Sigma^{-1}\Sigma = I$ .

Det følger af sætning 3.2 at modellen er en minimal og regulær eksponentiel familie med kanonisk parameter

$$(\eta, \Omega) = (\Sigma^{-1}\mu, \Sigma^{-1}) \in \mathbb{R}^2 \times \text{PD}_2,$$

hvor dimensionen af det kanoniske parameterrum er  $2 + 3 = 5$ . Hypotesen  $H_1$  giver ingen restriktioner på  $\eta$ , der fortsat frit kan variere i  $\mathbb{R}^2$ , mens  $\Omega = \Sigma^{-1} = \phi(\sigma_1, \sigma_2)$  for

$$(\sigma_1, \sigma_2) \in (0, \infty)^2.$$

Vi udregner de partielt afledte

$$\partial_{\sigma_1} \phi(\sigma_1, \sigma_2) = \frac{2}{3} \begin{pmatrix} -\frac{4}{\sigma_1^3} & \frac{1}{\sigma_1^2 \sigma_2} \\ \frac{1}{\sigma_1^2 \sigma_2} & 0 \end{pmatrix}$$

og

$$\partial_{\sigma_2} \phi(\sigma_1, \sigma_2) = \frac{2}{3} \begin{pmatrix} 0 & \frac{1}{\sigma_1 \sigma_2^2} \\ \frac{1}{\sigma_1 \sigma_2^2} & -\frac{4}{\sigma_2^3} \end{pmatrix}.$$

Det er klart at disse to matricer er lineært uafhængige da f.eks.  $\partial_{\sigma_1} \phi(\sigma_1, \sigma_2)_{11} \neq 0$  mens  $\partial_{\sigma_2} \phi(\sigma_1, \sigma_2)_{11} = 0$ . Jacobianten har derfor altid fuld rang. Endelig ser vi at  $\phi$  er en homeomorfi idet den har en invers på sit billede, der er restriktionen af den globalt kontinuerte afbildning

$$\text{PD}_2 \ni \Omega = \begin{pmatrix} \omega_1 & \omega_{12} \\ \omega_{12} & \omega_2 \end{pmatrix} \mapsto (\sqrt{4/(3\omega_1)}, \sqrt{4/(3\omega_2)}).$$

Vi har hermed vist at betingelserne i definition 2.25 i BMS er opfyldt med  $k = 5$  dimensionen af det kanoniske parameterum for den eksponentielle familie, og med  $B = \mathbb{R}^2 \times (0, \infty)^2$ , som er en åben delmængde af  $\mathbb{R}^4$ , hvorfor  $m = 4$ . Dermed specificerer  $H_1$  en krum eksponentiel familie af dimension 4 og orden 5.