

Projektopgave 2

“Sandsynlighedsregning og statistik”

Michael Andersen – michael@diku.dk
Henrik Jensen – henrikjensen@gmail.com
Ulrik Bonde – bonde@diku.dk
Julie Nielsen – julie@diku.dk

8. januar 2010

Del 1

Opgave 1

For at vise at Y er givet ved

$$f(y) = \frac{1}{y\sqrt{2\pi\sigma^2}} \exp - \frac{(\log y - \mu)^2}{2\sigma^2}$$

Benytter vi først formelen for normalfordeling $N \sim (\mu, \sigma^2)$ givet ved M.S. 5.3.5

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp - \frac{(y - \mu)^2}{2\sigma^2}$$

For at vise at ovenstående er lige med den logaritmiske normalfordeling benytter vi sætning 5.4.1 omkring transformation af kontinuerte funktioner på \mathbb{R} .

Som siger følgende

$$q(y) = \begin{cases} p(t^{-1}(y)) \left| \frac{d}{dy} t^{-1}(y) \right| & \text{hvis } y \in (v, h) \\ 0 & \text{hvis } y \notin (v, h) \end{cases}$$

Ovenstående er produktet af to led, derfor kan vi betragte de to led individuelt.

$$p(t^{-1}(y)) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp - \frac{(\log y - \mu)^2}{2\sigma^2}$$

$$\frac{d}{dy} |t^{-1}(y)| = y^{-1} = \frac{1}{y}$$

Vi ganger nu disse to led sammen og får derfor

$$q(y) = \begin{cases} \frac{1}{y\sqrt{2\pi\sigma^2}} \exp -\frac{(\log y - \mu)^2}{2\sigma^2} & \text{hvis } y \in (0, \infty) \\ 0 & \text{ellers} \end{cases}$$

Det er hermed vist at den logaritmiske normalfordeling er givet ved ovenstående formel.

Opgave 2

Vi skal her vise at middelværdien af Y er givet ved $E(Y) = e^{\mu + \sigma^2/2}$. Vi benytter her sætning 5.2.3, der siger at hvis der gælder at

$$\int_I |t(x)| p(x) dx < \infty$$

Da er den oprindelige stokastiske variable X er baseret på normalfordelingen og det vides at alle momenter for normalfordelingen eksistere, derfor vides det også at både middelværdi og varians må eksistere. Og derfor er middelværdien givet ved

$$E(t(X)) = \int_I t(x)p(x) dx$$

Vores transformerede stokastiske variable er e^X , hvor $X \sim N(\mu, \sigma^2)$

$$E(e^X) = \int_{-\infty}^{\infty} e^x \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-(x-\mu)^2}{2\sigma^2} dx$$

Vi foretager her en substitution $t = x - \mu$

$$\begin{aligned} &= \int_{-\infty}^{\infty} e^{t+\mu} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-t^2}{2\sigma^2} dt \\ &= e^{\mu} \int_{-\infty}^{\infty} e^t \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-t^2}{2\sigma^2} dt \end{aligned}$$

Integralet i ovenstående identificeres som værende e^T , hvor $T \sim N(0, \sigma^2)$

$$\begin{aligned} E(e^T) &= \int_{-\infty}^{\infty} e^t \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-t^2}{2\sigma^2} dt \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-t^2}{2\sigma^2} + t dt \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-t^2}{2\sigma^2} + t dt \end{aligned}$$

Her er vi interesseret i at omskrive eksponenten til noget vi kan genkende.

$$-\frac{t^2}{2\sigma^2} + t = -\frac{1}{2\sigma^2} (t^2 - 2\sigma^2 t)$$

$$= -\frac{1}{2\sigma^2} \left((t - \sigma^2)^2 - (\sigma^2)^2 \right)$$

$$= -\frac{(t - \sigma^2)^2}{2\sigma^2} + \frac{\sigma^2}{2}$$

Det sætter vi nu ind igen, og får følgende

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(t - \sigma^2)}{2\sigma^2} + \frac{\sigma^2}{2}$$

Dette kan vi nu omskrive til følgende

$$= \exp \frac{\sigma^2}{2} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(t - \sigma^2)}{2\sigma^2}$$

Integralet genkender man som værende tæthedsfunktionen for $T \sim N(\sigma^2, \sigma^2)$ som i henhold til 5.1.3 giver 1.

Sammen med resultatet fra $X \sim N(\mu, \sigma^2)$, giver dette.

$$\exp \mu \exp \frac{\sigma^2}{2} = \exp \mu + \frac{\sigma^2}{2}$$

Hermed er det vist at middelværdien for Y er givet ved ovenstående.

Opgave 3

Der skal bestemmes medianen af fordelingen for Y . Hvis x er givet ved fordelingen F_X , da gælder at

$$q_x = F_X^{-1}$$

Da $Y = \exp(X)$, gælder der

$$F_Y(y) = P(Y \leq y) = P(\exp X \leq y) = P(X \leq \log Y) = F_X(\log Y)$$

Og derfor gælder følgende

$$q_Y(z) = F_Y^{-1}(z) = (F_X \circ \log)^{-1}(z) = \exp q_X(z)$$

Derfor kan der nu indsættes $\frac{1}{2}$ som udgør medianen i fordeling.

$$q_Y\left(\frac{1}{2}\right) = \exp q_X\left(\frac{1}{2}\right)$$

Der findes ikke noget eksplicit udtryk for værdien af fordelingsfunktion for normalfordelingen. Men det der skal løses er

$$\Phi(x) = \int_{-\infty}^x \phi(y) dy$$

$$\frac{1}{2} = \int_{-\infty}^x \phi(y) dy$$

Da vi nu betragter $X \sim N(\mu, \sigma^2)$, og det vides at μ er middelværdien og findes der hvor $x = \mu$, må medianen være lige med

$$\exp \mu$$

Opgave 4

Vi benytter igen samme fremgangs måde som i opgave 2, og betragter derfor

$$(E(e^T))^2 = \int_{-\infty}^{\infty} (\exp t)^2 \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-(t-\mu)^2}{2\sigma^2} dt$$

Vi laver igen substitutionen $x = t - \mu$

$$\exp 2\mu \int_{-\infty}^{\infty} = \exp 2x \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-x^2}{2\sigma^2} dx$$

$$\exp 2\mu \int_{-\infty}^{\infty} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-x^2}{2\sigma^2} + 2x dx$$

Vi tager nu eksponenten ud igen og behandler den

$$\frac{-x^2}{2\sigma^2} + 2x = -\frac{1}{2\sigma^2} (x^2 - 2x2\sigma^2)$$

$$\begin{aligned} &= \frac{1}{2\sigma^2} \left(-(x - 2\sigma^2)^2 - (2\sigma^2)^2 \right) \\ &= -\frac{(x - 2\sigma^2)^2}{2\sigma^2} + \frac{(2\sigma^2)^2}{2\sigma^2} \end{aligned}$$

$$= -\frac{(x - 2\sigma^2)^2}{2\sigma^2} + \frac{4\sigma^4}{2\sigma^2} = -\frac{(x - 2\sigma^2)^2}{2\sigma^2} + 2\sigma^2$$

Dette kan nu indsættes og der fås derved

$$\exp 2\mu \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x - \sigma^2)^2}{2\sigma^2} + 2\sigma^2 dx$$

$$\exp 2\mu + 2\sigma^2 \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x - \sigma^2)^2}{2\sigma^2} dx$$

Det tilbageværende integrale genkendes som værende $X \sim N(\sigma^2, \sigma^2)$.

$$Var(Y) = E(Y^2) - (E(Y))^2$$

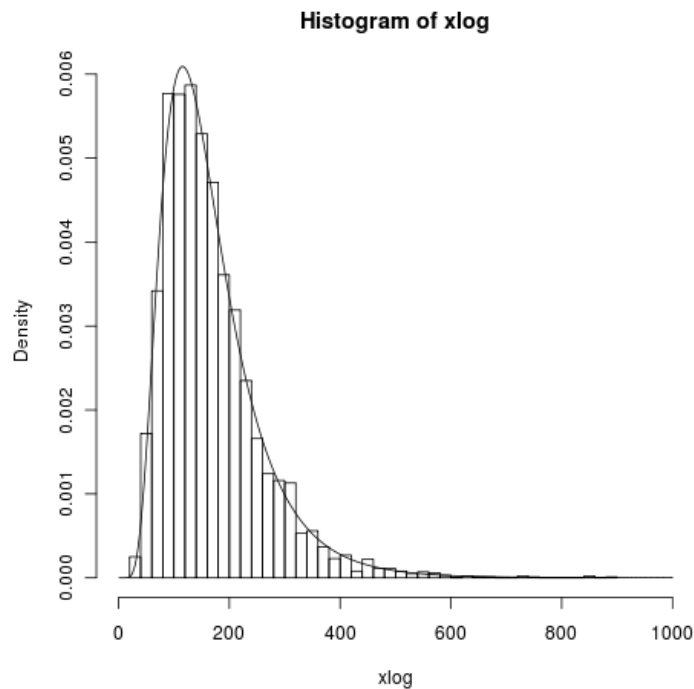
Med resultatet fra opgave 1 udregnes $E(Y^2)$ til

$$\exp \sigma^2 + 2\mu$$

Og vi har lige udregnet $(E(Y))^2$

$$(E(Y))^2 = \exp 2\mu + 2\sigma^2$$

$$Var(Y) = \exp \sigma^2 + 2\mu - \exp 2\mu + 2\sigma^2 = (\exp \sigma^2 - 1) \exp 2\mu + \sigma^2$$



Figur 1: Simulation af 5000 observationer fra den logaritmiske normalfordeling med parametre (5, 0.25).

Del 2

Vi har simuleret 5000 observationer fra den logaritmiske normalfordeling med parametre (5, 0.25). Plottet kan ses i figur 1.

Ved brug af **R** har vi fundet stikprøvens gennemsnit, varians, spredning og median. Disse ses i udskriften fra programmet i kodeboks 1.

Ved at bruge formlerne fra Del 1, har vi regnet de teoretiske værdier for gennemsnit, varians, spredning og median. De ses herunder.

$$\begin{aligned}
 E(S) &= e^{\mu + \sigma^2/2} \\
 &= e^{5 + 0.125} \\
 &= 168.1741
 \end{aligned} \tag{1}$$

$$\begin{aligned}
 Var(S) &= (e^{\sigma^2} - 1)e^{2\mu + \sigma^2} \\
 &= (e^{0.25} - 1)e^{10 + 0.25} \\
 &= 8032.96
 \end{aligned} \tag{2}$$

```

1 [1] "Assignment 7"
2 [1] "Middelvaerdi"
3 [1] 167.1192
4 [1] "Varians"
5 [1] 7976.314
6 [1] "Standardafvigelse"
7 [1] 89.31021
8 [1] "Median"
9 [1] 147.5676

```

Kodeboks 1: Udskrift fra R-program

$$\begin{aligned}\sqrt{\text{Var}(S)} &= \sqrt{8032.96} \\ &= 89.62678\end{aligned}\tag{3}$$

$$\begin{aligned}F^{-1}(0.5) &= e^{\mu} \\ &= e^5 \\ &= 148.4132\end{aligned}\tag{4}$$

Det ses, at de teoretiske værdier ligger meget tæt på dem som er blevet udregnet fra stikprøven. Dette kan også ses ved at kigge på plottet i figur 1, hvor tæthedsfunktionen følger de observerede værdier. Det var også at forvente, når man har taget så stor en stikprøve.

Del 3

I undersøgelsen indgår der 1079 mænd og 1145 kvinder. Vi ser i figur 2 og 3, at den originale data ikke er normalfordelt, men når vi transformerer denne logaritmisk, har vi en normalfordeling.

Vi siger derfor, at det logaritmisk transformerede indtag af A-vitamin for kvinder og mænd, er normalfordelt.

Vi estimerer $\mu = \bar{x} = 7.361356$ og $\sigma = s = 0.4486289$. Vi har den teoretiske fordeling

$$X \sim N(\mu, \sigma^2)\tag{5}$$

og får den estimerede fordeling, ved at indsætte vores estimer til μ og σ i (5).

$$X \sim N(7.361356, 0.2012679)\tag{6}$$

I figur 4 ses hvordan de observerede værdier forholder sig til den estimerede fordeling.

Vi undersøger nu, hvorvidt der er forskel på mænd og kvinders indtag af A-vitamin. Vi bruger R til at lave en t-test på vores data. Vi får outputet vist i kodeboks 2. Det 95 procents konfidensinterval aflæses direkte fra koden, mens estimatet for differensen fra t-testen fås ved at trække de beregnede middelværdier fra hinanden.

```

1      Welch Two Sample t-test
2
3      data:  logavitM and logavitK
4      t = 13.0815, df = 2206.488, p-value < 2.2e-16
5      alternative hypothesis: true difference in means is not equal to 0
6      95 percent confidence interval:
7      0.2041469 0.2761477
8      sample estimates:
9      mean of x mean of y
10     7.484993  7.244845

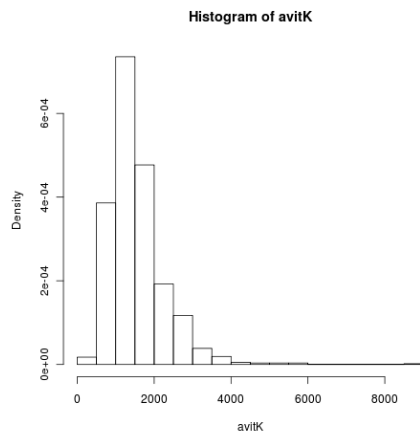
```

Kodeboks 2: Udskrift fra T-test

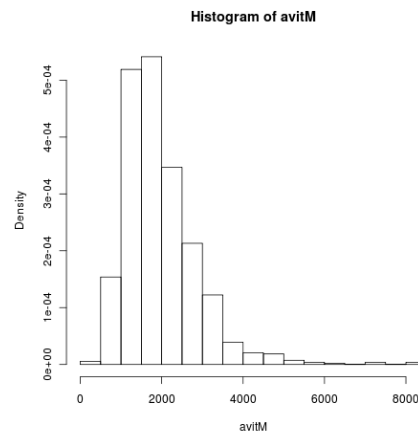
Vores 95-procents konfidensinterval for differensen mellem middelværdierne er $[0.2041469, 0.2761477]$. Dvs. at der er 95% sandsynlighed for at differensen mellem to middelværdier ligger i dette interval. Vores estimat af differens af middelværdierne er $7.484993 - 7.244845 = 0.240148$, hvor vi har at $0.240148 \in [0.2041469, 0.2761477]$. Mænd og kvinders logaritmiske indtag af A-vitamin kan derfor godt være det samme.

Ligesom vi kan finde medianen i den logaritmiske normalfordeling ved e^μ , kan vi bruge samme transformation for grænserne og differensen. På grund af regneregler for eksponentialfunktionen vil vi automatisk få et forhold når vi transformerer en differens, da $e^{a-b} = \frac{e^a}{e^b}$. Hvis vi transformerer estimatet og grænserne for konfidensintervallet for differensen, får vi et estimat og grænser for et konfidensinterval for forholdet, mellem medianerne i fordelingen af indtaget af A-vitamin, for mænd og kvinder. Konfidensintervallet for forholdet bliver $[1.226478, 1.318043]$ og differensen transformerer til 1.271437. Vi har at $1.271437 \in [1.226478, 1.318043]$.

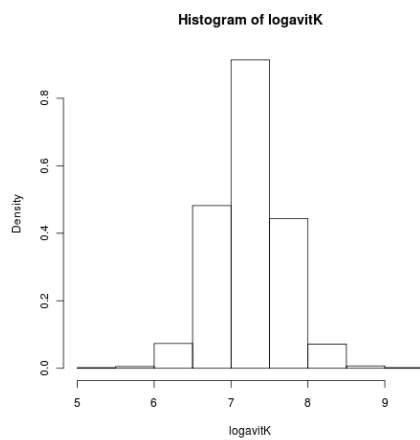
Da $1 \notin [1.226478, 1.318043]$ kan tæller og nævner i forholdet ikke være identiske, og der er derfor mindst 95% sandsynlighed for at kvinder og mænds indtag af A-vitamin er forskellige.



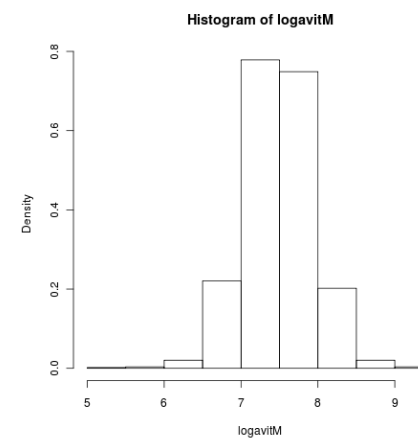
(a) A-vitaminindtag for kvinder



(b) A-vitaminindtag for mænd

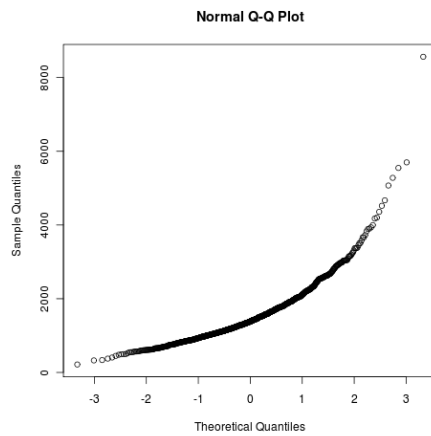


(c) Logaritmisk tranformeret A-vitaminindtag for kvinder

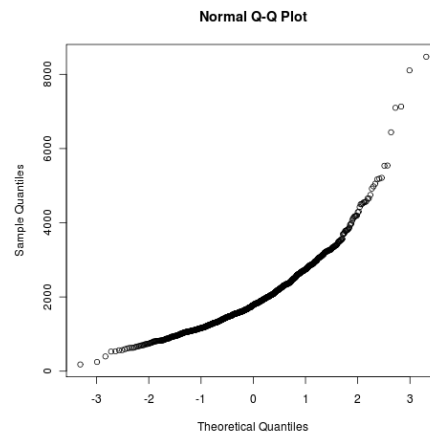


(d) Logaritmisk tranformeret A-vitaminindtag for mænd

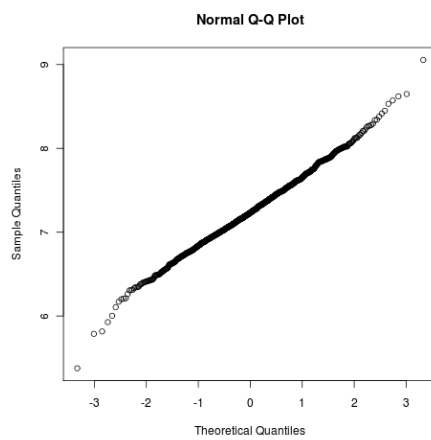
Figur 2: Kvinder og mænds indtag af A-vitamin. Den originale data ligner ikke en normalfordeling, men det gør den transformerede.



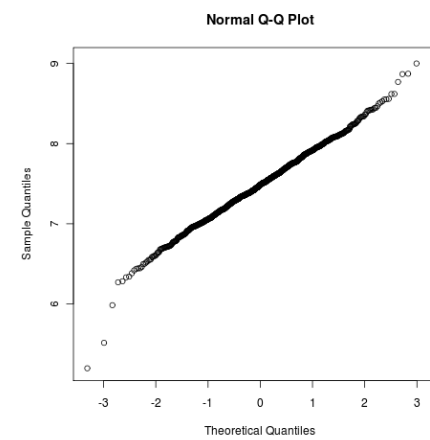
(a) QQ-plot for kvinders indtag af A-vitaminindtag



(b) QQ-plot for mænds A-vitaminindtag

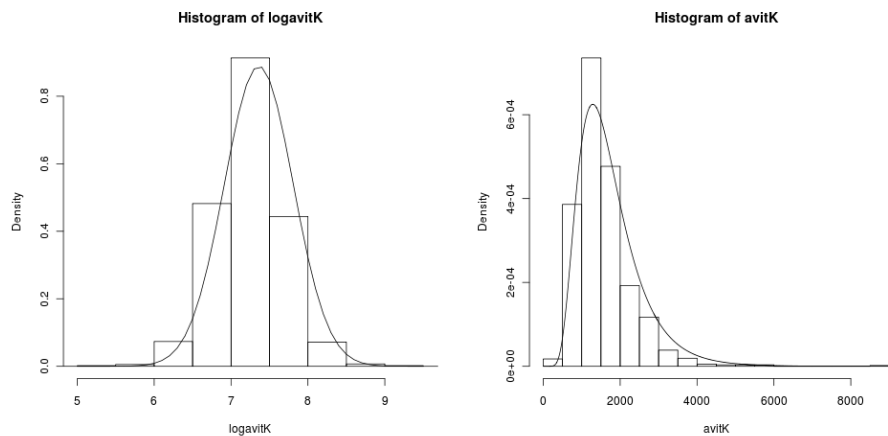


(c) QQ-plot af det logaritmisk tranformerede A-vitaminindtag for kvinder

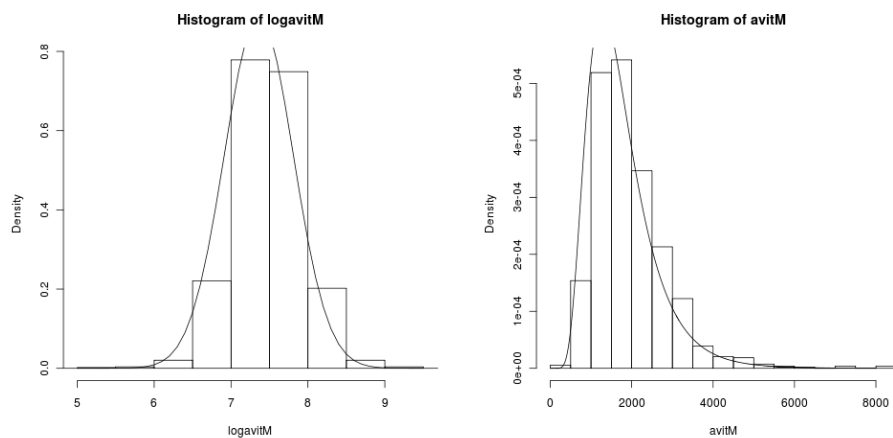


(d) QQ-plot af det logaritmisk tranformerede A-vitaminindtag for mænd

Figur 3: QQ-plots for indtag af A-vitamin. Det ses, at når data bliver logaritmisk transformeret, får vi en normalfordeling, da punktene, i disse qq-plots, følger en diagonal.



(a) Logaritmisk A-vitaminindtag for kvinder med estimeret tæthed (b) A-vitaminindtag for kvinder med estimeret tæthed



(c) Logaritmisk A-vitaminindtag for mænd med estimeret tæthed (d) A-vitaminindtag for mænd med estimeret tæthed

Figur 4: Det ses, at de observerede værdier nogenlunde følger den estimerede tæthed. Man kan også se, at der er flere kvinder end mænd i undersøgelsen, da værdierne for mænd ligger lidt under de estimerede værdier og vice versa.

A Source code

A.1 main.r

```
1  #!/usr/bin/env Rscript
2
3  # R statistics for the course "Sandsynlighedsregning og statistik"
4  # at the Department of Mathematical Sciences, University of Copenhagen
5
6  #####
7  #
8  # Project 2
9  # Part 2
10 #
11 #####
12
13 f <- function(y, mu, sigma) {
14
15     return (1/(sqrt(2 * pi * sigma^2) * y) * exp(-((log( y ) - mu)^2 / (2 * sigma ^ 2))))
16
17 }
18
19 # Assignment 1
20 assignment1 <- function() {
21
22     set.seed(42);
23
24     mymu <- 5;
25     mysigma <- sqrt(0.25);
26
27     xlog <- rlnorm(5000, meanlog=mymu, sdlog=mysigma);
28
29     # Set the filename for the plot and plot the graph
30     name <- "plot_1.png";
31     png(name);
32     hist(xlog, density=NULL, prob=TRUE, nclass=50, xlim=range(0,1000,100));
33
34     yval <- seq(1, 1000, 1);
35     fval <- f(yval, mu=mymu, sigma=mysigma);
36
37     points(yval, fval, type="l");
38
39     dev.off();
40
41     print("Assignment 7");
42     print("Median");
43     print(median(xlog));
44     print("Middelvaerdi");
45     print(mean(xlog));
46     print("Standardafvigelse");
47     print(sd(xlog));
48     print("Varians");
49     print(var(xlog));
50     print("");
51 }
52
53
54 # Assignment 3
55 assignment2 <- function() {
56
57     avitdata <- read.table("avit.txt", header=TRUE);
58     attach(avitdata);
59
60     avitM <- avit[sex==1];
61     avitK <- avit[sex==2];
62
63     logavitM <- log(avitM);
64     logavitK <- log(avitK);
65
66     print("Number of men in the experiment:");
67     print(length(avitM));
68     print("Number of females in the experiment:");
```

```

69   print(length(avitK));
70
71   #   print(length(avitK)+length(avitM))
72
73   name <- "plot_3_avitM.png";
74   png(name);
75   hist(avitM, prob=TRUE);
76
77   name <- "plot_3_avitK.png";
78   png(name);
79   hist(avitK, prob=TRUE);
80
81   name <- "plot_3_logavitK.png";
82   png(name);
83   hist(logavitK, prob=TRUE);
84
85   name <- "plot_3_logavitM.png";
86   png(name);
87   hist(logavitM, prob=TRUE);
88
89   name <- "plot_3_avitM_qq.png";
90   png(name);
91   qqnorm(avitM);
92
93   name <- "plot_3_avitK_qq.png";
94   png(name);
95   qqnorm(avitK);
96
97   name <- "plot_3_logavitM_qq.png";
98   png(name);
99   qqnorm(logavitM);
100
101   name <- "plot_3_logavitK_qq.png";
102   png(name);
103   qqnorm(logavitK);
104
105   logavit <- log(avit);
106   mloga <- mean(logavit);
107   print("Mean of the log");
108   print(mloga);
109
110   sdloga <- sd(logavit);
111   print("log standard deviation:");
112   print(sdloga);
113
114   dev.off();
115
116   # Exercise 11
117   # logavitM og logavitK ~ N(mean(log(avit)), sd(log(avit)))
118
119   # Exercise 12
120
121   name <- "plot_3_logavitM_with_normal.png";
122   png(name);
123   hist(logavitM, prob=TRUE);
124
125   yval <- seq(1,12,0.1);
126   fval <- dnorm(yval, mean=mloga, sd=sdloga);
127
128   points(yval, fval, type="l");
129
130   dev.off();
131
132   name <- "plot_3_avitM_with_lognormal.png";
133   png(name);
134   hist(avitM, prob=TRUE);
135
136   yval <- seq(0,8000,1);
137   fval <- f(yval, mloga, sdloga);
138
139   points(yval, fval, type="l");
140
141   dev.off();
142

```

```

143
144   name <- "plot_3_logavitK_with_normal.png";
145   png(name);
146   hist(logavitK, prob=TRUE);
147
148   yval <- seq(1,26,0.1);
149   fval <- dnorm(yval, mean=mloga, sd=sdloga);
150
151   points(yval, fval, type="l");
152
153   dev.off();
154
155   name <- "plot_3_avitK_with_lognormal.png";
156   png(name);
157   hist(avitK, prob=TRUE);
158
159   yval <- seq(0,8000,1);
160   fval <- f(yval, mloga, sdloga);
161
162   points(yval, fval, type="l");
163
164   dev.off();
165
166   t.test(logavitM, logavitK, var.equal=TRUE);
167
168 }
169
170 #####
171
172 #####
173
174 # Main
175
176 assignment1();
177 assignment2();
178

```