

# TMA4315: Project 2

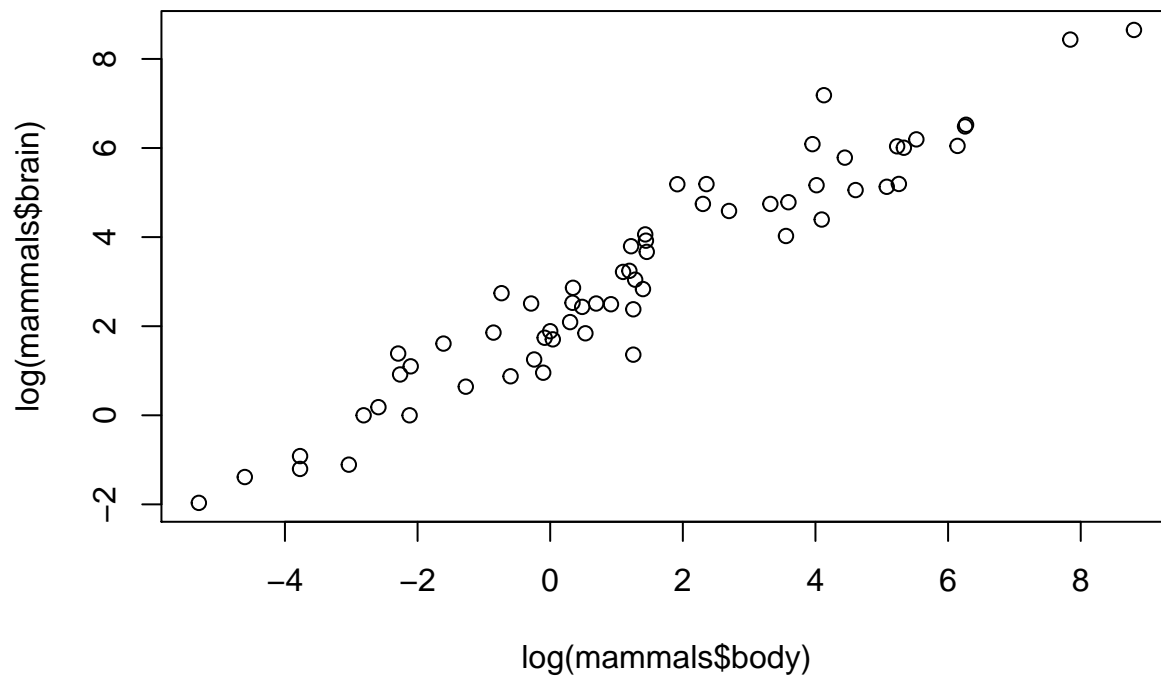
jototlan@stud.ntnu.no (10018), martigtu@stud.ntnu.no (10037)

## Problem 1

```
mammals <- read.table(  
  "https://www.math.ntnu.no/~jarlet/statmod/mammals.dat",  
  header=T)
```

a)

```
plot(log(mammals$body), log(mammals$brain)) # Seems pretty linear.
```



A log-log plot of the brain mass against body mass seems to reveal a linear trend. We thus fit the following model:

```
mod0 <- lm(log(brain) ~ log(body), data = mammals)  
summary(mod0)
```

```
##  
## Call:  
## lm(formula = log(brain) ~ log(body), data = mammals)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max
```

```
## -1.71550 -0.49228 -0.06162  0.43597  1.94829
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.13479    0.09604   22.23  <2e-16 ***
## log(body)    0.75169    0.02846   26.41  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6943 on 60 degrees of freedom
## Multiple R-squared:  0.9208, Adjusted R-squared:  0.9195
## F-statistic: 697.4 on 1 and 60 DF,  p-value: < 2.2e-16
```

b)

```
is.human = ifelse(mammals$species == "Human", 1, 0)
mammals$is.human = as.factor(is.human)

mod1 <- lm(log(brain) ~ log(body) + is.human, data = mammals)
summary(mod1)

##
## Call:
## lm(formula = log(brain) ~ log(body) + is.human, data = mammals)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.68392 -0.46764 -0.02398  0.47237  1.64949
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.11500    0.09030   23.421  < 2e-16 ***
## log(body)    0.74228    0.02687   27.622  < 2e-16 ***
## is.human1    2.00691    0.66083    3.037  0.00356 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6511 on 59 degrees of freedom
## Multiple R-squared:  0.9315, Adjusted R-squared:  0.9292
## F-statistic: 401.1 on 2 and 59 DF,  p-value: < 2.2e-16
```

Let  $\hat{\beta} = [\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2]^T$  be the coefficient estimates given in the summary above. Then the estimated effect on brain mass from being a human is  $\hat{\beta}_2 \approx 2.0069072$ . Since we have used a log-transform on both the brain mass and body mass, humans will according to the model be larger by a factor of  $e^{\hat{\beta}_2} = 7.4402704$ .

We use the notation  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  to represent the linear model. Here,  $\mathbf{X}$  is the  $n \times p$  design matrix, where  $n$  is the number of observations and  $p$  is the number of parameters used in the model. As usual,  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ . This (along with the other usual assumptions [how much detail is required here??](#)) gives the well known result:

$$\hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}).$$

Now we want to perform the hypothesis test

$$H_0 : \beta_2 = 0 \quad \text{vs.} \quad H_1 : \beta_2 > 0.$$

Under  $H_0$ , we obtain that (we also index from 0 in the design matrix)

$$\frac{\hat{\beta}_2}{\sigma \sqrt{(X^T X)_{2,2}^{-1}}} \sim \mathcal{N}(0, 1).$$

Combining this with the fact that

$$\frac{(n-p)s^2}{\sigma^2} \sim \chi_{n-p}^2,$$

where  $s^2 = RSS/(n-p)$ , we obtain the test statistic

$$\frac{\hat{\beta}_2}{s \sqrt{(X^T X)_{2,2}^{-1}}} \sim t_{n-p},$$

under  $H_0$ . We perform the calculations in R:

```
n <- nrow(mammals)
p <- 3
beta.2 <- mod1$coefficients[3]
s <- sqrt(deviance(mod1)/(n-p))
X <- model.matrix( ~ log(body) + is.human, data = mammals)
XtX.inv <- solve(t(X) %*% X)

T.stat <- beta.2/(s*sqrt(XtX.inv[3,3]))
p.val <- pt(T.stat, n - p, lower.tail = F)
p.val

## is.human1
## 0.001777696
```

The calculated p-value is 0.0017777.

c)

We now consider the linear model with only two parameters,  $\beta_0$  and  $\beta_1$ . Let  $Y_h = \beta_0 + \beta_1 x_h$  be the stochastic variable from which the log of the human brain mass is realized and  $\hat{Y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h$  be the corresponding estimator. Then we can find the pivotal quantity

$$T = \frac{Y_h - \hat{Y}_h}{s \sqrt{1 + 1/n + \frac{(x_h - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t_{n-2}.$$

We refer to the good old [subject-pages](#) (simple linear regression/prediction and prediction intervals in simple linear regression) for this result. Thus, we can find the one-sided prediction interval:

$$P(T \leq k) = 1 - \alpha \implies k = t_{n-2, \alpha}.$$

Rearranging, we arrive at

$$P\left(Y_h \leq t_{n-2, \alpha} \cdot s \sqrt{1 + 1/n + \frac{(x_h - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} + \hat{Y}_h\right) = 1 - \alpha$$

We denote the right hand side of the inequality above by  $U$  and calculate it with the observed values below

```
# HMMM trenegr ikke gjøre dette da egentlig?
x.h <- mammals$body[mammals$species == "Human"]
mammals.reduced <- mammals[mammals$species != "Human", ]
mod0.reduced <- lm(log(brain) ~ log(body), data = mammals.reduced)
summary(mod0.reduced)
pred.h <- predict(mod0.reduced, newdata = data.frame('body' = x.h))
s <- sqrt(deviance(mod0.reduced))
x.bar <- mean(mammals.reduced$body)
n <- nrow(mammals.reduced)
alpha = 0.05

U <- qt(1 - alpha, n) * s * sqrt(1 + 1/n + (x.h - x.bar)^2/sum((mammals.reduced$body - x.bar)^2) + pred
```

very uncertain on this one.

d) For a gamma-distributed random variable, the pdf takes the form

$$f(x | a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}.$$

Using the parametrization  $\mu = \frac{a}{b}$  and  $\nu = a$ , we construct the GLM with a log-link as follows. Let the mammalian brain size given body size be given as

$$y_i \sim \text{Gamma}(\mu_i, \nu),$$

where

$$-\frac{1}{\mu_i} = \mathbf{x}_i^T \boldsymbol{\beta} =: \eta_i.$$

Next, we fit the model:

```
mod.glm <- glm(brain ~ log(body) + is.human, family = Gamma(link = "log"), data = mammals)
summary(mod.glm)

##
## Call:
## glm(formula = brain ~ log(body) + is.human, family = Gamma(link = "log"),
##      data = mammals)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.4464  -0.6099  -0.2276   0.2725   1.8835
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.32733    0.10298  22.601  <2e-16 ***
## log(body)    0.74193    0.03064  24.212  <2e-16 ***
## is.human1    1.79601    0.75356   2.383   0.0204 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Gamma family taken to be 0.5512612)
##
##      Null deviance: 310.710  on 61  degrees of freedom
## Residual deviance:  25.849  on 59  degrees of freedom
## AIC: 523.38
```

```
##
## Number of Fisher Scoring iterations: 5
```

e)

We want to test whether the following relationship holds:

$$Y = Y_0 M^{3/4},$$

where  $Y$  is the brain mass,  $Y_0$  is a constant and  $M$  is the brain mass. Since this is equivalent to testing

$$\ln(Y) = \ln(Y_0) + \frac{3}{4} \ln(M),$$

we can, for the model in (b), simply perform the hypothesis test:

$$H_0 : \beta_1 = \frac{3}{4} \quad \text{vs.} \quad \beta_1 \neq \frac{3}{4}.$$

We follow the standard framework for a linear hypothesis test:

```
# Wald test:
C <- matrix(c(0, 1, 0), nrow = 1)
d <- 3/4
r <- 1
p <- 3
n <- nrow(mammals)
beta1 <- mod0$coefficients[2]
s2 <- deviance(mod0)
X <- model.matrix(~ log(body), data = mammals)
XtX.inv <- solve(t(X) %*% X)

F.stat <- (beta1-3/4)^2/(s2*XtX.inv[2,2])
p.val <- pf(F.stat, r, n - p, lower.tail = F)
p.val
```

```
## log(body)
## 0.9939246
```

**How to use LRT test on linear hypothesis? very large p-value!** For a generalized linear model, the Wald statistic can be written as

$$w = (C\hat{\beta} - d)^T [CF^{-1}(\hat{\beta})C^T]^{-1} (C\hat{\beta} - d),$$

which is asymptotically  $\chi^2$ -distributed with  $r = \text{rank}(C)$  degrees of freedom. We compute its value:

```
beta <- as.vector(mod1$coefficients)
denom <- solve(C %*% vcov(mod1) %*% t(C))
w <- (C %*% beta - d)^2*denom

p.val <- pchisq(w, r, lower.tail = F)
p.val
```

```
##           [,1]
## [1,] 0.7737976
```

We perform LRT tests by using an offset term.

```
mod0.o <- lm(log(brain) ~ 1, offset = 3/4*log(body), data = mammals)
anova(mod0.o, mod0, test= "Chisq")
```

```
## Analysis of Variance Table
##
## Model 1: log(brain) ~ 1
## Model 2: log(brain) ~ log(body)
##   Res.Df    RSS Df Sum of Sq Pr(>Chi)
## 1      61 28.924
## 2      60 28.923  1 0.0016912  0.9528
```

f)

## Problem 2

### Assumptions

In this problem we apply ordinal multinomial regression to data from Norway Chess 2021.

```
df <- read.csv('data/Norway\ Chess\ 2021.csv')
```

The response variable  $y_i$  is the outcome of the  $i$ 'th match. This can be considered an ordered categorical variable

$$y_i = \begin{cases} 1 & , \text{ white win} \\ 2 & , \text{ draw} \\ 3 & , \text{ black win,} \end{cases}$$

which may depend on relative strength of different players, which player plays white and black and the type of game played. The response can be determined by an underlying latent variable  $u_i$ , given by

$$u_i = -\mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i,$$

where  $\epsilon_i \stackrel{iid}{\sim} f$ , where  $f$  is some standard distribution with cdf  $F$ . In this model, the event  $y_i = r$  occurs if  $\theta_{r-1} < u_i \leq \theta_r$  for some parameters  $\{\theta_i\}_{i=0}^3$  satisfying

$$-\infty = \theta_0 < \theta_1 < \theta_2 < \theta_3 = \infty.$$

It follows that

$$P(y_i \leq r) = P(u_i \leq \theta_r) = P(\epsilon_i \leq \theta_r + \mathbf{x}_i^T \boldsymbol{\beta}) = F(\theta_r + \mathbf{x}_i^T \boldsymbol{\beta}),$$

so the probability of observing a particular outcome of the  $i$ 'th match becomes

$$\begin{aligned} \pi_{ir} = P(y_i = r) &= P(y_i \leq r) - P(y_i \leq r-1) \\ &= F(\theta_r + \mathbf{x}_i^T \boldsymbol{\beta}) - F(\theta_{r-1} + \mathbf{x}_i^T \boldsymbol{\beta}). \end{aligned}$$

This means that our model returns that white wins whenever  $u_i \leq \theta_1$ , draw if  $\theta_1 < u_i \leq \theta_2$  and black win for  $u_i > \theta_2$ .

### Models

#### Propositional odds model / Cumulative Logit

$$F(x) = \frac{e^x}{1 + e^x}, \quad \epsilon_i \sim \text{Logistic}(0, 1)$$

### Cumulative Probit

$$F(x) = \Phi(x), \quad \epsilon_i \sim N(0, 1)$$

### R

```
library(VGAM)
```

```
## Loading required package: stats4
```

```
## Loading required package: splines
```

```
fit <- vglm(y ~ factor(white) + factor(black),  
           family=cumulative(parallel = TRUE, link="logitlink"), data=df)
```

```
# P(u <= theta_1), P(u <= theta_2)  
p.less_or_equal <- plogis(predict(fit, df))
```

```
stats <- cbind('white'=df$white, 'black'=df$black,
              'P(white)'=round(p.less_or_equal[,1],2),
              'P(draw)'=round(p.less_or_equal[,2]-p.less_or_equal[,1],2),
              'P(black)'=round(1-p.less_or_equal[,2],2),
              'outcome'=c('white','draw','black')[df$y])
stats
```

##	white	black	P(white)	P(draw)	P(black)	outcome
## 1	"firouzja"	"carlsen"	"0.33"	"0.47"	"0.2"	"draw"
## 2	"firouzja"	"carlsen"	"0.33"	"0.47"	"0.2"	"draw"
## 3	"tari"	"rapport"	"0.1"	"0.38"	"0.52"	"black"
## 4	"nepomniachtchi"	"karjakin"	"0.36"	"0.46"	"0.17"	"white"
## 5	"nepomniachtchi"	"firouzja"	"0.31"	"0.48"	"0.22"	"draw"
## 6	"nepomniachtchi"	"firouzja"	"0.31"	"0.48"	"0.22"	"white"
## 7	"carlsen"	"tari"	"0.7"	"0.25"	"0.05"	"draw"
## 8	"carlsen"	"tari"	"0.7"	"0.25"	"0.05"	"white"
## 9	"karjakin"	"rapport"	"0.35"	"0.47"	"0.19"	"draw"
## 10	"karjakin"	"rapport"	"0.35"	"0.47"	"0.19"	"draw"
## 11	"firouzja"	"karjakin"	"0.55"	"0.36"	"0.09"	"draw"
## 12	"firouzja"	"karjakin"	"0.55"	"0.36"	"0.09"	"black"
## 13	"tari"	"nepomniachtchi"	"0.09"	"0.36"	"0.55"	"draw"
## 14	"tari"	"nepomniachtchi"	"0.09"	"0.36"	"0.55"	"black"
## 15	"rapport"	"carlsen"	"0.39"	"0.45"	"0.16"	"draw"
## 16	"rapport"	"carlsen"	"0.39"	"0.45"	"0.16"	"draw"
## 17	"tari"	"karjakin"	"0.13"	"0.43"	"0.44"	"draw"
## 18	"tari"	"karjakin"	"0.13"	"0.43"	"0.44"	"black"
## 19	"carlsen"	"nepomniachtchi"	"0.71"	"0.24"	"0.05"	"draw"
## 20	"carlsen"	"nepomniachtchi"	"0.71"	"0.24"	"0.05"	"white"
## 21	"rapport"	"firouzja"	"0.55"	"0.36"	"0.09"	"white"
## 22	"firouzja"	"nepomniachtchi"	"0.44"	"0.42"	"0.13"	"white"
## 23	"tari"	"carlsen"	"0.06"	"0.28"	"0.66"	"draw"
## 24	"tari"	"carlsen"	"0.06"	"0.28"	"0.66"	"white"
## 25	"rapport"	"karjakin"	"0.61"	"0.32"	"0.07"	"white"
## 26	"carlsen"	"firouzja"	"0.74"	"0.22"	"0.04"	"white"
## 27	"rapport"	"tari"	"0.49"	"0.4"	"0.11"	"white"
## 28	"karjakin"	"nepomniachtchi"	"0.32"	"0.47"	"0.2"	"draw"
## 29	"karjakin"	"nepomniachtchi"	"0.32"	"0.47"	"0.2"	"white"
## 30	"firouzja"	"nepomniachtchi"	"0.44"	"0.42"	"0.13"	"white"
## 31	"tari"	"carlsen"	"0.06"	"0.28"	"0.66"	"black"
## 32	"rapport"	"karjakin"	"0.61"	"0.32"	"0.07"	"white"
## 33	"nepomniachtchi"	"tari"	"0.26"	"0.48"	"0.26"	"black"
## 34	"carlsen"	"rapport"	"0.73"	"0.23"	"0.04"	"white"
## 35	"karjakin"	"firouzja"	"0.36"	"0.46"	"0.18"	"black"
## 36	"tari"	"firouzja"	"0.1"	"0.39"	"0.51"	"black"
## 37	"carlsen"	"karjakin"	"0.79"	"0.18"	"0.03"	"white"
## 38	"rapport"	"nepomniachtchi"	"0.51"	"0.39"	"0.11"	"draw"
## 39	"rapport"	"nepomniachtchi"	"0.51"	"0.39"	"0.11"	"black"
## 40	"firouzja"	"rapport"	"0.47"	"0.41"	"0.12"	"white"
## 41	"nepomniachtchi"	"carlsen"	"0.19"	"0.47"	"0.34"	"draw"
## 42	"nepomniachtchi"	"carlsen"	"0.19"	"0.47"	"0.34"	"black"
## 43	"karjakin"	"tari"	"0.31"	"0.48"	"0.21"	"draw"
## 44	"karjakin"	"tari"	"0.31"	"0.48"	"0.21"	"white"