# STK2100
# Oblig 1

**Egil Furnes**
**Studentnummer: 693784**

# Oppgave 1

## a)

### i)

The only model that can be written as a linear model as is, is model 2.

$$2. \quad Y = \beta_0 + \frac{\beta_1}{x} + \beta_2 x^2 + \epsilon$$

### ii)

For model 4 we can fix $\beta_2$ to a constant, say $c$ and get

$$4. \quad \beta_0 + \beta_1 x^c + \epsilon$$

### iii)

For model 5 we can log-transform such as this:

$$Y = \beta_0 x^{\beta_1} \varepsilon$$

$$\log(Y) = \log(\beta_0 x^{\beta_1} \varepsilon)$$

$$\log(Y) = \log(\beta_0) + \log(x^{\beta_1}) + \log(\varepsilon)$$

$$\log(Y) = \log(\beta_0) + \beta_1 \log(x) + \log(\varepsilon)$$

**b)**

1. $X = \begin{bmatrix} \frac{1}{1+x_i} & x_i^{1/2} \end{bmatrix}$, $\qquad\qquad\qquad \beta = \begin{bmatrix} \beta_0 \\ \beta_2 \end{bmatrix}$

2. $X = \begin{bmatrix} 1 & \frac{1}{x_i} & x_i^2 \end{bmatrix}$, $\qquad\qquad\qquad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$

3. $X = \begin{bmatrix} 1 & x_i & x_i^2 \end{bmatrix}$, $\qquad\qquad\qquad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$

4. $X = \begin{bmatrix} 1 & x_i \end{bmatrix}$, $\qquad\qquad\qquad \beta = \begin{bmatrix} \beta_0 \\ \beta_1' \end{bmatrix}$

5. $X = \begin{bmatrix} 1 & \log(x_i) \end{bmatrix}$, $\qquad\qquad\qquad \beta = \begin{bmatrix} \beta_0' \\ \beta_1 \end{bmatrix}$

# Oppgave 2

## a)

```r
# set seed for replication!
set.seed(1814)

# loading packages
library(tidyverse)

# reading data
nuclear <- read_delim("nuclear.dat")

# a)

# fitting the data using linear regression
lm1 <- lm(log(cost)~., nuclear)
lm1 %>% summary()

# creating a 95 percent confidence interval
confint(lm1, level = 0.95)[c("t1","t2","bw"),]
```

A 95% confidence interval is then found to be the following.

```
           2.5 %      97.5 %
t1 -0.041123331 0.05162679
t2 -0.003949911 0.01516185
bw -0.184211843 0.25780411
```

## b)

```r
# b)

# creating a new dataframe
df <- tibble(
  date = 70.0,
  t1 = 13,
  t2 = 50,
  cap = 800,
  pr = 1,
  ne = 0,
  ct = 0,
  bw = 1,
  cum.n = 8,
  pt = 1
)
```

```
17  # predicting on data using linear regression
18  pred1 <- predict(lm1, newdata = df, interval = "prediction", level =
        0.95)
19
20  # retrieving the coefficients
21  yhat <- pred1[1,"fit"]
22  lwry <- pred1[1,"lwr"]
23  upry <- pred1[1,"upr"]
24
25  # transforming y to find z
26  zfit <- exp(yhat)
27  zlwr <- exp(lwry)
28  zupr <- exp(upry)
29
30  # saving coefficients for z
31  predz <- data.frame(
32    fit = zfit,
33    lwr = zlwr,
34    upr = zupr
35  )
36
37  # presenting my findings
38  print(pred1)
39  print(predz)
```

The 95% prediction interval with the cost $Z$ is then found to be the following.

```
> print(pred1)
        fit      lwr      upr
1 5.964135 5.394248 6.534022
> print(predz)
        fit      lwr      upr
1 389.2163 220.1366 688.1607
```

## c)

```
1  # c)
2
3  # individual t-tests
4  sum_lm1 <- summary(lm1)
5  print(sum_lm1$coefficients[c("t1", "t2", "bw"),
6                             c("Estimate", "Pr(>|t|)"),
7                             drop = FALSE])
8
9  # joint F-tests
10 lm1red <- lm(log(cost)~date+cap+pr+ne+ct+cum.n+pt, data = nuclear)
11 anova(lm1red, lm1)
```

Find the output for the individual t-test, where we find that the p-value is larger than $0.5$ for all predictors t1, t2, and bw, such that we fail to reject the null-hypothesis $H_0 : \beta_j = 0$ that the predictors are significant.

```
        Estimate  Pr(>|t|)
t1 0.005251730 0.8160981
t2 0.005605968 0.2359862
bw 0.036796131 0.7326075
```

For the joint F-test, we find a p-value of $0.5173 > 0.5$ where we fail to reject the $H_0$ at a 5% confidence level.

```
Analysis of Variance Table

Model 1: log(cost) ~ date + cap + pr + ne + ct + cum.n + pt
Model 2: log(cost) ~ date + t1 + t2 + cap + pr + ne + ct + bw + cum.n +
    pt
  Res.Df     RSS Df Sum of Sq     F Pr(>F)
1     24 0.67195
2     21 0.60443  3   0.06752 0.782 0.5173
```

### d)

```
1  # d)
2
3  library(leaps)
4
5  fwd1 <- regsubsets(log(cost) ~ ., data = nuclear, method = "forward"
     , nvmax = 10)
6  sum1 <- summary(fwd1)
7  sum1$outmat %>% as_tibble()
8
9  bic1 <- sum1$bic
10 bicb <- names(which(sum1$which[which.min(bic1), ]))[-1]
11
12 all1 <- names(nuclear)[-1]
13 aic1 <- which.max(sum1$adjr2)  # Fixed: max adj R²
14 aicb <- names(which(sum1$which[aic1, ]))[-1]
15
16 bicb
17 aicb
18
19 fitaic <- lm(log(cost) ~ ., data = nuclear[, c("cost", aicb)])
20 fitbic <- lm(log(cost) ~ ., data = nuclear[, c("cost", bicb)])
21
22 aic2 <- AIC(fitaic)
23 bic2 <- BIC(fitbic)
```

```
24
25  aic2
26  bic2
```

The output of the order matrix is below, which gives the following order of variables included: `cap` →`date` →`ne` →`ct` →`bw` →`pr` →`t2` →`t1`.

```
# A tibble: 10 × 10
    date  t1    t2    cap   pr    ne    ct    bw
    <chr> <chr> <chr> <chr> <chr> <chr> <chr> <chr>
 1  " "   " "   " "   " "   " "   " "   " "   " "
 2  " "   " "   " "   "*"   " "   " "   " "   " "
 3  "*"   " "   " "   "*"   " "   " "   " "   " "
 4  "*"   " "   " "   "*"   " "   "*"   " "   " "
 5  "*"   " "   " "   "*"   " "   "*"   "*"   " "
 6  "*"   " "   " "   "*"   " "   "*"   "*"   " "
 7  "*"   " "   " "   "*"   " "   "*"   "*"   "*"
 8  "*"   " "   " "   "*"   "*"   "*"   "*"   "*"
 9  "*"   " "   "*"   "*"   "*"   "*"   "*"   "*"
10  "*"   "*"   "*"   "*"   "*"   "*"   "*"   "*"
#   2 more variables: cum.n <chr>, pt <chr>
```

Using forward selection I end up with variables `[1] "date" "cap" "ne" "ct" "pt"` selected on $BIC$, and with `"date" "t2" "cap" "pr" "ne" "ct" "bw" "cum.n" "pt"` selected on $AIC$. This yields respectively $AIC$ and $BIC$:

```
> aic2
[1] -14.11792
> bic2
[1] -5.034235
```

### e)

```
1   # e)
2
3   bwd2 <- regsubsets(log(cost) ~ ., data = nuclear, method = "backward
        ", nvmax = 10)
4   sum2 <- summary(bwd2)
5   sum2$outmat %>% as_tibble()
6
7   bic2 <- sum2$bic
8   bicc <- names(which(sum2$which[which.min(bic2), ]))[-1]
9
10  all2 <- names(nuclear)[-1]
11  aic2 <- which.max(sum2$adjr2)  # Fixed: max adj R²
12  aicc <- names(which(sum2$which[aic2, ]))[-1]
```

```
13
14  bicc
15  aicc
16
17  fitaic3 <- lm(log(cost) ~ ., data = nuclear[, c("cost", aicc)])
18  fitbic3 <- lm(log(cost) ~ ., data = nuclear[, c("cost", bicc)])
19
20  aic3 <- AIC(fitaic3)
21  bic3 <- BIC(fitbic3)
22
23  aic3
24  bic3
```

For backward selection we find the order cap →date →ct →t2 →pr →bw →t1.

```
# A tibble: 10 × 10
   date  t1    t2    cap   pr    ne    ct    bw
   <chr> <chr> <chr> <chr> <chr> <chr> <chr> <chr>
 1 " "   " "   " "   " "   " "   " "   " "   " "
 2 " "   " "   " "   "*"   " "   " "   " "   " "
 3 "*"   " "   " "   "*"   " "   " "   " "   " "
 4 "*"   " "   " "   "*"   " "   "*"   " "   " "
 5 "*"   " "   " "   "*"   " "   "*"   "*"   " "
 6 "*"   " "   " "   "*"   " "   "*"   "*"   " "
 7 "*"   " "   "*"   "*"   " "   "*"   "*"   " "
 8 "*"   " "   "*"   "*"   "*"   "*"   "*"   " "
 9 "*"   " "   "*"   "*"   "*"   "*"   "*"   "*"
10 "*"   "*"   "*"   "*"   "*"   "*"   "*"   "*"
#  2 more variables: cum.n <chr>, pt <chr>
```

Using $BIC$ we find the selected variables "date" "cap" "ne" "ct" "pt"  and for $AIC$ we find
"date" "t2" "cap" "pr" "ne" "ct" "cum.n" "pt" , and their respective $AIC$ and $BIC$:

```
> aic3
[1] -16.00549
> bic3
[1] -5.034235
```

## f)

```
1  # f)
2
3  faic <- aicb
4  fbic <- bicb
5  baic <- aicc
6  bbic <- bicc
```

```
7
8   cv <- function(vars, data, k=10){
9     n = nrow(data)
10    folds = sample(rep(1:k, length.out = n))
11    errors = numeric(k)
12    for (i in 1:k){
13      train = data[folds != i, ]
14      test = data[folds == i, ]
15      fit = lm(log(cost)~., data = train[, c("cost", vars)])
16      pred = predict(fit, newdata = test)
17      errors[i] = mean((log(test$cost) - pred)^2)
18    }
19    mean(errors)
20  }
21
22  cv_faic <- cv(faic, nuclear)
23  cv_fbic <- cv(fbic, nuclear)
24  cv_baic <- cv(baic, nuclear)
25  cv_bbic <- cv(bbic, nuclear)
26  cvs <- c(cv_faic, cv_fbic, cv_baic, cv_bbic)
27  names(cvs) <- c("fwd-aic", "fwd-bic", "bwd-aic", "bwd-bic")
28  bestcv <- names(cvs)[which.min(cvs)]
29  bestcv
30  min(cvs)
```

Using K-folds cross validation with $K = 10$ I find that the backward selection using $BIC$ gives the best model, and now prefer the model using date, cap, ne, ct, pt.

```
> bestcv <- names(cvs)[which.min(cvs)]
> bestcv
[1] "bwd-bic"
> min(cvs)
[1] 0.03565251
```

### g)

```
1   # g)
2
3   library(boot)
4
5   boot632 <- function(vars, data, B = 1000) {
6     set.seed(1814)
7     fit <- lm(log(cost) ~ ., data = data[, c("cost", vars)])
8     err_app <- mean(resid(fit)^2)
9     boot_fn <- function(data, indices) {
10      d <- data[indices, ]
11      fit_boot <- lm(log(cost) ~ ., data = d[, c("cost", vars)])
12      pred <- predict(fit_boot, newdata = data)
```

```
13        mean((log(data$cost) - pred)^2)
14    }
15    boot_res <- boot(data, boot_fn, R = B)
16    err_boot <- mean(boot_res$t)
17    0.368 * err_app + 0.632 * err_boot
18 }
19
20 boot_faic <- boot632(faic, nuclear)
21 boot_fbic <- boot632(fbic, nuclear)
22 boot_baic <- boot632(baic, nuclear)
23 boot_bbic <- boot632(bbic, nuclear)
24
25 boot_faic
26 boot_fbic
27 boot_baic
28 boot_bbic
```

Find the following $AIC$ and $BIC$, with the best being the backward model selected based on $AIC$.

```
> boot_faic
[1] 0.02699085
> boot_fbic
[1] 0.02738823
> boot_baic
[1] 0.02579897
> boot_bbic
[1] 0.02738823
```

## h)

```
1  # h)
2  library(glmnet)  # For ridge regression
3
4  x <- as.matrix(nuclear[, -1])  # All predictors, no cost
5  y <- log(nuclear$cost)          # Response
6
7  ridge_cv <- cv.glmnet(x, y, alpha = 0, nfolds = 10)
8  best_lambda <- ridge_cv$lambda.min
9
10 ridge_fit <- glmnet(x, y, alpha = 0, lambda = best_lambda)
11
12 cv <- function(vars, data, k = 10) {
13   set.seed(1814)
14   n <- nrow(data)
15   folds <- sample(rep(1:k, length.out = n))
16   errors <- numeric(k)
17   for (i in 1:k) {
18     train <- data[folds != i, ]
```

```
19        test <- data[folds == i, ]
20        fit <- lm(log(cost) ~ ., data = train[, c("cost", vars)])
21        pred <- predict(fit, newdata = test)
22        errors[i] <- mean((log(test$cost) - pred)^2)
23      }
24      mean(errors)
25    }
26
27    cv_ridge <- cv(all1, nuclear)  # all1 from d) is all predictors
28
29    best_lambda
30    cv_ridge
31
32    cv_faic
33    cv_fbic
34    cv_baic
35    cv_bbic
```

We can now assess how the ridge regression compares to the cross-validated subset models. The best lambda $\lambda$ we found was $\lambda \approx 0.0697$ which gave a cross-validated ridge-regression of $0.0513$, this is though higher than all of the subset models shown below, having values lower than $0.05$.

```
> best_lambda
[1] 0.06978489
> cv_ridge
[1] 0.05138956


> cv_faic
[1] 0.04475947
> cv_fbic
[1] 0.03703112
> cv_baic
[1] 0.04238047
> cv_bbic
[1] 0.03565251
```

## i)

```
1    # i)
2
3    library(glmnet)
4
5    x <- as.matrix(nuclear[, -1])
6    y <- log(nuclear$cost)
7
8    lasso_cv <- cv.glmnet(x, y, alpha = 1, nfolds = 10)
9    best_lambda_lasso <- lasso_cv$lambda.min
```

```
10  lasso_fit <- glmnet(x, y, alpha = 1, lambda = best_lambda_lasso)
11
12  # True Lasso CV error
13  cv_lasso <- lasso_cv$cvm[which(lasso_cv$lambda == best_lambda_lasso)
        ]
14
15  best_lambda_lasso
16  cv_lasso
17  cv_faic
18  cv_fbic
19  cv_baic
20  cv_bbic
21  cv_ridge
```

Here again we find a best lambda of $\lambda \approx 0.0245$, which gives a prediction error of $0.0387$, which in this case is better than that of ridge regression. This is also better than the subset models selected on $AIC$.

```
> best_lambda_lasso
[1] 0.02450281
> cv_lasso
[1] 0.03874789


> cv_faic
[1] 0.04837609
> cv_fbic
[1] 0.03693269
> cv_baic
[1] 0.04351076
> cv_bbic
[1] 0.03808641
> cv_ridge
[1] 0.05138956
```