# ISE789/OR791
# Homework 2

1. Let $y_i = \beta_0 + \beta_1(x_i - \bar{x}) + \epsilon_i$ for $i = 1, \cdots, n$. Using the formula $\hat{\boldsymbol{\beta}} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{y}$, show that

$$
\begin{aligned}
\hat{\beta}_0 &= \bar{y} \\
\hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.
\end{aligned}
$$

Now using the formula $var(\hat{\boldsymbol{\beta}}) = \sigma^2(\boldsymbol{X}'\boldsymbol{X})^{-1}$, show that

$$
\begin{aligned}
var(\hat{\beta}_0) &= \frac{\sigma^2}{n} \\
var(\hat{\beta}_1) &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.
\end{aligned}
$$

2. Prove the following: (i) $var(e_i) = \sigma^2 - var(\hat{y}_i)$. (ii) $var(e_i) = \sigma^2(1 - \boldsymbol{x}_i'(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{x}_i)$, where $\boldsymbol{x}_i$ is the $i$th row of the $\boldsymbol{X}$ matrix.

3. Suppose we have $n$ data from the model $y = \boldsymbol{x}'\boldsymbol{\beta} + \epsilon$, where the error satisfies the Gauss-Markov assumptions. Suppose, further, that we wish to predict the $(n+1)$st observation $y_{n+1}$ at $\boldsymbol{x}_{n+1}$. The predictor based on the least squares estimate of $\boldsymbol{\beta}$ is given by $\hat{y}_{n+1} = \boldsymbol{x}_{n+1}'\hat{\boldsymbol{\beta}}$.

   (a) Show that $E(\hat{y}_{n+1} - y_{n+1}) = 0$.

   (b) Suppose $\tilde{y}_{n+1} = \boldsymbol{a}'\boldsymbol{y}$ is another predictor of $y_{n+1}$ such that $E(\tilde{y}_{n+1} - y_{n+1}) = 0$. Show that $\boldsymbol{a}$ must satisfy $\boldsymbol{a}'\boldsymbol{X} = \boldsymbol{x}_{n+1}'$.

   (c) Show that $var(\hat{y}_{n+1}) \leq var(\tilde{y}_{n+1})$.

4. The dataset *teengamb* (available in the R library *faraway*) concerns a study of teenage gambling in Britain. Fit a regression model with the expenditure on gambling as the response and the sex (coded as male=0 and female=1), status, income, and verbal score as predictors. Present the output and answer the following questions.

1

(a) What percentage of variation in the response is explained by these predictors?

(b) Which observation has the largest (positive) residual? Give the case number.

(c) Compute the mean and median of the residuals.

(d) Compute the correlation of the residuals with the fitted values.

(e) Compute the correlation of the residuals with the income.

(f) For all other predictors held constant, what would be the difference in predicted expenditure on gambling for a male compared to a female.

(g) Which variables are statistically significant?

(h) Predict the amount that a male with average (given these data) status, income and verbal score would gamble along with an appropriate 95% CI. Repeat the prediction for a male with maximal values (for this data) of status, income, and verbal score. Which CI is wider and why is this expected?

(i) Fit a model with just income as a predictor and use an $F$-test to compare it to the full model.

Remarks:
1) Questions (f) and (h) will need to use the knowledge to be introduced in next class
2) Please find the dataset in the attachment. You may also find the data in R library *faraway*
3) Please choose the programming language you are comfortable with and submit your code along with your solution

1.

Given :

$$y_i = \beta_0 + \beta_1 (x_i - \bar{x}) + \varepsilon_i \quad \text{for } i = 1, \ldots, n \quad —①$$

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad —②$$

To prove :

$$\hat{\beta}_0 = \bar{y}.$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{n=1}^{i=n} (x_i - \bar{x})^2}$$

Solution :

we can write ① as . $(y_i - \bar{y}) = \hat{\beta}_i (x_i - \bar{x})$.

$$y = \beta \qquad x$$

$$\hat{\beta}_1 = \sum_{i=1}^{n} \left[ ((x_i - \bar{x})^T (x_i - \bar{x}))^{-1} (x_i - \bar{x})^T (y_i - \bar{y}) \right] - \text{putting.}$$

$$x = x_i - \bar{x} \quad \text{and} \quad y = (y_i - \bar{y}) \text{ into } ②.$$

$$\hat{\beta}_1 = \sum_{i=1}^{M} \left[ ((x_i - \bar{x})^2)^{-1} (x_i - \bar{x})(y_i - \bar{y}) \right]$$

$$\boxed{\hat{\beta}_1 = \frac{\sum_{n=1}^{i=n} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{n=1}^{i=n} (x_i - \bar{x})^2}}$$

Putting $y$ $\hat{\beta}_1$ value in eqn. ①.

$$\hat{\beta}_0 = y_i - \hat{\beta}_1 - \beta_1(x_i - \bar{x}) - \varepsilon_i$$

$$\hat{\beta}_0 = \sum_{i=1}^{n}\left[(y_i - \bar{y}) - \hat{\beta}_1(x_i - \bar{x})\right]$$

$I$

For $\hat{\beta}_0$ $\quad X = \begin{bmatrix} | \\ | \\ | \\ \vdots \\ n \end{bmatrix}_{n \times 1}$ $\quad$ as intercept at $x = 1$

So, $\hat{\beta}_0 = \sum_{i=1}^{n}\left(\begin{bmatrix} | \\ n \end{bmatrix}^T_{1 \times 1} \begin{bmatrix} | \end{bmatrix}\right)^{-1}\left(\begin{bmatrix} | \end{bmatrix}^T y\right)$

$$= \sum_{i=1}^{n}(n)^{-1}(y)$$

$$= \frac{\sum_{i=1}^{m} y}{n} = \bar{y}.$$

$$\boxed{\hat{\beta}_0 = \bar{y}}$$

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

$$\text{Var}(\hat{\beta}_0) = \sum_{i=1}^{n} \sigma^2 \left( \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}^T \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \right)^{-1} = \sigma^2 (n)^{-1} = \boxed{\dfrac{\sigma^2}{n}}$$

$$\text{Var}(\hat{\beta}_1) = \sigma^2 \sum_{i=} \left( (x_i - \bar{x})^T (x_i - \bar{x}) \right)^{-1}$$

$$\boxed{\text{Var}(\hat{\beta}_1) = \dfrac{\sigma^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}}$$

2) To prove:
$$\text{Var}(e_i) = \sigma^2 - \text{Var}(\hat{y}_i)$$

$$\text{Var}(e_i) = \text{Var}(y_i - \hat{y}_i) = \text{Var}(x_i \beta + \varepsilon_i - x_i \hat{\beta})$$

$$= \underbrace{\text{Var}(x_i \beta)}_{\text{constant}} + \underbrace{\text{Var}(\varepsilon_i)}_{\sigma^2} - \text{Var}(x_i \hat{\beta})$$

$$= 0 + \sigma^2 - \text{Var}(x_i \hat{\beta}) = \sigma^2 - \text{Var}(\hat{y}_i)$$

$$\boxed{\text{Var } e_i = \sigma^2 - \text{Var}(x_i \hat{\beta})} \qquad -①.$$

To prove: $\text{Var}(e_i) = \sigma^2 (1 - x_i^T (X^T X)^{-1} x_i)$

From ① $\quad \text{Var } e_i = \sigma^2 - \text{Var}(x_i \hat{\beta}) = \sigma^2 - x_i^T \text{Var}(\hat{\beta}_i)$

$$= \sigma^2 - x_i^T \sigma^2 (X^T X)^{-1} x_i$$

$$= \sigma^2 - \sigma^2 x_i^T (X^T X)^{-1} x_i$$

$$\boxed{\text{Var } e_i = \sigma^2 (1 - x_i^T (X^T X)^{-1} x_i)}$$

we have n data from model $y = x^T \beta + e$.

error is in Gauss - markov assumptions.

we wish to predict $(n+1)$st obs

$y_{n+1}$ at $x_{n+1}$.

and $\hat{y}_{n+1} = x^T_{n+1} \hat{\beta}$.

(a). $E(\hat{y}_{n+1} - y_{n+1}) = 0$.

$E(\hat{y}_{n+1} - y_{n+1}) = E(x^T_{n+1} \hat{\beta} - x^T_{n+1} \beta - \varepsilon)$

$= E(x^T_{n+1} \hat{\beta}) - x^T_{n+1} \beta - E(\varepsilon)$.

$= x^T_{n+1} E(\hat{\beta}) - x^T_{n+1} \beta - 0$

$= x^T_{n+1} \beta - x^T_{n+1} \beta$ (as $\hat{\beta}$ is unbiased estimator)

$= 0$.

(b). Suppose $\tilde{y}_{n+1} = a^T y$ another predictor of $y_{n+1}$

such that $E(\tilde{y}_{n+1} - y_{n+1}) = 0$.

$E(\tilde{y}_{n+1} - y_{n+1}) = E(a'y - x^T_{n+1} \beta - \varepsilon) = 0$

$\Rightarrow E(a'(X\beta + \varepsilon) - x^T_{n+1} \beta - \varepsilon) = 0$

$\Rightarrow E(a'X\beta - x^T_{n+1} \beta + a'\varepsilon - \varepsilon) = 0$

$\Rightarrow a'(($

$$\Rightarrow. \quad a'x\beta - x^T_{n+1}\beta + E\left[(a'-1)\varepsilon\right] = 0$$

$$\Rightarrow \quad a'x\beta - x^T_{n+1}\beta + a'-1 \, E[\varepsilon] \underbrace{= 0}_{\rightarrow 0},$$

$$\Rightarrow \quad \boxed{a'x\beta = x^T_{n+1}}$$

(a) $\mathcal{Y}_{ar}$ $Var(\hat{y}_{n+1}) \leq Var(\tilde{y}_{n+1})$

$$Var(\hat{y}_{n+1}) = Var(x'_{n+1}\beta)$$

$$= \sigma^2 \left((x'_{n+1})^T (x^Tx)^{-1}(x'_{n+1})\right)$$

$$= \sigma^2 \left(x_{n+1}(x^Tx)^{-1}(x'_{n+1})\right)$$

$$Var(\bar{y}_{n+1}) = Var(a'y) = Var(a'y).$$

$$= (a')' \, Var(x\beta)a'$$

$$= a x' Var(\beta)x a'$$

$$= a \, x' \sigma^2 (x^Tx)^{-1} x \, a'$$

$$= \sigma^2 \, a \, x'(x'x)^{-1} x a'$$

Since $a'x = x'_{n+1}$ then

$$Var(\hat{y}_{n+1}) = \sigma^2 (x_{n+1})(x^Tx)^{-1}(x'_{n+1})$$

$$\therefore \quad \text{var}(\breve{y}_{n+1}) = \text{var}(\bar{y}_{n+1})$$

otherwise $\quad \text{var}(\hat{y}_{n+1}) < \text{var}(\bar{y}_{n+1})$

$$\boxed{\therefore \quad \text{var}(\hat{y}_{n+1}) \leq \text{var}(\bar{y}_{n+1})}$$

# hw_2.R

4. The dataset teengamb (available in the R library faraway) concerns a study of teenage gambling in Britain. Fit a regression model with the expenditure on gambling as the response and the sex (coded as male=0 and female=1), status, income, and verbal score as predictors. Present the output and answer the following questions.

(a) What percentage of variation in the response is explained by these predictors?

(b) Which observation has the largest (positive) residual? Give the case number.

(c) Compute the mean and median of the residuals.

(d) Compute the correlation of the residuals with the fitted values.

(e) Compute the correlation of the residuals with the income.

(f) For all other predictors held constant, what would be the difference in predicted expenditure on gambling for a male compared to a female.

(g) Which variables are statistically significant?

(h) Predict the amount that a male with average (given these data) status, income and verbal score would gamble along with an appropriate 95% CI. Repeat the prediction for a male with maximal values (for this data) of status, income, and verbal score. Which CI is wider and why is this expected?

(i) Fit a model with just income as a predictor and use an F-test to compare it to the full model.

**Solution:**

<div align="center">

deepa

Tue Sep 24 20:50:26 2019

</div>

```
## Accessing the library Faraway to access the data

library(faraway)

## Warning: package 'faraway' was built under R version 3.5.3

##Showing the data and doing categorical distribution of sex variable

data(teengamb)
teengamb$sex <- factor(teengamb$sex)
attach(teengamb)
teengamb[1:3,]

##   sex status income verbal gamble
## 1   1    51    2.0      8      0
## 2   1    28    2.5      8      0
## 3   1    37    2.0      6      0

## Fitting linear model on gamble data on sex, status, income, verbal variabl
es

gamb.lm <- lm(gamble ~ sex+status+income+verbal)

## Showing statistics of the variable
summary(gamb.lm)

##
## Call:
## lm(formula = gamble ~ sex + status + income + verbal)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -51.082 -11.320  -1.451   9.452  94.252
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  22.55565   17.19680   1.312   0.1968
## sex1        -22.11833    8.21111  -2.694   0.0101 *
## status        0.05223    0.28111   0.186   0.8535
## income        4.96198    1.02539   4.839 1.79e-05 ***
## verbal       -2.95949    2.17215  -1.362   0.1803
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 22.69 on 42 degrees of freedom
```

```
## Multiple R-squared:  0.5267, Adjusted R-squared:  0.4816
## F-statistic: 11.69 on 4 and 42 DF,  p-value: 1.815e-06
```

**## A. Since, R squared is 0.5267 so, 52.67% of the response is explained by these predictors**

gamb.lm$residuals

```
##           1           2           3           4           5           6
##  10.6507430   9.3711318   5.4630298 -17.4957487  29.5194692  -2.9846919
##           7           8           9          10          11          12
##  -7.0242994 -12.3060734   6.8496267 -10.3329505   1.5934936  -3.0958161
##          13          14          15          16          17          18
##   0.1172839   9.5331344   2.8488167  17.2107726 -25.2627227 -27.7998544
##          19          20          21          22          23          24
##  13.1446553 -15.9510624 -16.0041386  -9.5801478 -27.2711657  94.2522174
##          25          26          27          28          29          30
##   0.6993361  -9.1670510 -25.8747696  -8.7455549  -6.8803097 -19.8090866
##          31          32          33          34          35          36
##  10.8793766  15.0599340  11.7462296  -3.5932770 -14.4016736  45.6051264
##          37          38          39          40          41          42
##  20.5472529  11.2429290 -51.0824078   8.8669438  -1.4513921  -3.8361619
##          43          44          45          46          47
##  -4.3831786 -14.8940753   5.4506347   1.4092321   7.1662399
```

max(gamb.lm$residuals)

```
## [1] 94.25222
```

**## B. 24th observation has maximum residual of 94.2522174**

mean(gamb.lm$residuals)

```
## [1] -3.065293e-17
```

median(gamb.lm$residuals)

```
## [1] -1.451392
```

**## C. Mean of the residual is approximately 0 while median is -1.451392.**

fitted_value <- gamble - gamb.lm$residuals
cor(gamb.lm$residuals, fitted_value)

```
## [1] -1.070659e-16
```

**## D. Correlation between residuals and fitted value is approximately zero.**

```r
cor(gamb.lm$residuals, income)
```

```
## [1] -7.242382e-17
```

## E. Correlation of residuals and income is almost 0

## F. Keeping everything constant, since coefficient is -22.11833 thus average female teen spend $ 22.1183 less than male teen.

## G. Since P-value of only income and sex is less than 5% so, only these are significantly important.

## H.

```r
male=data.frame(sex=0, status=mean(teengamb$status), income=mean(teengamb$income), verbal=mean(teengamb$verbal))
```

```r
predict(gamb.lm,male, se.fit=FALSE, interval='confidence')
        fit      lwr      upr
1 28.24252 18.78277 37.70227
```

```r
male_max=data.frame(sex=0, status=max(teengamb$status), income=max(teengamb$income), verbal=max(teengamb$verbal))
```

```r
predict(gamb.lm,male_max, se.fit=FALSE, interval='confidence')
        fit      lwr      upr
1 71.30794 42.23237 100.3835
```

Confidence interval of maximum values are wide because values of predictor is far away from regression line.

## I

```r
gamb_income.lm <- lm(gamble ~ income)
summary(gamb_income.lm)
```

```
Call:
lm(formula = gamble ~ income)

Residuals:
    Min      1Q  Median      3Q     Max
-46.020 -11.874  -3.757  11.934 107.120

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   -6.325      6.030  -1.049      0.3
income         5.520      1.036   5.330 3.05e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 24.95 on 45 degrees of freedom
Multiple R-squared:  0.387,    Adjusted R-squared:  0.3734
F-statistic: 28.41 on 1 and 45 DF,  p-value: 3.045e-06
```

F-Test for full model is 11.69 on 4 variables and F-Test for income is 28.41 on 1 variable so, income model is working better than the full model.