

STAT 3690 Lecture Note

Part VI: Linear model

Zhiyang Zhou (zhiyang.zhou@umanitoba.ca, zhiyanggeezhou.github.io)

2023/Mar/15 11:22:39

Multivariate linear regression

What is a linear model?

- Responses are linear functions with respect to unknown parameters.

Univariate/multiple linear regression (J&W Sec. 7.2–7.5)

- Model (population version):

$$Y \mid X_1, \dots, X_q \sim \left(\sum_{j=1}^q X_j \beta_j, \sigma^2 \right)$$

- Equiv. $Y = \sum_{j=1}^q X_j \beta_j + \varepsilon$ with $\varepsilon \perp\!\!\!\perp [X_1, \dots, X_q]^\top$ and $\varepsilon \sim (0, \sigma^2)$
- Univariate linear regression: $q = 2$ with $X_1 = 1$
- Multiple linear regression: $q > 2$ with $X_1 = 1$

- Model (sample version):

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

- $\mathbf{Y} = [Y_1, \dots, Y_n]^\top$
- Design matrix

$$\mathbf{X} = \begin{bmatrix} X_{11} & \cdots & X_{1q} \\ \vdots & \ddots & \vdots \\ X_{n1} & \cdots & X_{nq} \end{bmatrix}_{n \times q}$$

- * $\text{rk}(\mathbf{X}) = q$
 - $\boldsymbol{\beta} = [\beta_1, \dots, \beta_q]^\top$
 - $\boldsymbol{\varepsilon} = [\varepsilon_1, \dots, \varepsilon_n]^\top \sim (\mathbf{0}_n, \sigma^2 \mathbf{I}_n)$, independent of \mathbf{X}
-

- Least squares (LS) estimation (no need of normality)

- $\hat{\boldsymbol{\beta}}_{\text{LS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$
 - * $E(\hat{\boldsymbol{\beta}}_{\text{LS}} \mid \mathbf{X}) = \boldsymbol{\beta}$
 - $\hat{\sigma}_{\text{LS}}^2 = (n - q)^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{LS}})^\top (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{LS}}) = (n - q)^{-1} \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}$
 - * $n \times n$ hat matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$
 - * $E(\hat{\sigma}_{\text{LS}}^2 \mid \mathbf{X}) = \sigma^2$
-

- ML estimation (under normality)

- $\hat{\beta}_{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \hat{\beta}_{\text{LS}}$
 - * $\hat{\beta}_{\text{ML}} \mid \mathbf{X} \sim \text{MVN}_q(\beta, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$
 - $\hat{\sigma}_{\text{ML}}^2 = n^{-1} \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y} = n^{-1} (n - q) \hat{\sigma}_{\text{LS}}^2$
 - * Given \mathbf{X} , $n \hat{\sigma}_{\text{ML}}^2 / \sigma^2 = (n - q) \hat{\sigma}_{\text{LS}}^2 / \sigma^2 \sim \chi^2(n - q)$
-

- Inference (under normality)
 - To infer $\mathbf{a}^\top \beta$, given $\mathbf{a} \in \mathbb{R}^q$ (e.g., to compare β_1 and β_2 by checking $\mathbf{a}^\top \beta = \beta_1 - \beta_2$ with $\mathbf{a} = [1, -1, 0, \dots, 0]^\top$)
 - * Estimator: $\mathbf{a}^\top \hat{\beta}_{\text{ML}}$
 - * $100 \times (1 - \alpha)\%$ confidence interval for $\mathbf{a}^\top \beta$:

$$\mathbf{a}^\top \hat{\beta}_{\text{ML}} \pm \hat{\sigma}_{\text{LS}} \cdot t_{1-\alpha/2, n-q} \sqrt{\mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}$$

- To predict $Y_0 = \mathbf{X}_0^\top \beta + \varepsilon_0$ with \mathbf{X}_0 different from each row of \mathbf{X}
 - * Prediction: $\hat{Y}_0 = \mathbf{X}_0^\top \hat{\beta}_{\text{ML}}$
 - * $100 \times (1 - \alpha)\%$ prediction interval for Y_0

$$\mathbf{X}_0^\top \hat{\beta}_{\text{ML}} \pm \hat{\sigma}_{\text{LS}} \cdot t_{1-\alpha/2, n-q} \sqrt{1 + \mathbf{X}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_0}$$

Multivariate linear regression

- Model (population version):

$$Y_1, \dots, Y_p \mid X_1, \dots, X_q \sim ([X_1, \dots, X_q] \mathbf{B}, \Sigma)$$

- Equiv. $[Y_1, \dots, Y_p] = [X_1, \dots, X_q] \mathbf{B} + \varepsilon^\top$ with p -vector $\varepsilon \perp [X_1, \dots, X_q]$ and $\varepsilon \sim (\mathbf{0}_p, \Sigma)$
 - * Unknown coefficients

$$\mathbf{B} = \begin{bmatrix} b_{11} & \cdots & b_{1p} \\ \vdots & \ddots & \vdots \\ b_{q1} & \cdots & b_{qp} \end{bmatrix}_{q \times p} = \begin{bmatrix} \mathbf{b}_{1\cdot}^\top \\ \vdots \\ \mathbf{b}_{q\cdot}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{b}_{\cdot 1} & \cdots & \mathbf{b}_{\cdot p} \end{bmatrix}$$

- $\mathbf{b}_{i\cdot}^\top$: the i th row of \mathbf{B}
- $\mathbf{b}_{\cdot j}$: the j th column of \mathbf{B}

- Model (sample version):

$$\begin{matrix} \mathbf{Y} & = & \mathbf{X} & \mathbf{B} & + & \mathbf{E} \\ n \times p & = & n \times q & q \times p & + & n \times p \end{matrix}$$

- Response

$$\mathbf{Y} = \begin{bmatrix} Y_{11} & \cdots & Y_{1p} \\ \vdots & \ddots & \vdots \\ Y_{n1} & \cdots & Y_{np} \end{bmatrix}_{n \times p}$$

- Design matrix

$$\mathbf{X} = \begin{bmatrix} X_{11} & \cdots & X_{1q} \\ \vdots & \ddots & \vdots \\ X_{n1} & \cdots & X_{nq} \end{bmatrix}_{n \times q}$$

- * $\text{rk}(\mathbf{X}) = q < p + q \leq n$

- Error

$$\mathbf{E} = \begin{bmatrix} e_{11} & \cdots & e_{1q} \\ \vdots & \ddots & \vdots \\ e_{n1} & \cdots & e_{nq} \end{bmatrix}_{n \times q} = \begin{bmatrix} \mathbf{e}_{1\cdot}^\top \\ \vdots \\ \mathbf{e}_{n\cdot}^\top \end{bmatrix}$$

- * $\mathbf{e}_{i\cdot} \perp\!\!\!\perp [X_{i1}, \dots, X_{iq}]$
- * $\mathbf{e}_{i\cdot} \stackrel{\text{iid}}{\sim} (\mathbf{0}_p, \mathbf{\Sigma})$

- Relationship with MANOVA
 - MANOVA models can be expressed as multivariate linear regression with a carefully selected \mathbf{X} .
- Exercise 6.1: rephrase the following one-way MANOVA model

$$\mathbf{Y}_{ij} = \boldsymbol{\mu} + \boldsymbol{\tau}_i + \mathbf{E}_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, m$$

into a multivariate linear regression model, where $\mathbf{E}_{ij} \stackrel{\text{iid}}{\sim} \text{MVN}_p(\mathbf{0}, \mathbf{\Sigma})$ and $\sum_i \boldsymbol{\tau}_i = \mathbf{0}$.

- LS estimation (no need of normality)
 - $\hat{\mathbf{B}}_{\text{LS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$
 - * $E(\hat{\mathbf{B}}_{\text{LS}} | \mathbf{X}) = \mathbf{B}$
 - $\hat{\mathbf{\Sigma}}_{\text{LS}} = (n - q)^{-1} (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}}_{\text{LS}})^\top (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}}_{\text{LS}}) = (n - q)^{-1} \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}$
 - * $E(\hat{\mathbf{\Sigma}}_{\text{LS}} | \mathbf{X}) = \mathbf{\Sigma}$
- ML estimation (under normality)
 - $\hat{\mathbf{B}}_{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \hat{\mathbf{B}}_{\text{LS}}$
 - $\hat{\mathbf{\Sigma}}_{\text{ML}} = n^{-1} \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y} = n^{-1} (n - q) \hat{\mathbf{\Sigma}}_{\text{LS}}$
 - * Given \mathbf{X} , $n \hat{\mathbf{\Sigma}}_{\text{ML}} \sim W_p(\mathbf{\Sigma}, n - q)$
- Inference (under normality)
 - To infer $\mathbf{B}^\top \mathbf{a}$, given $\mathbf{a} \in \mathbb{R}^q$ (e.g., to compare the 1st and 2nd rows of \mathbf{B} , i.e., \mathbf{b}_1 and \mathbf{b}_2 ., by checking $\mathbf{B}^\top \mathbf{a} = \mathbf{b}_1 - \mathbf{b}_2$ with $\mathbf{a} = [1, -1, 0, \dots, 0]^\top$)
 - * Estimator: $\hat{\mathbf{B}}_{\text{ML}}^\top \mathbf{a}$
 - * $100 \times (1 - \alpha)\%$ confidence region for $\mathbf{B}^\top \mathbf{a}$

$$\left\{ \mathbf{u} \in \mathbb{R}^p : (\mathbf{u} - \hat{\mathbf{B}}_{\text{ML}}^\top \mathbf{a})^\top \hat{\mathbf{\Sigma}}_{\text{LS}}^{-1} (\mathbf{u} - \hat{\mathbf{B}}_{\text{ML}}^\top \mathbf{a}) \leq \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} \cdot \frac{p(n - q)}{n - p - q + 1} F_{1-\alpha, p, n-p-q+1} \right\}$$
 - To predict $\mathbf{Y}_0 = \mathbf{B}^\top \mathbf{X}_0 + \mathbf{E}_0$ with newly observed $\mathbf{X}_0 \in \mathbb{R}^q$
 - * Prediction: $\hat{\mathbf{Y}}_0 = \mathbf{B}_{\text{ML}}^\top \mathbf{X}_0$
 - * $100 \times (1 - \alpha)\%$ prediction region for \mathbf{Y}_0

$$\left\{ \mathbf{u} \in \mathbb{R}^p : (\mathbf{u} - \hat{\mathbf{Y}}_0)^\top \hat{\mathbf{\Sigma}}_{\text{LS}}^{-1} (\mathbf{u} - \hat{\mathbf{Y}}_0) \leq \{1 + \mathbf{X}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_0\} \cdot \frac{p(n - q)}{n - p - q + 1} F_{1-\alpha, p, n-p-q+1} \right\}$$
 - To infer $\mathbf{a}^\top \mathbf{Y}_0 = \mathbf{a}^\top (\mathbf{B}^\top \mathbf{X}_0 + \mathbf{E}_0)$, given $\mathbf{a} \in \mathbb{R}^p$ and newly observed $\mathbf{X}_0 \in \mathbb{R}^q$
 - * Prediction: $\mathbf{a}^\top \hat{\mathbf{Y}}_0 = \mathbf{a}^\top \mathbf{B}_{\text{ML}}^\top \mathbf{X}_0$
 - * $100 \times (1 - \alpha)\%$ prediction interval for $\mathbf{a}^\top \mathbf{Y}_0$

$$\mathbf{a}^\top \hat{\mathbf{Y}}_0 \pm \sqrt{\mathbf{a}^\top \hat{\mathbf{\Sigma}}_{\text{LS}} \mathbf{a} \cdot \{1 + \mathbf{X}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_0\} \cdot t_{1-\alpha/2, n-q}}$$
 - $100 \times (1 - \alpha)\%$ simultaneous prediction intervals for $\mathbf{a}_k^\top \mathbf{Y}_0$, $k = 1, \dots, m$, given $\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^p$ and newly observed $\mathbf{X}_0 \in \mathbb{R}^q$
 - * (Bonferroni)

$$\mathbf{a}_k^\top \hat{\mathbf{Y}}_0 \pm \sqrt{\mathbf{a}_k^\top \hat{\mathbf{\Sigma}}_{\text{LS}} \mathbf{a}_k \cdot \{1 + \mathbf{X}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_0\} \cdot t_{1-\alpha/(2m), n-q}}$$
 - * (Scheffé's)

$$\mathbf{a}_k^\top \hat{\mathbf{Y}}_0 \pm \sqrt{\mathbf{a}_k^\top \hat{\mathbf{\Sigma}}_{\text{LS}} \mathbf{a}_k \cdot \{1 + \mathbf{X}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_0\} \cdot \frac{p(n - q)}{n - p - q + 1} F_{1-\alpha, p, n-p-q+1}}$$

```

options(digits = 4)
tear <- c(
  6.5, 6.2, 5.8, 6.5, 6.5, 6.9, 7.2, 6.9, 6.1, 6.3,
  6.7, 6.6, 7.2, 7.1, 6.8, 7.1, 7.0, 7.2, 7.5, 7.6
)
gloss <- c(
  9.5, 9.9, 9.6, 9.6, 9.2, 9.1, 10.0, 9.9, 9.5, 9.4,
  9.1, 9.3, 8.3, 8.4, 8.5, 9.2, 8.8, 9.7, 10.1, 9.2
)
rate <- factor(gl(2,10,length=length(tear)), labels=c("Low", "High"))

# Model fitting
fit <- lm(cbind(tear, gloss) ~ rate)
summary(fit)

# Prediction
(Obs_new <- data.frame(rate = factor(c("High"), levels = c("Low", "High"))))
(prediction <- t(predict(fit, newdata = Obs_new)))

# Prediction region
n = nrow(model.matrix(fit))
p = ncol(coef(fit))
q = ncol(model.matrix(fit))-1
(X <- model.matrix(fit))
(X0 <- t(model.matrix(~rate, Obs_new)))
(SigmaHatLS <- crossprod(resid(fit))/(n-q))
quad_form <- drop(t(X0) %*% solve(crossprod(X)) %*% X0)
fvalue = p*(n-q)/(n-p-q+1)*qf(0.95, p, n-p-q+1)

# 95% prediction region for Y0
c1 = sqrt((1 + quad_form)*fvalue)
car::ellipse(center = as.vector(prediction), shape = SigmaHatLS, radius = c1, add = F,
  xlab = "tear", ylab = "gloss", col = 'blue')

# 95% confidence region for t(B)X0
c2 = sqrt(quad_form*fvalue)
car::ellipse(center = as.vector(prediction), shape = SigmaHatLS, radius = c2, add = T,
  xlab = "tear", ylab = "gloss", col = 'red')

# 95% Scheffé's simultaneous prediction intervals for entries of Y0
a1 = c(1,0)
c(
  t(a1) %*% prediction - (t(a1) %*% SigmaHatLS %*% a1)^.5 * c1,
  t(a1) %*% prediction + (t(a1) %*% SigmaHatLS %*% a1)^.5 * c1
) # for tear
a2 = c(0,1)
c(
  t(a2) %*% prediction - (t(a2) %*% SigmaHatLS %*% a2)^.5 * c1,
  t(a2) %*% prediction + (t(a2) %*% SigmaHatLS %*% a2)^.5 * c1
) # for gloss

```

Testing for nested models

- $H_0 : E(\mathbf{Y} \mid \mathbf{X}) = \mathbf{X}_{(0)}\mathbf{B}_{(0)}$ (nested model) vs. $H_1 : E(\mathbf{Y} \mid \mathbf{X}) = \mathbf{X}_{(0)}\mathbf{B}_{(0)} + \mathbf{X}_{(1)}\mathbf{B}_{(1)}$ (larger model)
 - When $\mathbf{X}_{(0)}$ has only the column of ones, the model under H_0 is the empty/null model (i.e., only the intercept).
 - When $\mathbf{X}_{(1)}$ only involves one explanatory variable (i.e., is of $n \times 1$), we are testing the significance of that variable.
- Likelihood ratio

$$\lambda = \left(\frac{\det \hat{\Sigma}_{\text{ML}, H_0}}{\det \hat{\Sigma}_{\text{ML}}} \right)^{-n/2} = \left[\det \left\{ (\hat{\Sigma}_{\text{ML}, H_0} - \hat{\Sigma}_{\text{ML}}) \hat{\Sigma}_{\text{ML}}^{-1} + \mathbf{I} \right\} \right]^{-n/2}$$

- Test statistics alternative to the likelihood ratio
 - Wilks' lambda: $\prod_i (1 + \eta_i)^{-1}$
 - Pillai's trace: $\sum_i \{\eta_i (1 + \eta_i)^{-1}\}$
 - Hotelling-Lawley trace: $\sum_i \eta_i$
 - Roy's largest root: $\eta_1 (1 + \eta_1)^{-1}$
 - * Suppose $\eta_1 \geq \dots \geq \eta_p$ are eigenvalues of $(\hat{\Sigma}_{\text{ML}, H_0} - \hat{\Sigma}_{\text{ML}}) \hat{\Sigma}_{\text{ML}}^{-1}$
 - * When $\mathbf{X}_{(1)}$ has only one column (i.e., is of $n \times 1$), all the four tests are equivalent;
 - * As n increases, all these tests give similar results.

```
options(digits = 4)
tear <- c(
  6.5, 6.2, 5.8, 6.5, 6.5, 6.9, 7.2, 6.9, 6.1, 6.3,
  6.7, 6.6, 7.2, 7.1, 6.8, 7.1, 7.0, 7.2, 7.5, 7.6
)
gloss <- c(
  9.5, 9.9, 9.6, 9.6, 9.2, 9.1, 10.0, 9.9, 9.5, 9.4,
  9.1, 9.3, 8.3, 8.4, 8.5, 9.2, 8.8, 9.7, 10.1, 9.2
)
opacity <- c(
  4.4, 6.4, 3.0, 4.1, 0.8, 5.7, 2.0, 3.9, 1.9, 5.7,
  2.8, 4.1, 3.8, 1.6, 3.4, 8.4, 5.2, 6.9, 2.7, 1.9
)
rate <- factor(gl(2,10,length=length(tear)), labels=c("Low", "High"))
additive <- factor(gl(2,5,length=length(tear)), labels=c("Low", "High"))

# Testing the necessity of interaction
fit0 <- lm(cbind(tear, gloss, opacity) ~ -1)
fit1 = lm(cbind(tear, gloss, opacity) ~ -1+rate+additive+rate:additive)
anova(fit1, fit0, test='Wilks')
anova(fit1, fit0, test='Pillai')
anova(fit1, fit0, test='Hotelling')
anova(fit1, fit0, test='Roy')
```

Information criteria

- Akaike's information criterion (AIC)

$$- \ln \text{Likelihood} + 2 \times \text{number of parameters to estimate}$$

$$- \text{Number of parameters to estimate in } \mathbf{B} \text{ and } \mathbf{\Sigma}: pq + p(p+1)/2$$

- The smaller, the better.
- Bayesian information criterion (BIC)
 - $-\ln \text{Likelihood} + \ln n \times \text{number of parameters to estimate}$
- Model selection using information criteria proceeds as follows
 - Select models of interest, say M_1, \dots, M_K , which do NOT need to be nested.
 - * Candidate models should be selected using domain-specific expertise, if possible. Or, you can go through all possible models.
 - Compute the specific information criterion for each model.
 - Select the model with the smallest value of the information criterion.

```
options(digits = 4)
tear <- c(
  6.5, 6.2, 5.8, 6.5, 6.5, 6.9, 7.2, 6.9, 6.1, 6.3,
  6.7, 6.6, 7.2, 7.1, 6.8, 7.1, 7.0, 7.2, 7.5, 7.6
)
gloss <- c(
  9.5, 9.9, 9.6, 9.6, 9.2, 9.1, 10.0, 9.9, 9.5, 9.4,
  9.1, 9.3, 8.3, 8.4, 8.5, 9.2, 8.8, 9.7, 10.1, 9.2
)
opacity <- c(
  4.4, 6.4, 3.0, 4.1, 0.8, 5.7, 2.0, 3.9, 1.9, 5.7,
  2.8, 4.1, 3.8, 1.6, 3.4, 8.4, 5.2, 6.9, 2.7, 1.9
)
rate <- factor(gl(2,10,length=length(opacity)), labels=c("Low", "High"))
additive <- factor(gl(2,5,length=length(opacity)), labels=c("Low", "High"))

fit0 <- lm(cbind(tear, gloss, opacity) ~ rate)
logLik(fit0)
AIC(fit0)
BIC(fit0)

logLik.mlm <- function(object, ...) {
  resids <- residuals(object)
  Sigma_ML <- crossprod(resids)/nrow(resids)
  ans <- sum(mvtnorm::dmvnorm(resids, sigma = Sigma_ML, log = TRUE))
  df <- prod(dim(coef(object))) + choose(ncol(Sigma_ML) + 1, 2)
  attr(ans, "df") <- df
  class(ans) <- "logLik"
  return(ans)
}
logLik(fit0)
AIC(fit0)
BIC(fit0)

fit1 <- lm(cbind(tear, gloss, opacity) ~ additive)
AIC(fit1)
BIC(fit1)
```

Multivariate influence measures

- Hat matrix $\mathbf{H} = [h_{ij}]_{n \times n} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

- Leverage: the influence of the i th observation (i.e., the i th row of \mathbf{Y}), say $\mathbf{Y}_{i\cdot}^\top$, on $\hat{\mathbf{Y}}_{i\cdot}$ ($= h_{ii}\mathbf{Y}_{i\cdot} + \sum_{j \neq i} h_{ij}\mathbf{Y}_{j\cdot}$); specifically, $\mathbf{Y}_{i\cdot}$ is said to have a high leverage if h_{ii} is large compared to the other diagonal entries of hat matrix \mathbf{H}
- (Externally) Studentized residuals

$$T_i^2 = \frac{\hat{\mathbf{e}}_{i\cdot}^\top \hat{\boldsymbol{\Sigma}}_{\text{LS},(-i)}^{-1} \hat{\mathbf{e}}_{i\cdot}}{1 - h_{ii}}$$

- $\hat{\mathbf{e}}_{i\cdot}^\top$: the i th row of residual matrix $\hat{\mathbf{E}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$
- $\hat{\mathbf{E}}_{(-i)}^\top$: the remaining part of $\hat{\mathbf{E}}$ with Row i removed
- $\hat{\boldsymbol{\Sigma}}_{\text{LS},(-i)} = (n - q - 1)^{-1} \hat{\mathbf{E}}_{(-i)}^\top \hat{\mathbf{E}}_{(-i)}$: LS estimator of $\boldsymbol{\Sigma}$ after removing Row i from the residual matrix
- The i th observation may be considered as a potential outlier if

$$T_i^2 > \frac{p(n - q - 1)}{n - p - q} F_{1-\alpha, p, n-q-1}$$

* $F_{1-\alpha, p, n-q-1}$: the $1 - \alpha$ quantile of $F(p, n - q - 1)$

- (Multivariate) Cook's distance

$$D_i = \frac{h_{ii}}{q(1 - h_{ii})^2} \hat{\mathbf{e}}_{i\cdot}^\top \hat{\boldsymbol{\Sigma}}_{\text{LS}}^{-1} \hat{\mathbf{e}}_{i\cdot}$$

- The Cut-off is far from unique even for multiple linear regression (i.e., the case with $p = 1$)
- Pay attention to a small set of observations that have substantially higher values than the remaining observations

```
install.packages(c("car"))
options(digits = 4)
tear <- c(
  6.5, 6.2, 5.8, 6.5, 6.5, 6.9, 7.2, 6.9, 6.1, 6.3,
  6.7, 6.6, 7.2, 7.1, 6.8, 7.1, 7.0, 7.2, 7.5, 7.6
)
gloss <- c(
  9.5, 9.9, 9.6, 9.6, 9.2, 9.1, 10.0, 9.9, 9.5, 9.4,
  9.1, 9.3, 8.3, 8.4, 8.5, 9.2, 8.8, 9.7, 10.1, 9.2
)
opacity <- c(
  4.4, 6.4, 3.0, 4.1, 0.8, 5.7, 2.0, 3.9, 1.9, 5.7,
  2.8, 4.1, 3.8, 1.6, 3.4, 8.4, 5.2, 6.9, 2.7, 1.9
)
rate <- factor(gl(2,10,length=length(opacity)), labels=c("Low", "High"))
additive <- factor(gl(2,5,length=length(opacity)), labels=c("Low", "High"))

fit0 <- lm(cbind(tear, gloss, opacity) ~ rate*additive)
resids <- residuals(fit0)

# Leverage
X <- model.matrix(fit0)
H <- X %*% solve(crossprod(X)) %*% t(X)
(Hii = diag(H))
hist(Hii, 50)

# Externally Studentized residuals
n <- nrow(X)
```

```

p = ncol(resids)
T_square = numeric(n)
for (i in 1:n){
  SigmaHatLS_i <- crossprod(resids[-i,])/(n-1-ncol(X))
  T_square[i] = t(resids[i,]) %*% solve(SigmaHatLS_i) %*% resids[i,]
}
hist(T_square, 50)
which(T_square > p*(n-1-ncol(X))/(n-p-ncol(X))*qf(.95, p, n-1-ncol(X)))

# Cook's distance
SigmaHatLS <- crossprod(resids)/(n - ncol(X))
cookD <- Hii/((1 - Hii)^2*ncol(X)) * diag(resids %*% solve(SigmaHatLS) %*% t(resids))
hist(cookD, 50)
which(cookD>0.4)

```

Normality of residuals

- Check the normality of residuals following Lecture Note Part 3
 - Apply Box-Cox transformation to columns of \mathbf{Y}
-