

1. In research, it was investigated how the tensile strength of a paper depends on the percentage of the hardwood portion in raw material mixture  $X_1$ =Hardwood and (mechanical) scrubbing pressure  $X_2$ =pressure during the manufacturing process of paper. Below is a part of the material available in the study. The entire material can be found in dataset [paper.txt](#).

	strength	hardwood	pressure
1	196.6	2	400
2	197.7	2	500
3	199.8	2	650
4	198.4	2	400
.			
35	197.8	8	500
36	199.8	8	650

Denote explanatory variables as  $X_1$ =hardwood and  $X_2$ =pressure. Consider modeling the response variable  $Y$ =strength by following two different models:

$$\mathcal{M}_1 : Y_i = \beta_0 + \beta_1 x_{i1} + \varepsilon_i,$$

$$\mathcal{M}_{1|2} : Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i,$$

where in each model the random error term  $\varepsilon_i$  is assumed to follow normal distribution  $\varepsilon_i \sim N(0, \sigma^2)$ .

- Under the model  $\mathcal{M}_{1|2}$ , calculate the maximum likelihood estimate for the parameter  $\beta_2$ .  
(1 point)
- Under the model  $\mathcal{M}_{1|2}$ , find the restricted maximum likelihood estimate, i.e., an unbiased estimate  $\tilde{\sigma}^2$  for the variance parameter  $\sigma^2$ . (1 point)
- Under the model  $\mathcal{M}_{1|2}$ , calculate the fitted value  $\hat{\mu}_1$  for the first observation  $i = 1$  in the data set.  
(1 point)
- Under the model  $\mathcal{M}_{1|2}$ , calculate maximum likelihood estimate for the expected value  $\mu_{i*}$ , when  $x_{i*1} = 7$  and  $x_{i*2} = 500$ .  
(1 point)
- Under the model  $\mathcal{M}_{1|2}$ , calculate the 80% **prediction interval** for the new observation  $Y_{i*}$ , when  $x_{i*1} = 7$  and  $x_{i*2} = 500$ . Particularly, what is your estimate for lower bound of the prediction interval?  
(1 point)
- Consider the following hypotheses**

$H_0$  : Model  $\mathcal{M}_1$  is the true model,

$H_1$  : Model  $\mathcal{M}_{1|2}$  is the true model.

**Select the appropriate test statistic to test the above hypotheses. Calculate the value of the test statistic.** (1 point)

2. To all karate enthusiasts, it would be nice to find such a punching board (makiwara board) that will withstand the blows but which would not be so rigid or hard that training would then harm hands. The makiwara board can be made in different kinds of wood. In study, it was examined how much a makiwara board bends (in millimeters) of the force of the strike in different tree species. The makiwara boards used in study were made in two different ways. Data set is given in file [makiwaraboard.txt](#).

	WoodType	BoardType	Deflection
1	1	1	144.3
2	1	1	125.9
3	1	1	263.2
4	1	1	114.6
5	1	1	242.5
6	1	1	141.9
.			
.			
335	4	2	73.3
336	4	2	44.9

Description: Results of experiments measuring deflection (mm) of makiwara boards of two types (stacked and tapered) and of four wood types (Cherry, Ash, Fir, and Oak).

Wood Type: 1=Cherry, 2=Ash, 3=Fir, 4=Oak

Board Type: 1=Stacked, 2=Tapered

Source: P.K. Smith, T. Niiler, and P.W. McCullough (2010). "Evaluating Makiwara Punching Board Performance," Journal of Asian Martial Arts, Vol 19, #2, pp. 34-45.

Denote explanatory variables as  $X_1$ =WoodType and  $X_2$ =BoardType. Consider modeling the response variable  $Y$ =Deflection by following two different models:

$$\begin{aligned}\mathcal{M}_{1|2}: \quad Y_i &\sim N(\mu_{jh}, \sigma^2), \\ \mu_{jh} &= \beta_0 + \beta_j + \alpha_h, \\ \mathcal{M}_{12}: \quad Y_i &\sim N(\mu_{jh}, \sigma^2) \\ \mu_{jh} &= \beta_0 + \beta_j + \alpha_h + \gamma_{jh},\end{aligned}$$

where index  $j$  is related to the categories of the variable  $X_1$ =WoodType and index  $h$  is related to the categories of the variable  $X_2$ =BoardType.

- (a) Under the model  $\mathcal{M}_{1|2}$ , calculate the maximum likelihood estimate for the expected value  $\mu_{jh}$ , when the explanatory variables  $X_1, X_2$  are set on values

$$\begin{aligned}X_1 &= \text{Oak} = 4, \\ X_2 &= \text{Tapered} = 2.\end{aligned}$$

That is, find the maximum likelihood estimate for the expected value  $\mu_{42}$ .

(2 points)

- (b) Let us assume that the model  $\mathcal{M}_{1|2}$  fits sufficiently enough to the given data set. In which tree species the estimate of the expected value  $\mu_{jh}$  is in highest level?
- Cherry,
  - Ash,
  - Fir,
  - Oak.
- (c) Under the model  $\mathcal{M}_{12}$ , calculate the maximum likelihood estimate for the parameter  $\gamma_{32}$ .  
(1 point)
- (d) Under the model  $\mathcal{M}_{12}$ , calculate the residual  $e_i$  for the last observation  $i = 336$  in the data set.  
(1 point)
- (e) Under the model  $\mathcal{M}_{12}$ . Consider the following hypotheses

$$H_0 : \gamma_{jh} = 0, \quad H_1 : \gamma_{jh} \neq 0.$$

Select the appropriate test statistic to test the above hypotheses. Calculate the value of the test statistic.

(1 point)

3. (a) Consider the following small data, where  $X_1$  is a numerical explanatory variable and  $X_2$  is categorical explanatory variable having class values  $\{a, b, c\}$ .

	X1	X2	Y
1	3	a	46.0
2	3	b	55.4
3	3	c	57.9
4	6	a	55.5
5	6	b	66.7
6	6	c	68.6
7	9	a	65.3
8	9	b	76.5
9	9	c	78.3

Consider modeling the response variable  $Y$  by the following linear model:

$$\mathcal{M}_{12} : Y_i = \beta_0 + \beta_1 x_{i1} + \alpha_j + \gamma_j x_{i1} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2),$$

where index  $j$  is related to the categories of  $X_2$ . The model  $\mathcal{M}_{1|2}$  can be written in matrix form as

$$\mathcal{M}_{1|2} : \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{I}.$$

Write in details what kind forms the model matrix  $\mathbf{X}$  and parameter vector  $\boldsymbol{\beta}$  have in case of given data is modeled by the model  $\mathcal{M}_{12}$ .

(2 points)

- (b) Below is the R estimation output of the `lm`-function related to the particular linear model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ .

```

Residuals:
    Min       1Q   Median       3Q      Max
-5.8753 -1.8275 -0.0943  2.1809  7.2335

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  18.0026      1.5837   11.367 6.02e-14 ***
factor(x1)2  -1.7752      1.5259   -1.163  0.2517
factor(x1)3  -2.9361      1.4336   -2.048  0.0473 *
factor(x2)2   1.2917      1.3836    0.934  0.3563
factor(x2)3   0.3726      1.4915    0.250  0.8040
factor(x2)4  -3.8796      1.5543   -2.496  0.0169 *
---

Residual standard error: 3.388 on 39 degrees of freedom
Multiple R-squared:  0.3396,    Adjusted R-squared:  0.2549
F-statistic: 4.011 on 5 and 39 DF,  p-value: 0.004953

```

- i. What kind of linear model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  the output is related to?
  - A.  $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i$ ,
  - B.  $Y_i = \beta_0 + \beta_1 x_{i1} + \alpha_j + \varepsilon_i$ ,
  - C.  $Y_i = \beta_0 + \beta_j + \alpha_h + \varepsilon_i$ ,
  - D.  $Y_i = \beta_0 + \beta_j + \alpha_h + \gamma_{jh} + \varepsilon_i$ .
- ii. Calculate the maximum likelihood estimate of the expected value  $\mu$  when the explanatory variables  $X_1, X_2$  are set on the values

$$x_1 = 3,$$

$$x_2 = 3.$$

(2 points)

- (c) Consider the linear model

$$\mathbf{y} \sim N(\boldsymbol{\mu}, \sigma^2 \mathbf{I}),$$

$$\boldsymbol{\mu} = \mathbf{1}\beta_0,$$

where  $\mathbf{1}$  is a vector of ones  $\mathbf{1} = (1, 1, \dots, 1)'$ . The sample mean

$$\bar{y} = \frac{y_1 + y_2 + \dots + y_n}{n} = \frac{1}{n} \mathbf{1}' \mathbf{y}$$

is the maximum likelihood estimator for the parameter  $\beta_0$ , i.e.,  $\hat{\beta}_0 = \bar{y}$ . Make yourself familiar with Theorem 1.1 in section 1.3.2 Multivariate Normal Distribution and then calculate the expected value  $E(\hat{\beta}_0)$  and the variance  $\text{Var}(\hat{\beta}_0)$ .

(2 points)