

STAT 510 Homework 12
Due Date: 11:00 A.M., Wednesday, May 1

1. An experiment was conducted at 15 research stations around the country to determine how dose of a chemical mixture affects the leaf area of a certain type of plant. Instructions for creating the chemical mixture and for carrying out the experiment were sent to the managers at each of the 15 research stations. At each station, a completely randomized design was used to assign 5 doses of the chemical (0, 25, 50, 75, and 100 mL/day) to 20 plants with 4 plants per dose. Each plant grew in its own pot and received its assigned chemical dose each day of the experiment. At three weeks of age, the leaf area of each plant was recorded. The data for this experiment are available in the file

<http://dnett.github.io/S510/LeafArea.txt>

For $i = 1, \dots, 15$, $j = 1, \dots, 5$, and $k = 1, \dots, 4$, let y_{ijk} be the leaf area for the k th plant that received dose j in research station i , and suppose

$$y_{ijk} = (\beta_1 + b_{1i}) + (\beta_2 + b_{2i})x_j + e_{ijk} \quad (1)$$

In this model (1), β_1 and β_2 are unknown parameters, $x_1 = 0, x_2 = 25, x_3 = 50, x_4 = 75, x_5 = 100$, the e_{ijk} terms are *iid* $N(0, \sigma_e^2)$, and the b_{1i} and b_{2i} terms are normal random effects independent of the e_{ijk} terms. More specifically, let

$$\mathbf{b}_i = \begin{bmatrix} b_{1i} \\ b_{2i} \end{bmatrix} \text{ for all } i = 1, \dots, 15.$$

We assume $\mathbf{b}_1, \dots, \mathbf{b}_{15} \stackrel{iid}{\sim} N(\mathbf{0}, \Sigma_b)$ for some positive definite 2×2 variance matrix Σ_b .

Model (1) is a special case of what is sometimes referred to as a *random coefficient model* because the regression coefficients are assumed to be random variables rather than fixed parameters. It is straightforward to fit such a model in R using code like the following.

```
d = read.delim("http://dnett.github.io/S510/LeafArea.txt")
library(lme4)
o = lmer(LeafArea ~ Dose + (1 + Dose | ResearchStation), data = d)
```

The approximate BLUEs of β_1 and β_2 can be obtained with code like

```
fixef(o)
```

As usual, the estimated variance of the estimator of the fixed effects parameters is given by `vcov(o)`.

The empirical BLUPs of b_{1i} and b_{2i} can be obtained with code like

```
ranef(o)
```

Typing `summary(o)` provides you with enough information to determine the REML estimates of σ_e^2 and Σ_b . The estimate of the matrix Σ_b is not provided directly, but you can compute it from the given estimates of the variances and the provided estimate of the correlation between b_{1i} and b_{2i} labeled `Corr` in the `Random effects` portion of the output.

- (a) Provide the REML estimate of σ_e^2 .
 - (b) Provide the REML estimate of Σ_b .
 - (c) Make a scatterplot of leaf area vs. dose for the data from the 7th research station. Add a black line to the plot that shows the estimate of the regression function $\beta_1 + \beta_2 x$ for $x \in (0, 100)$.
 - (d) Find the prediction of the regression function for the 7th research station; i.e., predict $(\beta_1 + b_{17}) + (\beta_2 + b_{27})x$ for $x \in (0, 100)$.
 - (e) Using only the data from the 7th research station, find the ordinary least squares estimate of the regression function for the simple linear regression of leaf area on dose of the chemical.
 - (f) To the plot in part (c), add a red line that shows the regression function predicted in part (d) and a blue line that shows the regression function estimated in part (e).
 - (g) Compute the likelihood ratio statistic for testing $H_0 : \beta_2 = 0$.
 - (h) Find AIC for the fit of model (1) to the data from all 15 research stations.
 - (i) Find AIC for a simplified version of model (1) that assumes there is one slope coefficient common to all research stations.
 - (j) Find AIC for a simplified version of model (1) that assumes there is one intercept coefficient common to all research stations and one slope coefficient common to all research stations.
 - (k) According to AIC, which model is preferred among model (1), the model considered in part (i), and the model considered in part (j).
2. Model (1) from problem 1 can be written in linear mixed-effects model form as $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$. Define \mathbf{X} , $\boldsymbol{\beta}$, \mathbf{Z} , \mathbf{u} , $\mathbf{G} = \text{Var}(\mathbf{u})$, and $\mathbf{R} = \text{Var}(\mathbf{e})$ using terms from model (1). Assume the response vector \mathbf{y} is ordered as in the dataset LeafArea.txt.
 3. In 1846, a group of pioneers traveling west became stranded in the eastern Sierra Nevada mountains. By the time the last survivor was rescued in the spring of 1847, 40 of 87 members in the original group had died from starvation and exposure to extreme cold. The group became known as the Donner Party. The dataset

<http://dnett.github.io/S510/Donner.txt>

contains the age, sex, and status (survived or died) of the members of the group that were 15 years of age or older. Conduct an analysis of this data set to determine how age and sex are associated with the probability of survival. Support your answer with appropriate tests and/or confidence intervals. State your conclusions in ways that will be easily interpretable by nonstatisticians.

4. Consider an experiment with three treatments (A , B , and C). Suppose there are 10 experimental units for each of treatments A and B . Suppose there are 50 experiment units for treatment C . Imagine that the response for each experimental unit has a binomial distribution with $m = 20$ trials (same for all experiment units) and a success probability that depends on treatment. Suppose that (unknown to the researcher) the success probabilities for treatments A , B , and C are 0.5, 0.5, and 0.95, respectively. Rather than using logistic regression for analysis, a researcher decides to use a standard three-treatment ANOVA assuming a normal response $[\text{lm}(\mathbf{y} \sim \text{trt})]$. The researcher is primarily interested in a comparison of treatments A and B , so he examines the R output for the coefficients to see if the “trtB” (the name R would use) coefficient is significant because he knows that provides a test for the difference in treatment A and B means due to the set first to zero constraints. Explain why this might not be a safe analysis strategy.