# Homework 7
Due Thursday, 3/25/21 on Canvas.

1. Let $p_1, \ldots, p_m$ be p-values for null hypotheses $H_{0,1}, \ldots, H_{0,m}$, where $p_j \sim U[0, 1]$ if $H_{0,j}$ is true. Define $\pi_0 \in [0, 1]$ to be the fraction of the $m$ hypotheses that are true.

    (a) Show that the Bonferroni procedure controls the FWER regardless of the dependence between $p_1, \ldots, p_m$. That is, show that for any $\alpha \in (0, 1)$,

    $$P\left(\exists j = 1, \ldots, m \text{ such that } H_{0,j} \text{ is true and } p_j \leq \alpha/m\right) \leq \pi_0 \alpha.$$

    (b) If $p_1, \ldots, p_m$ are independent, show that the Bonferroni procedure provides exact control of the FWER for small $\alpha$. That is, show that for any $m \geq 1$ and significance level $\alpha \in (0, 1)$,

    $$P\left(\exists j = 1, \ldots, m \text{ such that } H_{0,j} \text{ is true and } p_j \leq \alpha/m\right) = \pi_0 \alpha \{1 + o(1)\}$$

    as $\alpha \to 0$.

    (c) Is (b) necessarily true if $p_1, \ldots, p_m$ are dependent? If yes, prove it. If not, find a counterexample.

2. (Preliminaries for estimation in random effects models) Let $Y \in \mathbb{R}^n$ be a random vector with $\mathbb{E}(Y) = 0$ and $\text{Var}(Y) = \Sigma(\theta) \in \mathbb{R}^{n \times n}$, where $\theta \in \mathbb{R}^p$ is an unknown parameter. Assume that $\Sigma(\theta)$ is continuously differentiable as a function of $\theta$, i.e. there exist continuous matrix functions $M_j(\theta) \in \mathbb{R}^{n \times n}$ for $j = 1, \ldots, p$ such that

    $$\Sigma(\theta + \delta) - \Sigma(\theta) = \sum_{j=1}^{p} \delta_j M_j(\theta) + o\left(\|\delta\|_2\right).$$

    Let

    $$\ell(\theta) = -\frac{1}{2} \log\left[\det\{\Sigma(\theta)\}\right] - \frac{1}{2} Y^T \{\Sigma(\theta)\}^{-1} Y$$

    be the log-likelihood for the normal distribution (up to additive constants that do not depend on $\theta$). Note that we are NOT assuming $Y$ is normally distributed.

    (a) Let $V \in \mathbb{R}^{n \times n}$ be a symmetric, positive definite matrix. For any symmetric matrix $A \in \mathbb{R}^{n \times n}$, prove the following:

    $$\log\{\det(V + \epsilon A)\} - \log\{\det(V)\} = \epsilon \,\text{Tr}\left(AV^{-1}\right) + o(\epsilon)$$
    $$u^T (V + \epsilon A)^{-1} u - u^T V^{-1} u = -\epsilon u^T V^{-1} A V^{-1} u + o(\epsilon)$$

    for $\epsilon > 0$ and $u \in \mathbb{R}^n$.

(b) Use part (a) to show that

$$[\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})]_j = -\frac{1}{2} \text{Tr}\left[M_j(\boldsymbol{\theta})\{\Sigma(\boldsymbol{\theta})\}^{-1}\right] + \frac{1}{2} Y^T \{\Sigma(\boldsymbol{\theta})\}^{-1} M_j(\boldsymbol{\theta})\{\Sigma(\boldsymbol{\theta})\}^{-1} Y.$$

Conclude that a root of $\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})$ is a suitable estimator for $\boldsymbol{\theta}$. That is, show that

$$\mathbb{E}\{\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})\} = 0$$

regardless of whether or not $Y$ is normally distributed.

3. A study was done to compare the yields of 56 varieties of wheat in a randomized complete block design (RCBD) with four blocks of size 56. The data for the experiment are in the file "wheat56.txt". The four blocks are observations 1-56, 57-112, 113-168 and 169-224. The varieties, yields, latitudes and longitudes of each plot (latitudes and longitudes in unstated units) are given. Although the units are unstated, keep in mind that agricultural field trials like this are carried out at a single farm so that the weather is essentially the same at all plots. The labeling of the varieties as 1-56 "in order" in Block 1 is for convenience; you may assume that in fact the variety assignments were properly randomized in all four blocks.

(a) Estimate the variety effects using the standard model for an RCBD treating blocks and varieties as fixed effects. Using appropriate tables and/or plots, summarize your findings about the differences between varieties. As part of your analysis, include an F-test for the hypothesis of no variety effects.

(b) Find a 95% confidence interval for the mean yield of varieties 1-20 minus the mean yield of varieties 21-56.

(c) Plot the residuals as a function of the geographic coordinates of the plots. Discuss any patterns you see and comment on the reasonableness of the assumptions underlying the analyses in (a). Can you identify any varieties whose yields (relative to other varieties) might be over or underestimated because of the plots to which they were assigned? Comment.

(d) Reanalyze the data including a linear function of the coordinates in your mean function. What effect does this change have on your inferences about variety effects? In particular, which estimated variety effects change the most from the analysis in (c)? Plot the residuals as a function of the geographic coordinates of the plots. To what extent are any problems you noted with the residual plot in (c) fixed?

(e) Answer the same questions as in (d), but this time including a quadratic function of the coordinates (i.e., a second order polynomial in latitude and longitude) in your mean function.

(f) Do you think the design used for this study was well-chosen? Discuss any problems you see and describe how the study might have been designed differently to avoid or reduce these problems.

4. **PhD problem**: Proof of Sheffé. Let $X \in \mathbb{R}^{n \times p}$ be a full rank design matrix and $Y \sim N\left(X\beta, \sigma^2 I_n\right)$ for some $\beta \in \mathbb{R}^p$. Let $\mathcal{S} = \left\{c \in \mathbb{R}^p : \sum_{j=1}^p c_j = 0\right\}$.

(a) If $\hat{\beta}$ is the the ordinary least squares estimate for $\beta$, show that for any $\alpha \in (0, 1)$,

$$P\left[\left\{\frac{c^T \hat{\beta} - c^T \beta}{se\left(c^T \hat{\beta}\right)}\right\}^2 \leq q_{1-\alpha} \forall c \in \mathcal{S}\right] \leq \alpha$$

where $se\left(c^T \hat{\beta}\right) = \sqrt{\hat{\sigma}^2 c^T \left(X^T X\right)^{-1} c}$ and $q_{1-\alpha}$ is the $1-\alpha$ quantile of the $(p-1)F_{p-1, n-p}$ distribution.

(b) Use (a) to derive simultaneous $1-\alpha$ confidence intervals for every point in $\left\{c^T \beta : c \in \mathcal{S}\right\}$.