

Assignment 3 — Due Oct 18, 2018

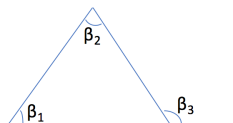
1. Consider the model $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ where $\varepsilon_i \sim \text{iid } N(0, \sigma_i^2)$ for observations $i = 1, \dots, n$. Suppose $X_i > 0$ and $\sigma_i^2 = \sigma^2 X_i$. Let $\hat{Q}(\beta) = \sum_{i=1}^n \sigma_i^{-2} \{Y_i - (\beta_0 + \beta_1 X_i)\}^2$ denote the weighted error sum of squares. Define

$\mathbf{Y} = (Y_1, \dots, Y_n)'$: $n \times 1$ vector of response variables

\mathbf{X} : $n \times 2$ design matrix with 1's in the first column and $(X_1, \dots, X_n)'$ in the second column.

$\beta = (\beta_0, \beta_1)'$: 2×1 vector of regression coefficients

- Let $\tilde{\beta} = (\tilde{\beta}_0, \tilde{\beta}_1)'$ denote the weighted least squares estimates of β . Derive the distribution (i.e., the type of distribution, mean vector, and the variance-covariance matrix) of $\tilde{\beta}$.
 - Despite the weighted variances, derive the ordinary least squares estimates of β by minimizing the (unweighted) error sum of squares $Q(\beta) = \sum_{i=1}^n \{Y_i - (\beta_0 + \beta_1 X_i)\}^2$ all in matrix terms.
 - Let $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)'$ denote the (ordinary) least squares estimates of β you found in part (b), and let $\hat{e} = (\hat{e}_1, \dots, \hat{e}_n)' = \mathbf{Y} - \mathbf{X}\hat{\beta}$ denote the raw residuals under ordinary least squares estimates. Find the distribution of $\hat{\beta}$ and the variance of \hat{e}_i for $i = 1, \dots, n$.
 - Draw a connection between $\text{Var}(\hat{\beta}_1)$ and $\text{Var}(\tilde{\beta}_1)$.
2. A chemist studied the concentration (Y) of a solution over time (X). Fifteen (15) identical solutions were prepared. The 15 solutions were randomly divided into five sets of three, and the five sets were measured, respectively, after 1, 3, 5, 7, and 9 hours. The data file is `solution.txt`. online submission
- Identify ingredients of the experiment/study including population vs. sample, study/experimental units, and whether cause-and-effect relationships can be established.
 - (Hand calculation with calculator) Perform a simple linear regression analysis with concentration (Y) as the response variable and time (X) as the explanatory variable. Obtain the regression coefficients estimates along with their standard errors, an unbiased estimate of error variance, and the coefficient of determination R^2 .
 - (Hand calculation with calculator) Perform a hypothesis test to determine whether there is evidence that the mean concentrations are different for the different hours since the solutions are prepared. Indicate the assumptions underlying the test.
 - Indicate the model underlying the regression analysis in (b) and assess the model assumptions by using suitable graphical techniques. How reasonable are the assumptions? What remedial measures are desirable, if any?
 - What kind of variable transformations would you recommend to the chemist? Give reasoning.
 - Repeat (b) but now with the concentration after the transformation in (e).
 - Compare the results in (b) and (f).
3. Suppose we wish to measure the three angles β_1 , β_2 , and β_3 as depicted in the diagram below.



Elementary geometry shows that $\beta_1 + \beta_2 = \beta_3$. Suppose, as a check on the accuracy of the results, we decide to measure all three angles, with measurement error. Let b_j be the actual measurement for β_j , $j = 1, 2, 3$. Due to measurement error, $b_1 + b_2$ might not be equal to b_3 . Assume that the measurement errors are independent and follow normal distribution with mean 0 and variance σ^2 . Formulate this as a multiple linear regression model, and derive the least squares estimates $\hat{\beta}_1$, $\hat{\beta}_2$, and $\hat{\beta}_3$, with their standard deviations.

4. Consider the multiple regression model

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_{p-1} X_{p-1} + \varepsilon,$$

where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. Suppose σ^2 is unknown and consider n independent observations from the model. Let $\mathbf{Y} = (Y_1, \dots, Y_n)'$ denote the observed responses; let $\mathbf{X} = [\mathbf{x}_0, \dots, \mathbf{x}_{p-1}]$ denote the $n \times p$ design matrix, where \mathbf{x}_i corresponds to the $(i+1)^{th}$ column of \mathbf{X} ; and let $\boldsymbol{\beta} = (\beta_0, \dots, \beta_{p-1})'$ denote the vector of regression coefficients. Assume that $\text{rank}(\mathbf{X}) = p$. Let $\hat{\beta}_j$ denote the least squares estimate of β_j , and define $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \dots, \hat{\beta}_{p-1})'$.

- Suppose $\mathbf{c} = (c_0, \dots, c_{p-1})'$ is a vector of known constants. What is the distribution of $\sum_{j=0}^{p-1} c_j \hat{\beta}_j$? Justify your answer, and find its mean and variance using matrix notation.
- Consider the hypothesis $H_0 : \sum_{j=0}^{p-1} c_j \beta_j = h$ and $H_1 : \sum_{j=0}^{p-1} c_j \beta_j \neq h$, where h is a given constant. Explain how to test these hypotheses at significance level α . Construct a suitable test statistic, find its distribution, and specify the rejection region.
- Suppose we wish to predict the value of a future observation Y_{n+1} . Let $\mathbf{z} = (1, z_1, \dots, z_{p-1})'$ denote its corresponding vector of predictor variables (i.e., $X_i = z_i$, for $i = 1, \dots, p-1$), and consider the prediction $\hat{Y}_{n+1} = \mathbf{z} \hat{\boldsymbol{\beta}}$. Find the distribution of $Y_{n+1} - \hat{Y}_{n+1}$.
- Show that the MSE (mean square error) of the prediction $\hat{Y}_{n+1} = \mathbf{z} \hat{\boldsymbol{\beta}}$ is strictly greater than σ^2 .
- Given the vector \mathbf{z} of predictor variables for the future observation Y_{n+1} , find an interval \mathcal{I} such that $\mathbb{P}(Y_{n+1} \in \mathcal{I}) = 1 - \alpha$.
- Suppose an additional predictor variable X_p is added to the model to obtain

$$Y = \gamma_0 + \gamma_1 X_1 + \cdots + \gamma_{p-1} X_{p-1} + \gamma_p X_p + \varepsilon.$$

Suppose that the augmented design matrix $\widetilde{\mathbf{X}} = [\mathbf{X}, \mathbf{x}_p]$ has rank $p+1$. Let $\hat{\mathbf{x}}_p = a_0 \mathbf{x}_0 + \cdots + a_{p-1} \mathbf{x}_{p-1}$ denote the least-squares projection of \mathbf{x}_p onto the subspace of \mathbb{R}^n spanned by the columns of \mathbf{X} . Find an expression for the residual vector $\mathbf{r}_p = \mathbf{x}_p - \hat{\mathbf{x}}_p$ in terms of \mathbf{X} and \mathbf{x}_p .

- Express the least squares estimates $\hat{\gamma}_0, \dots, \hat{\gamma}_p$ in terms of only $\hat{\beta}_0, \dots, \hat{\beta}_{p-1}, a_0, \dots, a_{p-1}, \mathbf{r}_p$, and \mathbf{Y} .