

STA 302f13 Assignment Nine¹

This assignment assumes you are using the [Formula sheet](#). There is a link on the course home page in case the one in this document does not work. The formula sheet (or part of it) will be provided with the quiz. **Bring your printouts for Question 9 to the quiz, including the plots.**

1. For the general linear regression model, show that the square of the sample correlation between Y and \hat{Y} values is equal to R^2 .
2. This question compares the error terms ϵ_i to the residuals $\hat{\epsilon}_i$. Answer True or False to each statement. For statements about the residuals, show a calculation that proves your answer. You may use anything on the formula sheet.
 - (a) $E(\epsilon_i) = 0$
 - (b) $E(\hat{\epsilon}_i) = 0$
 - (c) $Var(\epsilon_i) = 0$
 - (d) $Var(\hat{\epsilon}_i) = 0$
 - (e) ϵ_i has a normal distribution.
 - (f) $\hat{\epsilon}_i$ has a normal distribution.
 - (g) $\epsilon_1, \dots, \epsilon_n$ are independent.
 - (h) $\hat{\epsilon}_1, \dots, \hat{\epsilon}_n$ are independent.
3. One of these statements is true, and the other is false. Pick one, and show it is true with a quick calculation. Start with something from the formula sheet.
 - $\hat{Y} = \mathbf{X}\hat{\beta} + \hat{\epsilon}$
 - $\mathbf{Y} = \mathbf{X}\hat{\beta} + \hat{\epsilon}$

As the saying goes, “Data equals fit plus residual.”

4. The *deleted residual* is $\hat{\epsilon}_{(i)} = Y_i - \mathbf{x}'_i \hat{\beta}_{(i)}$, where $\hat{\beta}_{(i)}$ is defined as usual, but based on the $n - 1$ observations with observation i deleted.
 - (a) Guided by an expression on the formula sheet, write the formula for the Studentized deleted residual. You don't have to prove anything. You will need the symbols $\mathbf{X}_{(i)}$ and $s_{(i)}$, which are defined in the obvious way.
 - (b) If the model is correct, what is the distribution of the Studentized deleted residual? Make sure you have the degrees of freedom right.
 - (c) Why are numerator and denominator independent?
5. For the general linear regression model, are \hat{Y} and $\hat{\epsilon}$ independent?
 - (a) Answer Yes or No and prove your answer.
 - (b) What does this imply about the plot of predicted values against residuals?

¹Copyright information is at the end of the last page.

6. For the general linear regression model, are \mathbf{Y} and $\hat{\mathbf{Y}}$ independent? Answer Yes or No and prove your answer.
7. For the general linear regression model, are \mathbf{Y} and $\hat{\boldsymbol{\epsilon}}$ independent? Answer Yes or No and prove your answer.
8. Prove that the sample correlation between residuals and independent variable values must equal exactly zero. Does this result depend on the correctness of the model?
9. Lecture slide set 6 used the `trees` data. Typing `help(trees)` at the R prompt gives more information. For this question, bring your R printouts to the quiz, *including the plots*.
 - (a) Fit an ordinary model with two independent variables. How much of the variability in Volume is explained? You have to admit, that's pretty good.
 - (b) Now let's look at the deleted Studentized residuals. One student made an excellent suggestion, which was to look at boxplots. Try `boxplot(vaname)`, where `vaname` is the name of the deleted Studentized residual. If you don't know what a boxplot is, look in the Wikipedia. This part is interesting, but it will not be on the quiz. Do you see one possible high outlier.
 - (c) Now treat the deleted Studentized residuals as t -test statistics, with a Bonferroni correction to achieve a *joint* significance level of 0.05. What is the critical value? It's a number on your R printout. This *could* be on the quiz.
 - (d) Is there evidence of outliers? Answer yes or No.
 - (e) Now plot predicted values against standardized residuals. Put a title on the plot. See `help(title)`. Do you see anything fishy, or perhaps wavy?
 - (f) Now plot the independent variables in the model against the standardized residuals. It's a bit subjective, but when I do this I see a curvilinear trend for one independent variable, but not for the other. Which one?

Then I thought about it for a while. Finally, combining a bit of geometry with what little I know about trees, I came up with a model. This model has *one* independent variable, a function of Height and Girth, and it explains almost 98% of the variation in volume. The residual plots look pretty clean. Can you guess my model?

This assignment was prepared by [Jerry Brunner](#), Department of Statistical Sciences, University of Toronto. It is licensed under a [Creative Commons Attribution - ShareAlike 3.0 Unported License](#). Use any part of it as you like and share the result freely. The L^AT_EX source code is available from the course website: <http://www.utstat.toronto.edu/~brunner/oldclass/302f13>