

Problem Set 02
Linear Models -- Fall 2015
Due date: 2015-09-30

Problems sets are due at the start of class on the due date. Please hand write or type up and print the solutions; we will not accept e-mail solution sets except in exceptional circumstances. You may discuss problem sets with others, but must write up your own solutions. This means that you should have no need to look at other's final written solutions. Many of these problems come from a variety of textbooks, which are referenced in the problems. These are for citation purposes and not because you will need to consult the text itself (though you may feel free to do so).

1. Consider the case of a simple linear regression (no intercept) with a random design; specifically, assume:

$$x_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1) \quad (1)$$

$$\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2) \quad (2)$$

$$\mathbb{E}(\epsilon_i|x) = 0 \quad (3)$$

$$y_i = x_i\beta + \epsilon \quad (4)$$

Define $\tilde{\beta} = \frac{1}{n} \sum_i x_i y_i$.

(a) Give both an argument and counterargument for using $\tilde{\beta}$ in lieu of $\hat{\beta}$. Hint: What is $\frac{1}{n} \sum_i x_i^2$ an estimate of?

(b) Calculate $\mathbb{E}(\tilde{\beta}|X)$. Is this estimator unbiased when conditioned on X ? Is it unbiased when calculating the unconditional expectation?

(c) Compute the unconditional variance of $\tilde{\beta}$ and compare to the unconditional variance of $\hat{\beta}$. Which estimator would you rather use?

2. [Fumio Hayashi, pg 32, #6] Prove that, under assumptions I-IV, $Cov(\hat{\beta}, r|X) = 0$. The covariance can be written as:

$$Cov(\hat{\beta}, r|X) = \mathbb{E} \left\{ (\hat{\beta} - \beta)(r - \mathbb{E}(r|X))^t \middle| X \right\} \quad (5)$$

Hint: Use the relations $M\epsilon = r$ and $\hat{\beta} - \beta = A\epsilon$.

3. [Fumio Hayashi, pg 46, #7] Prove that, under assumptions I-V, the estimator s^2 has variance:

$$\mathbb{V}(s^2|X) = \frac{2\sigma^4}{n-p} \quad (6)$$

4. [Fumio Hayashi, pg 74, #5] In the restricted least squares, the sum of squared residuals is minimized subject to the constraint implied by the null hypothesis $D\beta = d$. Form the Lagrangian

as:

$$\mathcal{L}(b, \lambda) = \frac{1}{2}(y - Xb)^t(Y - Xb) + \lambda^t(Db - d) \quad (7)$$

Where λ is a k -dimensional vector. Let $\tilde{\beta}$ be the solution to the restricted regression.

(a) Let $\hat{\beta}$ be the standard (unrestricted) OLS estimator. Show:

$$\tilde{\beta} = \hat{\beta} - (X^t X)^{-1} D^t [D(X^t X)^{-1} D^t]^{-1} (D\hat{\beta} - d) \quad (8)$$

$$\lambda = [D(X^t X)^{-1} D^t]^{-1} (D\hat{\beta} - d) \quad (9)$$

(b) Let $\tilde{e} = y - X\tilde{\beta}$, the residuals from the restricted regression. Show that:

$$SSR_R - SSR_U = (\hat{\beta} - \tilde{\beta})^t (X^t X) (\hat{\beta} - \tilde{\beta}) \quad (10)$$

$$= (D\hat{\beta} - d)^t [D(X^t X)^{-1} D^t]^{-1} (D\hat{\beta} - d) \quad (11)$$

$$= \lambda^t D(X^t X)^{-1} D^t \lambda \quad (12)$$

$$= \tilde{e}^t P \tilde{e} \quad (13)$$

In other words, justify each step.

(c) Show how this verifies that the two F-statistics are the same:

$$\frac{(SSR_R - SSR_U)/k}{SSR_U/(n - p)} = \frac{(D\hat{\beta} - d)^t [D(X^t X)^{-1} D^t]^{-1} (D\hat{\beta} - d)/k}{s^2} \quad (14)$$

5. [Sheather 2009, pg 147] Chateau Latour is widely acknowledged as one of the world's greatest wine estates with a rich history dating back to at least 1638. The Grand Vin de Chateau Latour is a wine of incredible power and longevity. At a tasting in New York in April 2000, the 1863 and 1899 vintages of Latour were rated alongside the 1945 and the 1961 vintages as the best in a line-up of 39 vintages ranging from 1863 to 1999 (*Wine Spectator*, August 31, 2000). Quality of a particular vintage of Chateau Latour has a huge impact on price. For example, in March 2007, the 1997 vintage of Chateau Latour could be purchased for as little as \$159 per bottle while the 2000 vintage of Chateau Latour costs as least \$700 per bottle (<http://www.wine-searcher.com>).

While many studies have identified that the timing of the harvest of the grapes has an important effect on the quality of the vintage, with quality improving the earlier the harvest. A less explored issue of interest is the effect of unwanted rain at vintage time on the quality of icon wine like Chateau Latour. This question addresses this issue.

The Chateau Latour web site (<http://www.chateau-latour.com>) provides a rich source of data. In particular, data on the quality of each vintage, harvest dates and weather at harvest time were obtained from the site for the vintages from 1961 to 2004. An example of the information on weather at harvest time is given below for the 1994 vintage:

Harvest began on the 13th September and lasted on the 29th, frequently interrupted by storm showers. But quite amazingly the dilution effect in the grapes was very limited ...(<http://www.chateau-latour.com/commentaires/1994uk.html>)

Each vintage was classified as having had "unwanted rain at harvest" (e.g., the 1994 vintage) or not (e.g., the 1996 vintage) on the basis of information like that reproduced above. Thus, the data consist of:

Vintage = year the grapes were harvested

Quality – on a scale from 1 (worst) to 5 (best) with some half points

End of harvest – measured as the number days since August 31

Rain – a dummy variable for unwanted rain at harvest = 1 if yes.

The data can be found at euler.stat.yale.edu/~tba3/psets/pset02/data/Latour.txt.
The first model considered was:

$$\text{Quality} = \beta_0 + \beta_1 \cdot \text{End of Harvest} + \beta_2 \text{Rain} + \beta_3 \text{End of Harvest} \times \text{Rain} + \epsilon \quad (15)$$

(a) Show that the coefficient of the interaction term in model ?? is statistically significant at the 0.05 level. In other words, show that the rate of change in quality rating depends on whether there has been any unwanted rain at vintage.

(b) Estimate the number of days of delay to the end of harvest it takes to decrease the quality rating by 1 point when there is:

(i) No unwanted rain at harvest

(ii) Some unwanted rain at harvest

(c) Construct a confidence interval for the (total) end of harvest slope variable when there is unwanted rain.