

# Homework 9

Nathan Bick

## 1. Problems 13.F:1

1. Solve the system of two equations in the two unknowns  $R_0$  and  $R_1$  in the proof of Lemma 13.3 and show that the solutions are as given in lemma 13.3.

$$\text{Lemma 13.3: } \hat{R}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n y_i \sum_{i=1}^n x_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}$$

$$\text{and } \hat{R}_0 = \bar{y} - \hat{R}_1 \bar{x}.$$

The system of equations is:

$$nR_0 + \left(\sum_{i=1}^n x_i\right)R_1 = \sum_{i=1}^n y_i$$

$$\left(\sum_{i=1}^n x_i\right)R_0 + \left(\sum_{i=1}^n x_i^2\right)R_1 = \sum_{i=1}^n x_i y_i.$$

first we see

$$\begin{cases} \hat{R}_0 = \sum y_i - (\sum x_i) \hat{R}_1 = \bar{y} - \hat{R}_1 \bar{x} \\ (\sum x_i) \hat{R}_0 + (\sum x_i^2) \hat{R}_1 = \sum x_i y_i \end{cases}$$

$$\Rightarrow (\sum x_i) \left( \frac{\sum y_i - (\sum x_i) \hat{R}_1}{n} \right) + (\sum x_i^2) \hat{R}_1 = \sum x_i y_i.$$

$$\Rightarrow \frac{\sum x_i \sum y_i - (\sum x_i)^2 \hat{R}_1}{n} + \frac{n}{n} (\sum x_i^2) \hat{R}_1 = \sum x_i y_i$$

$$\Rightarrow \frac{n(\sum x_i^2) R_1 - (\sum x_i)^2 R_1}{n} = \sum x_i y_i - \frac{\sum x_i \sum y_i}{n}$$

$$\Rightarrow (n \sum x_i^2 - (\sum x_i)^2) R_1 = n \sum x_i y_i - \sum x_i \sum y_i$$

$$\Rightarrow R_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$



## 2. Exercises 14.3.1

1. Compare Cases A, B, and C of example 14.2. Which case is best? why?

In example 14.2, we see three cases of one-stage clustering. In all three cases we get the same  $E(\bar{y}_{clus}) = \frac{9}{2} = \bar{y}_u$ , so each case shows an unbiased and effective estimator. However, each case has very different  $V(\bar{y}_{clus})$ .

$$A: V(\bar{y}_{clus}) = 9, \quad B: V(\bar{y}_{clus}) = \frac{1}{4}, \quad C: V(\bar{y}_{clus}) = 0.$$

C is better than B, which is better than A. We see if the clusters are each internally heterogeneous among the  $y$  values,  $V(\bar{y}_{clus})$  will be less than when each cluster is homogeneous.

C is best because our sampling variance is low.