

ECON 178 WI 2021: Homework 1

Due: Feb 1, 2021 (by 12:30pm PT)

Instructions:

- The homework has a total of 40 points. The TAs will randomly pick one problem to grade and this problem is worth 30 points (you will get 30 points if your answers are correct or almost correct). The remaining 10 points will be graded on completion of this assignment.
- There will be two separate submissions: one for your R code and one for your writeup. Please submit both on **Gradescope** (more details for the submission of the R part are given in “Applied questions”).
- Please follow the policy stated in the syllabus about academic integrity.
- You must read, understand, agree and sign the integrity pledge (<https://academicintegrity.ucsd.edu/forms/form-pledge.html>) before completing any assignment for ECON178. After you sign the pledge form, a receipt will be emailed to you. Please include this receipt in the submission of your writeup (not the code) on **Gradescope**.

Conceptual questions

The following questions involve reviewing exercises about expectations, conditional expectations, biases and variances, and basic properties of Normal (also called Gaussian) distributions.

Question 1

Suppose that we have a model $y_i = \beta x_i + \epsilon_i$ ($i = 1, \dots, n$) where $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = 0$, $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 0$, and ϵ_i is distributed normally with mean 0 and variance σ^2 ; that is, $\epsilon_i \sim N(0, \sigma^2)$. Furthermore, $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are independently distributed, and the x_i s ($i = 1, \dots, n$) are non-random.

- (a) The OLS estimator for β minimizes the Sum of Squared Residuals:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left[\sum_{i=1}^n (y_i - \beta x_i)^2 \right]$$

Take the first-order condition to show that

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}.$$

- (b) Assume $\mathbb{E}[\epsilon_i | \beta] = 0$ for all $i = 1, \dots, n$. Show that

$$\hat{\beta} = \beta + \frac{\sum_{i=1}^n x_i \epsilon_i}{\sum_{i=1}^n x_i^2}$$

What is $\mathbb{E}[\hat{\beta} | \beta]$ and $\text{Var}(\hat{\beta} | \beta)$? Use this to show that, conditional on β , $\hat{\beta}$ has the following distribution:

$$\hat{\beta} | \beta \sim N\left(\beta, \frac{\sigma^2}{\sum_{i=1}^n x_i^2}\right).$$

- (c) Suppose we believe that β is distributed normally with mean 0 and variance $\frac{\sigma^2}{\lambda}$; that is, $\beta \sim N(0, \frac{\sigma^2}{\lambda})$. Additionally assume that β is independent of ϵ_i for all $i = 1, \dots, n$. Compute the mean and variance of $\hat{\beta}$. That is, what is $\mathbb{E}[\hat{\beta}]$ and $\text{Var}(\hat{\beta})$?

(Hint you might find useful: $\mathbb{E}[w_1] = \mathbb{E}[\mathbb{E}[w_1 | w_2]]$ and $\text{Var}(w_1) = \mathbb{E}[\text{Var}(w_1 | w_2)] + \text{Var}(\mathbb{E}[w_1 | w_2])$ for any random variables w_1 and w_2 .)

- (d) Since everything is normally distributed, it turns out that

$$\mathbb{E}[\beta | \hat{\beta}] = \mathbb{E}[\beta] + \frac{\text{Cov}(\beta, \hat{\beta})}{\text{Var}(\hat{\beta})} \cdot (\hat{\beta} - \mathbb{E}[\hat{\beta}]).$$

Let $\hat{\beta}^{RR} = \mathbb{E}[\beta | \hat{\beta}]$. Compute $\text{Cov}(\beta, \hat{\beta})$ and use the value of $\mathbb{E}[\beta]$ along with the values of $\mathbb{E}[\hat{\beta}]$, $\text{Cov}(\beta, \hat{\beta})$, and $\text{Var}(\hat{\beta})$ you have computed to show that

$$\hat{\beta}^{RR} = \mathbb{E}[\beta | \hat{\beta}] = \frac{\sum_{i=1}^n x_i^2}{\sum_{i=1}^n x_i^2 + \lambda} \cdot \hat{\beta}$$

(Hint: $\text{Cov}(w_1, w_2) = \mathbb{E}[(w_1 - \mathbb{E}[w_1])(w_2 - \mathbb{E}[w_2])]$ and $\mathbb{E}[w_1 w_2] = \mathbb{E}[w_1 \mathbb{E}[w_2 | w_1]]$ for any random variables w_1 and w_2)

- (e) Does $\hat{\beta}^{RR}$ increase or decrease as λ increases? How does this relate to β being distributed $N(0, \frac{\sigma^2}{\lambda})$?

Question 2

Let us consider the linear regression model $y_i = \beta_0 + \beta_1 x_i + u_i$ ($i = 1, \dots, n$), which satisfies Assumptions MLR.1 through MLR.5 (see Slide 7 in “Linear_regression_review” under “Modules” on Canvas)¹. The x_i s ($i = 1, \dots, n$) and β_0 and β_1 are nonrandom. The randomness comes from u_i s ($i = 1, \dots, n$) where $\text{var}(u_i) = \sigma^2$. Let $\hat{\beta}_0$ and $\hat{\beta}_1$ be the usual OLS estimators (which are unbiased for

β_0 and β_1 , respectively) obtained from running a regression of $\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{n-1} \\ y_n \end{pmatrix}$ on $\begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix}$ (the intercept

column) and $\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{pmatrix}$. Suppose you also run a regression of $\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{n-1} \\ y_n \end{pmatrix}$ on $\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{pmatrix}$ only

(excluding the intercept column) to obtain another estimator $\tilde{\beta}_1$ of β_1 .

¹The model is a simple special case of the general multiple regression model in “Linear_regression_review”. Solving this question does not require knowledge about matrix operations.

- a) Give the expression of $\tilde{\beta}_1$ as a function of y_i s and x_i s ($i = 1, \dots, n$).
- b) Derive $\mathbb{E}(\tilde{\beta}_1)$ in terms of β_0 , β_1 , and x_i s. Show that $\tilde{\beta}_1$ is unbiased for β_1 when $\beta_0 = 0$. If $\beta_0 \neq 0$, when will $\tilde{\beta}_1$ be unbiased for β_1 ?
- c) Derive $\text{Var}(\tilde{\beta}_1)$, the variance of $\tilde{\beta}_1$, in terms of σ^2 and x_i s ($i = 1, \dots, n$).
- d) Show that $\text{Var}(\tilde{\beta}_1)$ is no greater than $\text{Var}(\hat{\beta}_1)$; that is, $\text{Var}(\tilde{\beta}_1) \leq \text{Var}(\hat{\beta}_1)$. When do you have $\text{Var}(\tilde{\beta}_1) = \text{Var}(\hat{\beta}_1)$? (Hint you might find useful: use $\sum_{i=1}^n x_i^2 \geq \sum_{i=1}^n (x_i - \bar{x})^2$ where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$.)
- e) Choosing between $\hat{\beta}_1$ and $\tilde{\beta}_1$ leads to a tradeoff between the bias and variance. Comment on this tradeoff.

Question 3

Let \hat{v} be an estimator of the truth v . Show that $\mathbb{E}(\hat{v} - v)^2 = \text{Var}(\hat{v}) + [\text{Bias}(\hat{v})]^2$ where $\text{Bias}(\hat{v}) = \mathbb{E}(\hat{v}) - v$. (Hint: The randomness comes from \hat{v} only and v is nonrandom).

Applied questions (with the use of R)

For this question you will be asked to use tools from R for coding.

Installation

- To install R, please see <https://www.r-project.org/>.
- Once you install R, please install also R Studio <https://rstudio.com/products/rstudio/download/>.
- You will need to use R Studio to solve the problem set.

Download

from Canvas \Rightarrow Assignments

- data_ps1.csv;
- template_ps1.R .

Submission

- Open the template_ps1.R file that we provided on Canvas \Rightarrow Assignments.
- All your solutions and code need to be saved in a single file named template_ps1_YOURFIRSTANDLASTNAME.R file. Please use the template_ps1.R provided in Canvas to structure your answers.
- **Any file that is not an .R will not be accepted, and the grade for this exercise will be zero.**
- Please submit your code on **Gradescope**.
- Please follow the policy stated in the syllabus about academic integrity.

Useful readings

In addition to the lectures provided by the instructor and the TAs, you might find the following readings useful:

- Chapter 2.3 and 3.6 in the textbook "An introduction to statistical learning with applications in R".

Question 4

We want to predict which variables are the most correlated with the balance in a bank account. To do so we use the credit data set (Dua, D. and Graff, C., 2019, UCI Machine Learning Repository) available on Canvas \Rightarrow Assignments.

1. Download the dataset from Canvas and open it using the command "read.csv".
2. Open the data and report how many columns and rows the dataset has;
3. See the names of the variables (see online the command "names");
4. Run a linear regression with Balance as a function of Income using the command "lm";
5. Report the summary of your results (see online the command "summary")
7. Plot a scatter plot of the regression (Hint: use abline() to draw the regression line)
8. Write down the interpretation of the coefficients as a comment in your .R script (Hint: see template file).

Please write all your answer and code in template_ps1.R file and submit that file on Gradescope as described in the "Submission" section.