

Homework 2

1 Directions:

- **Due: Thursday February 10, 2022, at 9pm.** Late submissions will be accepted for 24 hours after that time, with a 15% penalty. (The enforcement is strict, beginning at 9:01pm, except for extreme situations; having a poor wifi connection or minor computer problems is not sufficient for the penalty to be waived.)
- Upload the homework to Canvas as a single pdf file.
- If the graders cannot easily read your submission (writing is illegible, image is too dark, or if the contrast is too low) then you might receive a zero or only partial credit.
- Any non-administrative questions must be asked in office hours or (if a brief response is sufficient) Piazza.

2 Problems

Problem 1. [13 points total (6,3,4)]

Book problem Chapter 3, problem 3 “Suppose we have a data set with five predictors,
 $X_1 = \text{GPA}, \dots$ ”

Note: for interactions, use products, e.g., $X_4(i) = \text{GPA}(i) \times \text{IQ}(i)$

$$Y^\wedge = 50 + 20 * \text{GPA} + 0.07 * \text{IQ} + 35 * \text{Level} + 0.01 * \text{GPA} * \text{IQ} - 10 * \text{GPA} * \text{Level}$$

$$\therefore \text{For High School: } Y^\wedge = 50 + 20 * \text{GPA} + 0.07 * \text{IQ} + 0.01 * \text{GPA} * \text{IQ} \dots *$$

$$\text{For College: } Y^\wedge = 85 + 10 * \text{GPA} + 0.07 * \text{IQ} + 0.01 * \text{GPA} * \text{IQ} \dots **$$

(a) When High School > College then,

$$50 + 20 * \text{GPA} > 85 + 10 * \text{GPA}$$

$$\Leftrightarrow \text{GPA} > 3.5$$

\therefore if GPA is high enough high school graduates earn more

iii is correct

(b) college/IQ=110/GAP=4

$\therefore **$

$$Y^\wedge = 85 + 10 * 4 + 0.07 * 110 + 0.01 * 4 * 110$$

$$= 85 + 40 + 7.7 + 4.4$$

$$= \$137.1k$$

(c) False because to verify it, we need see GPA/IQ has an impact. Also, It is possible to have many evidences although the coefficient is very small.

Problem 2. [12 points total (3 points each)]

Book problem Chapter 3, problem 4 “I collect a set of data . . .”

- (a) Here the true relationship is linear, then the cubic regression would have more noise although it is hard to tell without knowing the actual training data.
∴ We would expect the RSS for the linear regression will be lower than the one for cubic regression.
- (b) We do not have enough information to conclude this because the test RSS depends on the test data. However, we might be able to expect cubic regression will have a higher test RSS due to the overfit from the training because the true relationship is linear.
- (c) We would expect cubic regression will have lower RSS than the linear one because the true relationship is more flexible.
- (d) We do not have enough information to conclude this because we do not know how far it is from linear. Therefore, it is not clear enough to know what level of flexibility would fit better.

Problem 3. [5 points]

Suppose we have a data set with one feature X to predict another feature Y . Let n denote the number of samples. Let \bar{X} and \bar{Y} denote the average values of X and Y respectively in the data set

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X(i) \qquad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y(i).$$

Let β_0^* and β_1^* denote the coefficients for the ordinary least squares (O.L.S.) solution,

$$\{\beta_0^*, \beta_1^*\} = \arg \min_{\{\beta_0, \beta_1\}} \frac{1}{n} \sum_{i=1}^n \left(Y(i) - (\beta_0 + \beta_1 X(i)) \right)^2.$$

Their values are

$$\beta_0^* = \bar{Y} - \beta_1^* \bar{X} \qquad \beta_1^* = \frac{\sum_{i=1}^n (X(i) - \bar{X})(Y(i) - \bar{Y})}{\sum_{i=1}^n (X(i) - \bar{X})^2}$$

Using those formulas, calculate the O.L.S. model's prediction for $X = \bar{X}$ (i.e., the prediction \hat{Y} for a new sample whose X feature has the value \bar{X}).

$$\beta_0^* = \bar{Y} - \beta_1^* \bar{X} \dots *$$

$$Y = \beta_0^* + \beta_1^* X$$

Plug in *

$$= \bar{Y} - \beta_1^* \bar{X} + \beta_1^* X$$

$$= \bar{Y} (\because X = \bar{X})$$

$$\Leftrightarrow \bar{Y} = Y$$

Problem 4. [16 points total (3,3,10)]

Book problem Chapter 6, problem 1 “We perform best subset ...”

Notes regarding the book’s pseudocode for Algorithms 6.1-6.3:

- “RSS” stands for “residual sum of squares” which is the (un-normalized) MSE,

$$\text{RSS} = \sum_{i=1}^n \left(Y(i) - \hat{Y}(i) \right)^2$$

In the book’s pseudo-codes, “RSS” refers to *training set* RSS.

- “Cross-validated prediction error” – you can read this as “Validation set MSE.”
 - “ R^2 ” and “adjusted R^2 ” – you can ignore these for this homework
- a) The best subset with k predictors is the model with the smallest RSS because it considers all possible models with k predictor.
- b) With given information, it is not easy to answer this question. Best subset might overfits if n is smaller than p and other two might choose better on the test set.
- c)
- i. true
 - ii. true
 - iii. false
 - iv. false
 - v. false