

Homework 2

Due Date

March 2nd at 4pm

Tricky Questions

State whether you agree or disagree with the following statements, and explain your reasoning.

- Removing an outlier or high leverage point always increases R^2 .
- If the correlation matrix between all independent variables in a regression model has an off-diagonal element near 1, then that indicates that at least one pair of independent variables are *collinear* with each other
- The numerical values chosen for a dummy variable do not impact the performance of the regression model

RABE 3.3

A teacher has created a dataset containing the scores on a final examination F , as well as the scores in two preliminary examinations P_1 and P_2 for 22 students in a statistics course. The data can be found on Canvas under **Files>data>exams.csv**.

- Fit each of the following models to the data:

$$\text{Model 1: } F_i = \beta_0 + \beta_1 P_{1i} + \epsilon_i$$

$$\text{Model 2: } F_i = \beta_0 + \beta_2 P_{2i} + \epsilon_i$$

$$\text{Model 3: } F_i = \beta_0 + \beta_1 P_{1i} + \beta_2 P_{2i} + \epsilon_i$$

- Which variable individually, P_1 or P_2 is a better predictor of F ?
- Which of the three models would you use to predict the final examination scores for a student who scored $P_1 = 78$ and $P_2 = 85$? What is your prediction in this case?

RABE 3.14 + 4.7

A national insurance organization wanted to study the consumption of cigarettes in all 50 states and the District of Columbia. The data from 1970 are available on Canvas under **Files>data>cigarettes.csv**, and the variable definitions are given in the table below. For parts (a) and (b) below, specify the null and alternative hypotheses, the test used, and your conclusion using a significance $\alpha = .05$.

| Variable | Definition |
|----------|--|
| AGE | Median of the state's population |
| HS | Percentage of people over 25 years of age in a state who had completed high school |
| INCOME | Per capita personal income for a state (in dollars) |
| FEMALE | Percentage of population identified as "female" |
| PRICE | Average price (in cents) of a pack of cigarettes in the state |

| Variable | Definition |
|----------|--|
| SALES | Number of packs of cigarettes sold in a state per capita |

- Test the hypothesis that the variable FEMALE is not needed in the regression equation relating SALES to the five predictor variables
- Test the hypothesis that the variables FEMALE and HS are not needed in the above regression equation
- Compute that 95% Confidence Interval for the true regression coefficient of the variable INCOME.
- What percentage of the variation in SALES can be accounted for by the three variables PRICE, AGE, and INCOME?
- Using an added variable plot, show the effect of including the INCOME variable
- What percentage of the variation in SALES can be accounted for when INCOME is removed from the above regression?
- Compute the pairwise correlation coefficients matrix and construct the corresponding scatter plot matrix.
- Are there any disagreements between the pairwise correlation coefficients and the corresponding scatter plot matrix?
- Is there any difference between your expectations in part (a) and what you see in the pairwise correlation coefficients matrix, or in the scatter plot matrix?

RABE 4.1a

Using the milk production dataset (on canvas under **Files>data>milk_production.csv**, described in RABE pages 3-4), fit the following model:

$$\text{CURRMILK} = \beta_0 + \beta_1 \text{PREVIOUS} + \beta_2 \text{FAT} + \beta_3 \text{PROTEIN} + \beta_4 \text{DAYS} + \beta_5 \text{LACTAT} + \beta_6 \text{I79}$$

Now, for your fit model determine:

- If the regression assumptions (linearity and iid normal errors) are met
- If any outliers are present in the data
- If any linear dependence exists between the independent variables