

Statistics 201
PRACTICE MIDTERM EXAM KEY

Note that pages have been condensed on this key to fit on 3 pages, to save paper if you print it.

Open notes. Calculator required. There are 5 problems, with a total of 14 parts. Each part of each problem (a, b, etc) is worth 7 points, except Problem 5a, which is worth 9 points.

1. Problem 1.28 on page 37 of your textbook describes data from 84 medium-sized counties in the US. For each county, X = percentage of adults in the county having at least a high-school diploma, and Y = crime rate (crimes reported per 100,000 residents) last year. Here is some R output from fitting a simple linear regression model to the data:

```
>lm(formula = Crime ~ Diploma, data = Midterm)

            Estimate Std. Error t value Pr(>|t|)
(Intercept) 20517.60   3277.64    6.260 1.67e-08 ***
Diploma      -170.58    41.57   -4.103 9.57e-05 ***

Residual standard error: 2356 on 82 degrees of freedom
Multiple R-squared:  0.1703,    Adjusted R-squared:  0.1602
F-statistic: 16.83 on 1 and 82 DF,  p-value: 9.571e-05

>anova(Midterm)
Analysis of Variance Table

Response: Crime
          Df Sum Sq Mean Sq F value    Pr(>F)
Diploma    1 93462942 93462942  16.834 9.571e-05 ***
Residuals 82 455273165  5552112
```

a. Write the *population* version of the regression model.

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \text{or} \quad E\{Y_i\} = \beta_0 + \beta_1 X_i$$

b. Write the *estimated* (sample) regression function.

$$\hat{Y}_i = 20517.6 - 170.58 X_i \quad \text{or} \quad Y_i = 20517.6 - 170.58 X_i + e_i$$

c. Interpret the slope in the context of this situation.

If two counties differ by 1% in adults with a high-school diploma, on average the crime rate would be 170.58 lower (per 100,000 residents) for the county with the higher percentage of high school diplomas.

d. According to the last US Census, 82.7% of Orange County adults have a high school diploma. Round this number to 83%, and obtain a point estimate for the crime rate in Orange County.

$$\hat{Y}_h = 20517.6 - 170.58 X_h = 20517.6 - 170.58 (83) = 6359.46.$$

This is the predicted number of crimes reported per 100,000 residents.

Problem 1, continued...

- e. A 95% confidence interval for $E\{Y_h\}$ when $X_h = 70$ is 7702 to 9453. Interpret this interval in words, in the context of this situation.

For the population of all medium-sized counties in which 70% of adults have a high-school diploma, we are 95% certain that the population mean crime rate last year was between 7702 and 9453 crimes reported per 100,000 residents.

- f. The results indicate that counties with higher percentages of high-school graduates tend to have lower crime rates. Can we conclude from this study that having a high school diploma causes people to be less likely to commit crimes, in other words, that higher high-school graduation rates cause crime to be lower? Explain your answer.

No. The data clearly come from an observational study, because you can't randomly assign different counties to have specific high-school graduation rates. There are many possible confounding variables that could influence crime rates, such as income levels, county services provided, etc. You can't make cause and effect conclusions in observational studies.

2. Define \mathbf{I} to be an $n \times n$ identity matrix, and \mathbf{H} to be the usual hat matrix. A matrix that plays a useful role in regression inference is $(\mathbf{I} - \mathbf{H})$. Show using matrix algebra that $(\mathbf{I} - \mathbf{H})$ is idempotent. You can use the fact that \mathbf{H} is idempotent.

First, you're told that you can use the fact that \mathbf{H} is idempotent, so $\mathbf{H}\mathbf{H} = \mathbf{H}$. Also, $\mathbf{I}\mathbf{I} = \mathbf{I}$, and anything multiplied by \mathbf{I} is itself, so $\mathbf{H}\mathbf{I} = \mathbf{I}\mathbf{H} = \mathbf{H}$.

So, $(\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H}) = \mathbf{I}\mathbf{I} - \mathbf{H}\mathbf{I} - \mathbf{I}\mathbf{H} + \mathbf{H}\mathbf{H} = \mathbf{I} - \mathbf{H} - \mathbf{H} + \mathbf{H} = \mathbf{I} - \mathbf{H}$.

3. A company offers a training course for the Math SAT. They give their students a test at the end of the course, graded from 0 to 100. They would like to use that test in the future to predict how well students will score on the Math SAT. They have scores on their test and the Math SAT for a sample of students. Thus, X = score on the company's test and Y = score on the Math SAT, which ranges from 200 to 800. They plan to use the usual simple linear regression model.

- a. Would the intercept have a useful meaning in this example? Explain your answer.

It might. If some students score at or near 0 on the company's test, then the intercept would be the predicted SAT score for students who score 0 on the company test. If all students score much higher, then the intercept would not have a useful meaning.

- b. One of the company analysts states that the intercept should be fixed at 200, because that's the lowest the SAT Math score can be. Suppose the intercept is set to 200 for this situation. Write the population model.

$$Y_i = 200 + \beta_1 X_i + \varepsilon_i \quad \text{or} \quad E\{Y_i\} = 200 + \beta_1 X_i$$

c. Write the full and reduced models to test whether or not it makes sense to set the intercept to be 200.

$$\text{Full model: } Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \text{or} \quad E\{Y_i\} = \beta_0 + \beta_1 X_i$$

$$\text{Reduced model: } Y_i = 200 + \beta_1 X_i + \varepsilon_i \quad \text{or} \quad E\{Y_i\} = 200 + \beta_1 X_i \quad \text{or} \quad Y_i - 200 = \beta_1 X_i + \varepsilon_i$$

d. Write the sum that is to be minimized to get the least squares regression line, if the model you wrote in Part b is used.

$$Q = \sum(\varepsilon_i)^2 = \sum(Y_i - 200 - \beta_1 X_i)^2$$

4. What assumption is being examined by looking at a normal probability plot? Be specific.

The assumption that the error terms ε in the population are normally distributed.

5. A regression equation is to be fit for predicting Y = resting pulse rate using the predictor variables X_1 = number of minutes of exercise per week and X_2 = gender, with 1 = male and 0 = female. Here are the X values results for 6 individuals:

Exercise/week	200	10	420	50	350	140
Gender	Male	Female	Female	Male	Male	Female

a. (9 points) Write down the \mathbf{X} matrix that would be used for this situation, filling in numerical values.

$$\mathbf{X} = \begin{bmatrix} 1 & 200 & 1 \\ 1 & 10 & 0 \\ 1 & 420 & 0 \\ 1 & 50 & 1 \\ 1 & 350 & 1 \\ 1 & 140 & 0 \end{bmatrix}$$

b. Explain in words what the coefficient attached to X_2 represents.

The coefficient would be present in the model for males, and absent in the model for females, so it represents the average difference in pulse rates for males and females with number of minutes of exercise per week held constant. Notice that the model forces us to assume that the difference between males and females is constant across the range of exercise amounts.