

Stat 21 Homework 7

Name

Due: Sunday, March 27th by midnight

Contents

Part I: Concept problems	1
Problem 2	2
Problem 3	2
Problem 4	2
Problem 5	2
Part II: R Problems	3
Problem 6	3
Problem 7	4
Problem 8	4
Problem 9	4
Problem 10	5

Use this file as the template for your submission. Do not delete anything from this template unless you are prompted to do so (e.g. where to write your name above, where to write your solutions or code below). Make sure you have installed the following packages in your version of RStudio: **tidyverse**, **knitr** **before** you attempt to knit this document.

Your completed assignment should be submitted as a single **PDF** using the link under Week 8 titled “Submit HW 7 to Gradescope”. You must use R markdown to write up your solutions. For any homework problems that involve coding in R, you must provide **both** the code and a written answer interpreting the output within the context of the problem. You are allowed to work with your classmates on this homework assignment but you are expected to write up your own solutions. Every answer must be supported by a written statement unless otherwise specified. *A good rule of thumb is to make sure your answer is understandable to someone who hasn’t read the problem question (or code output associated with it).*

Additionally, make sure that when you upload your solutions to Gradescope, you select which pages correspond with which questions. Also, check to make sure that your knitted homework document is not uploaded as an extra-long single page document. Failure to do these things will result in a penalty on your homework grade. Finally, I strongly recommend that you address and resolve any knitting or R coding issues before Saturday as solutions to any R-coding questions that are not knitted properly will not receive any credit.

Part I: Concept problems

For problems 1-2 consider this regression model was fit to a sample of breakfast cereals. The response variable Y is calories per serving. The predictor variables are X_1 = grams of sugar per serving, and X_2 = grams of fiber per serving. The fitted regression model is

$$\hat{Calories} = 109.3 + 1.0Sugar - 3.7Fiber.$$

Problem 1

- (a) How many calories would you predict for a breakfast cereal that had 1 gram of fiber and 11 grams of sugar per serving?
- (b) Frosted Flakes is a breakfast cereal that has 1 gram of fiber and 11 grams of sugar per serving. It also has 110 calories per serving. Compute the residual for Frosted Flakes and explain what this value means.

Solution:

- (a) [Write your answer here]
- (b) [Write your answer here]

Problem 2

- (a) Does the prediction equation for number of calories per serving suggest that the amount of sugar has a weaker relationship with the number of calories than the amount of fiber? Explain why or why not.
- (b) In the context of this setting, interpret -3.7 the coefficient of X_2 . That is, describe how fiber is related to calories per serving, in the presence of the sugar variable.

Solution:

- (a) [Write your answer here]
- (b) [Write your answer here]

For problems 3-5 read the article, “Scientists rise up against statistical significance” at <https://www.nature.com/articles/d41586-019-00857-9>.

Problem 3

The article claims, “. . . researchers have been warned that a statistically non-significant result does not ‘prove’ the null hypothesis (the hypothesis that there is no difference between groups or no effect of a treatment on some measured outcome).” Explain why failing to reject the null hypothesis does not prove that there is no effect. What does failing to reject the null hypothesis really mean instead?

Solution: [Write your answer here]

Problem 4

- (a) In the graphic “Beware false conclusions”, results are shown from two studies: one that found “significant” results, and another that found “non-significant” results. The article claims that it is “ludicrous” to say that the second study found “no association.” Briefly explain why this is the case.
- (b) Regarding the same two studies in part (a), the article claims that it is “absurd” to say that the two studies are in conflict, even though one was “significant” and the other was “not significant”. Briefly explain why this is the case.

Solution:

- (a) [Write your answer here]
- (b) [Write your answer here]

Problem 5

In the section titled “Quit categorizing”, the article claims that, “Statistically significant estimates are biased. . . Consequently, any discussion that focuses on estimates chosen for their significance will be biased.” Briefly explain why this is the case.

Solution: [Write your answer here]

Part II: R Problems

Problem 6

In 2016 Hillary Clinton won the Democratic nomination for president over Bernie Sanders. A paper was circulated that claimed to show evidence of election fraud based, among other things, on Clinton doing better in states that don't have a paper trail for votes cast in a primary election than she did in states that have a paper trail. The file `ClintonSanders` has data from that paper for the 31 states that held primaries before June.

```
data(ClintonSanders)
ClintonSanders %>% head
```

##	State	Delegates	PaperTrail	PopularVote	AfAmPercent
## 1	Alabama	83.02	Paper Trail	77.8	26.2
## 2	Arizona	56.00	Paper Trail	56.3	4.1
## 3	Arkansas	68.75	No Paper Trail	66.1	15.4
## 4	Connecticut	50.91	Paper Trail	51.8	10.1
## 5	Delaware	57.14	No Paper Trail	59.8	21.4
## 6	Florida	65.89	No Paper Trail	64.4	16.0

The variable `Delegates` gives the percentage of delegates won by Clinton for each state. The variable `AfAmPercent` gives the percentage of residents in the state who are African American. `PaperTrail` indicates whether or not the voting system in the state includes a paper trail.

- Conduct a regression of `Delegates` on `PaperTrail`. What does this regression say about how Clinton did in states with and without a paper trail?
- Conduct a regression of `Delegates` on `PaperTrail` and `AfAmPercent`. What does this regression say about how Clinton did in states with and without a paper trail? What is the effect of `AfAmPercent`?
- Repeat parts (a) and (b) but in place of `Delegates` as the response variable, use `PopularVote`, which is the percentage of the popular vote that Clinton received. Do any important conclusions change when using `PopularVote` as the response variable instead?

Solution:

```
## Use this space for your solution to part (a)
```

- [Write your solution to part (a) here.]

```
## Use this space for your solution to part (b)
```

- [Write your solution to part (b) here.]

```
## Use this space for your solution to part (c)
```

- [Write your solution to part (c) here.]

It seems reasonable to predict the number of calories (per serving) in breakfast cereals using the amount of sugar (grams per serving). The file `Cereal` also has a variable showing the amount of fiber (grams per serving) for each of the 36 cereals. Use this data below for problems 7-8.

```
data(Cereal)
Cereal %>% head
```

##	Cereal	Calories	Sugar	Fiber
----	--------	----------	-------	-------

## 1	Common Sense Oat Bran	100	6	3
## 2	Product 19	100	3	1
## 3	All Bran Xtra Fiber	50	0	14
## 4	Just Right	140	9	2
## 5	Original Oat Bran	70	5	10
## 6	Heartwise	90	5	6

Problem 7

Fit a multiple regression model to predict **Calories** based on both predictors: **Sugar** and **Fiber**. Examine each of the measures below and identify which (if any) of the cereals you might classify as possibly “unusual” in that measure. Include specific numerical values and justification for each case.

- (a) Standardized residuals
- (b) Studentized residuals

Solution:

Use this space for your solution to part (a)

- (a) [Write your solution to part (a) here.]

Use this space for your solution to part (b)

- (b) [Write your solution to part (b) here.]

Problem 8

Fit a multiple regression model to predict **Calories** based on both predictors: **Sugar** and **Fiber**. Examine each of the measures below and identify which (if any) of the cereals you might classify as possibly “unusual” in that measure. Include specific numerical values and justification for each case.

- (a) Leverage
- (b) Cook’s D

Solution:

Use this space for your solution to part (a)

- (a) [Write your solution to part (a) here.]

Use this space for your solution to part (b)

- (b) [Write your solution to part (b) here.]

Problem 9

Two types of dementia are Dementia with Lewy Bodies and Alzheimer’s disease. Some people are afflicted with both of these. The file **LewyBody2Groups** includes the variable **Type**, which has two levels: “DLB/AD” for the 20 subjects with both types of dementia and “DLB” for the 19 subjects with only Lewy Body dementia. The variable **APC** gives annualized percentage change in brain gray matter. The variable **MMSE** measures change in functional performance on the Mini Mental State Examination.

```
data("LewyBody2Groups")
LewyBody2Groups %>% head
```

##	Type	APC	MMSE
## 1	DLB	0.85	2.22
## 2	DLB	0.49	0.37

```
## 3 DLB 0.12 -0.10
## 4 DLB 0.00 -2.99
## 5 DLB -0.22 0.66
## 6 DLB -0.35 -2.47
```

- Fit an interaction model that produces two regression lines for predicting MMSE from APC, one for each of the two levels of **Type**. Write down the fitted prediction equation for each level of **Type**.
- Use a t-test to test the null hypothesis that the interaction term is not needed and parallel regression lines are adequate. Specify the null and alternative, the p-value and your chosen significance level in addition to the conclusion of the test.
- Use a nested F-test to test the null hypothesis that neither of the terms involving **Type** is needed and a common regression line for both levels of **Type** is adequate for modeling how MMSE depends on APC. Specify the null and alternative, the p-value and your chosen significance level in addition to the conclusion of the test.

Solution:

```
## Use this space for your solution to part (a)
```

- [Write your solution to part (a) here.]

```
## Use this space for your solution to part (b)
```

- [Write your solution to part (b) here.]

```
## Use this space for your solution to part (c)
```

- [Write your solution to part (c) here.]

Problem 10

Consider the data introduced in Exercise 1.45 on the time (in minutes) it took to play a sample of Major League Baseball games. The datafile **BaseballTimes2017** contains four quantitative variables (**Runs**, **Margin**, **Pitchers**, and **Attendance**) that might be useful in predicting the game times (**Time**).

```
data("BaseballTimes2017")
BaseballTimes2017 %>% head
```

```
##      Game League Runs Margin Pitchers Attendance Time
## 1 CHC-ARI     NL   11      5       10      39131   203
## 2 KCR-CHW     AL    9      3        7      18137   169
## 3 MIN-DET     AL   13      5       10      29733   201
## 4 SDP-LAD     NL    7      1        6      52898   179
## 5 COL-MIA     NL    9      3       10      20096   204
## 6 CIN-MIL     NL   21      1       10      34517   235
```

From among these four predictors choose a model for each of the goals below.

- Maximize the adjusted coefficient of determination.
- Minimize Mallows's C_p .
- After considering the models for parts (a) and (b), which model would you choose to predict baseball game times? Explain your choice.

Solution:

```
# Use this space for part (a)
```

- [Write your solution to part (a) here]

Use this space for part (b)

(a) [Write your solution to part (b) here]

(b) [Write your solution to part (c) here]