

14.310x Data Analysis for Social Scientists
Causality, Analyzing Random Experiments, and Nonparametric Regression

Welcome to your seventh homework assignment! You will have about one week to work through the assignment. We encourage you to get an early start, particularly if you still feel you need more experience using R. We have provided this PDF copy of the assignment so that you can print and work through the assignment offline. You can also go online directly to complete the assignment. If you choose to work on the assignment using this PDF, please go back to the online platform to submit your answers based on the output produced. Some of the questions we are asking are not easily solvable using math so we recommend you to use your R knowledge and the content of previous homework assignments to find numeric solutions.

Good luck!

Please find a glossary of R terms that will be useful for this week's homework here.

The following problems are based on the paper:

Duflo, Esther, Rema Hanna, and Stephen P. Ryan. 2012. "Incentives Work: Getting Teachers to Come to School." *American Economic Review*, 102(4): 1241-78.

In this experiment, the researchers set out to test whether providing teachers with cameras to take photos to prove their attendance could be effective in reducing teacher absenteeism. First, read the abstract of the paper using the link above. You can refer back to the paper as necessary.

Note: The dataset used to generate the Lecture 15 slides relating to this paper is slightly different than the dataset we have provided, so do not be alarmed if your answers are slightly different!

In order to complete this exercise, we are providing you with the code. The code has some missing parts that you have to fill in order to run it. The dataset that you will need is `teachers_final.csv`

Let's start thinking through how Fisher's ideas can be applied to evaluate this program in this context.

Question 1

First, consider the case where we have 8 schools. Our aim is to calculate the Fisher's exact p-value. Under the assumption that we will have the same number of treated and control units, how many potential treatment assignments across these 8 units are possible?

- 50
- 60
- 70
- 80

Suppose that after the treatment has been assigned and the experiment has been carried out, the researcher has the following data. The variable **open** corresponds to the fraction of days that the school was opened when random visits were made.

For Questions 2-4, we will look at these 8 schools found in `teachers_final.csv`:

treatment	open
0	0.462
1	0.731
0	0.571
0	0.923
0	0.333
1	0.750
1	0.893
1	0.692

Assume that we define as our statistic the absolute difference in means by treatment status.

To help you compute the test statistic for the observed data, we have provided you with the R code to load in this table and generate different permutations, although it is missing some parts that you will need to fill in. We make use of the package `perm`, specifically the function `ChooseMatrix`. Be sure to look up the documentation to make sure you understand what it is doing.

Question 2

For this observed data, what would be the value of our statistic?

We recommend you compute this test statistic on your own and then check your answer using the code provided. *Please round your answer to two decimal spaces.*

Question 3

According to your results, among the test statistics computed for all treatment assignments, how many are larger than the observed test statistic?

- 11
- 16
- 21
- 26
- 31
- 36

Question 4

What would be the Fisher's Exact p-value in this case? *Please round your answer to two decimal places.*

Question 5

Now load the data set `teachers_final.csv` in R and name it `schools`. This is done in line 31 of the provided R code.

With 49 schools treated, what are the number of possible assignments in this case?

- $\binom{49}{8}$
- $\binom{100}{51}$
- $\binom{1001}{49}$
- $\binom{100}{49}$

Question 6

A solution to this problem with a large number of observations is to simulate different random assignments and calculate the proportion of simulations in which the statistic exceeds the value of the observed data. We have provided you with the code that performs this exercise on the data `teachers_final.csv` with 100 simulations. However, we have replaced line 46 with blanks for you to fill in (XXXX).

Fill in line 46 with the correct code. What is the result of that line?

The figure you calculated in Question 6 represents an approximation to Fisher's p-value. You can explore with changing the number of simulations and the number of schools to see if they change the p-value.

Question 7

Since we are working in a very large sample, we can now consider Neyman's methods of inference. What is the Average Treatment Effect (ATE) on the observed data set? You will need to use R to compute this answer. *Please round your answer to three decimal places.*

Question 8

What is the upper bound of the standard error of this point estimate using Neyman's method?

(Hint: Use the conservative estimator of sampling standard deviation, $\sqrt{\hat{V}_{reymar}}$, as your upper bound.) *Please round your answer to three decimal places.*

Question 9

What is the t-statistic if we want to test the null hypothesis that ATE is equal to zero? *Please round your answer to two decimal places.*

Question 10

Is the associated p-value to this test similar to the one we found for the sharp null hypothesis in Question 6?

- Yes
- No

Question 11

The 95% confidence interval is given by (A, B). What are the values of A and B? *Please round to three decimal places. For instance, if your answer is .6789, please round to .679.*

Now, imagine that you are considering a randomized experiment similar to the camera experiment. The exception is that you plan to give teachers lower incentives: half the monetary amount that was given in the previous experiment.

Question 12

Imagine that the relationship between incentives and the variable **open** is linear. What would be the expected ATE of this new intervention? *Please round your answer to the third decimal place, i.e. if it is 0.3414, please round to 0.341.*

Question 13

Assume that the value from Question 12 is the minimum ATE such that the intervention is cost-effective. What is the sample size required to have a power of at least 90% with the following properties?

- with a significance level of 5%
- an equal number of treated and control units
- σ^2 is the average of the variance of the control and treatment group in the existing data

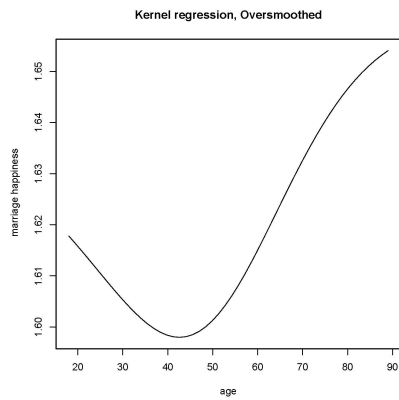
Hint: Recall that the formula for sample size is:

$$N = \frac{(\Phi^{-1}(1 - \beta) + \Phi^{-1}(1 - \frac{\alpha}{2}))^2}{\frac{\tau^2}{\sigma^2} \gamma(1 - \gamma)}$$

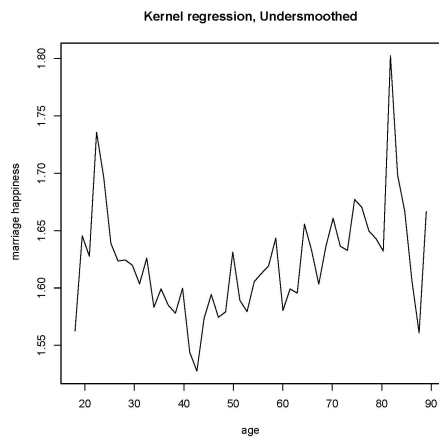
Where β is the operating characteristic and $1 - \beta$ is the desired power.

Now we are going to consider nonparametric regressions. The following plots show three different nonparametric regressions that relates the level of happiness in a marriage with age (where 2 corresponds to “very happy”, 1 to “pretty happy”, and 0 to “not too happy”).

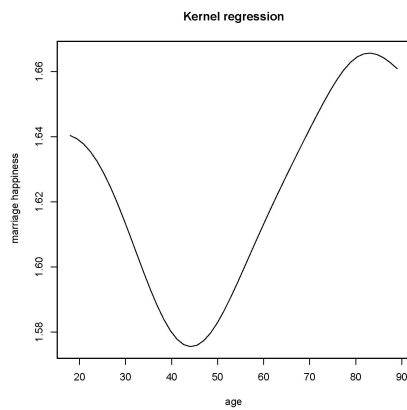
Plot A:



Plot B:



Plot C:



Question 14

Rank the three plots from the one with the narrowest to the widest bandwidth.

- ☐ a, c, b
- ☐ b, a, c

- c, b, a
- b, c, a
- c, a, b
- a, b, c

Going back to the data from `teachers_incentives.csv`, we are now going to focus on two variables:

- `pctpostwritten`, which denotes the mean student test scores after the intervention
- `open`

We want to see what the relationship between the fraction of days the school is open and student achievement. Use the code below (from lecture) to plot the kernel regression between these two variables using the R package `np`:

```
attach(schools)
plot <- npreg(xdat=XXX, ydat=XXX, bws=XXX,
bandwidth.compute=FALSE)
plot(plot)
```

Question 15

Use your code to generate plots for the following bandwidths. Which of them seems most appropriate given the data?

- 0.001
- 0.04
- 1
- 20

Question 16

Suppose we are interested in testing whether or not the distribution of the share of days a school is found to be open in the treatment group is statistically distinguishable from the distribution for that of the control group. Which of the following would be most useful for these purposes?

- Joint density plot
- Histogram of the variable by group
- Kolmogorov-Smirnov test
- Kernel regression
- None of the above

Question 17

Let $i \in T, C$ index the cohort school i assigned. m_i denotes the sample mean of a variable (e.g. student scores) for group i , μ_i denotes the population mean of the variable, and F_i denotes the CDF for group i

For each hypothesis test below, indicate which of the following methods is most useful for testing that hypothesis. **Enter N for using Neyman's method of inference, F for Fisher's exact test, and K for the KS test.**

- A. $H_0: \mu_T - \mu_C = 0$ vs. $H_1: \mu_T - \mu_C \neq 0$
- B. $H_0: \mu_T - \mu_C > 0$ vs. $H_1: \mu_T - \mu_C \leq 0$
- C. $H_0: m_T - m_C < 0$ vs. $H_1: m_T - m_C \geq 0$
- D. $H_0: F_T = F_C$ vs. $H_1: F_T \neq F_C$
- E. $H_0: F_T > G$ vs. $H_1: F_T \leq G$ where $G \sim N(0,1)$

Question 18

Use the R command `stat_ecdf()` to generate a plot of the CDFs for each cohort to see those results visually. Does the distribution of open in the treatment group FOSD that of the control group?

- Yes
- No

Note: the command to run a KS test in R is `ks.test()`. Look up the help file for this function and use it to test whether the distribution of test scores (`pctpostwritten`) in the treatment group first order stochastically dominates the distribution of test scores in the control group. Though the test fails when you have ties (so we are unable to use it in the case of test scores), you may find it useful in other applications.

Module 7

Matrix functions

- **chooseMatrix(n,m)**

Function to create a matrix of *choose(n,m)* rows and *n* columns. Each row has *m* 1s and the other elements are 0, arranged in a manner reflecting one possible selection.

- **NROW(x)/NCOL(x)**

Returns the number of rows or columns in matrix *x*.

Set seed

- **set.seed(x)**

Sets the seed of R's random number generator, which is useful for creating simulations or random objects that can be reproduced. The seed *x* is a number that is used to set the starting point for generating a series of random numbers. Each unique seed returns a unique random number sequence and once a seed is set, the same sequence can be produced repeatedly.

Question 1

0.0/1.0 point (graded)

First, consider the case where we have 8 schools. Our aim is to calculate the Fisher's exact p-value. Under the assumption that we will have the same number of treated and control units, how many potential treatment assignments across these 8 units are possible?

☐ 50

☐ 60

☒ 70 ✓

☐ 80

Explanation

As it was discussed in the lecture there are 8 units and 4 of them are going to be assigned to the treatment group. So in this case we will have that the total number of potential treatment assignments is $\binom{8}{4}$ which is equal to 70.

[Show answer](#)

Suppose that after the treatment has been assigned and the experiment has been carried out, the researcher has the following data. The variable **open** corresponds to the fraction of days that the school was opened when random visits were made.

For Questions 2-4, we will look at these 8 schools found in [teachers_final.csv](#):

treatment	open
0	0.462
1	0.731
0	0.571
0	0.923
0	0.333
1	0.750
1	0.893
1	0.692

Assume that we define as our statistic the absolute difference in means by treatment status.

To help you compute the test statistic for the observed data, we have provided you with the [R code](#) to load in this table and generate different permutations, although it is missing some parts that you will need to fill in. We make use of the package `perm`, specifically the function `ChooseMatrix`. Be sure to look up the documentation to make sure you understand what it is doing.

Question 2

0.0/1.0 point (graded)

For this observed data, what would be the value of our statistic?

We recommend you compute this test statistic on your own and then check your answer using the code provided.

Please round your answer to two decimal spaces.

Answer: 0.19425

Explanation

We have that:

$$\hat{\tau} = |\bar{\mathbf{Y}}_T^{obs} - \bar{\mathbf{Y}}_C^{obs}| = \left| \frac{Y_2^{obs} + Y_6^{obs} + Y_7^{obs} + Y_8^{obs}}{4} - \frac{Y_1^{obs} + Y_3^{obs} + Y_4^{obs} + Y_5^{obs}}{4} \right| = 0.19425$$

[Show answer](#)

Submit

You have used 0 of 2 attempts

i Answers are displayed within the problem

Question 3

0.0/1.0 point (graded)

Chrome

According to your results, among the test statistics computed for all treatment assignments, how many are larger than or equal to the observed test statistic?

☐ 11

☒ 16 ✓

☐ 21

☐ 26

☐ 31

☐ 36

Explanation

We can use a conditional function in R to test whether the statistic in other assignments exceeds or is equal to the value in the observed data. In particular we have that by running the code:

```
larger_than_observed <- (test_statistic >= observed_test) sum(larger_than_observed)
```

There are 16 assignments in which this is the case.

Question 4

0.0/1.0 point (graded)

What would be the Fisher's Exact p-value in this case?

Please round your answer to two decimal places.

Answer: 16/70

Explanation

In this case we know that the p-value is given by $16/70$ that equals ≈ 0.229 .

How should we interpret this? In general, we want to know whether the camera intervention had an effect, and whether the treated schools were open more frequently than the control schools. The mean of the treatment group is higher than the mean of the control group, indicating the teacher camera intervention may have indeed had an effect. However, under the sharp null hypothesis that there is no treatment effect in any of the schools in our sample, we have that if we randomly allocate 4 units to treatment, 23% of the time, the treatment and control groups would have looked at least as different as what we observed here, or even more different.

[Show answer](#)

Submit

You have used 0 of 2 attempts

Question 5

0.0/1.0 point (graded)

Now load the data set `teachers_final.csv` in R and name it `schools`. This is done in line 31 of the provided R code.

With 49 schools treated, what are the number of possible assignments in this case?

☐ $\binom{49}{8}$

☒ $\binom{100}{51}$ ✓ ✓ ✓

☐ $\binom{1001}{49}$

☒ $\binom{100}{49}$ ✓ ✓ ✓

Explanation

There are 100 schools in the data and 49 schools are treated. Thus, the answer is $\binom{100}{49}$.

[Show answer](#)

Question 6

0.0/1.0 point (graded)

A solution to this problem with a large number of observations is to simulate different random assignments and calculate the proportion of simulations in which the statistic exceeds the value of the observed data. We have provided you with the code that performs this exercise on the data `teachers_final.csv` with 100 simulations. However, we have replaced line 46 with blanks for you to fill in (XXXX).

Fill in line 46 with the correct code. What is the result of that line?

Answer: 0

Explanation

The correct code is

```
sum(abs(simul_stat) >= actual_stat)/NROW(simul_stat)
```

By running this code, you can see that by performing the 100 simulations, the p-value that we obtain is very close to 0.

[Show answer](#)

Submit

You have used 0 of 2 attempts

Question 7

0.0/1.0 point (graded)

Since we are working in a very large sample, we can now consider Neyman's methods of inference. What is the Average Treatment Effect (ATE) on the observed data set? You will need to use R to compute this answer.

Please round your answer to three decimal places.

Answer: 0.1969

Explanation

In this case we have that the ATE is given by:

$$\overline{\mathbf{Y}}_T^{obs} - \overline{\mathbf{Y}}_C^{obs} = 0.1969$$

Which is also the value of **actual_stat** in the code.

[Show answer](#)

Submit

You have used 0 of 2 attempts

Question 8

0.0/1.0 point (graded)

What is the upper bound of the standard error of this point estimate using Neyman's method? (Hint: Use the conservative estimator of sampling standard deviation, $\sqrt{\hat{V}_{\text{neyman}}}$, as your upper bound.)

Please round your answer to the three decimal places.

Answer: 0.031

Explanation

We need the estimated standard error, and will use the conservative estimator of sampling variance \hat{V}_{neyman} ,

$$\hat{V}_{\text{neyman}} = \frac{s_c^2}{N_c} + \frac{s_t^2}{N_t}$$

where

$$s_c^2 = \frac{1}{N_c - 1} \sum_{i: W_i = 0} (Y_i^{\text{obs}} - \bar{Y}_c^{\text{obs}})^2$$

and

$$s_t^2 = \frac{1}{N_t-1} \sum_{i:W_i=1} (Y_i^{obs} - \bar{Y}_t^{obs})^2$$

So,

$$\hat{V}_{\text{neyman}} = \frac{s_c^2}{N_c} + \frac{s_t^2}{N_t} = 0.03055^2$$

This operation is also performed with the R code we have provided.

[Show answer](#)

Submit

You have used 0 of 2 attempts

i Answers are displayed within the problem

Question 9

0/0/1 0 point (graded)

Question 9

0.0/1.0 point (graded)

What is the t-statistic if we want to test the null hypothesis that ATE is equal to zero?

Please round your answer to two decimal places.

Answer: 6.45

Explanation

In this case we will have that the t-statistic is:

$$t = \frac{\bar{Y}_T^{obs} - \bar{Y}_C^{obs}}{\sqrt{\hat{V}_{neyman}}} = \frac{0.1969}{0.03055} = 6.45$$

[Show answer](#)

Submit

You have used 0 of 2 attempts

Is the associated p-value to this test similar to the one we found for the sharp null hypothesis in Question 6?

☒ Yes ✓

☐ No

Explanation

The associated p-value to this test is $2 * (1 - \Phi(6.45)) \approx 0$, which is the same as the one we found in question 6.

[Show answer](#)

Submit

You have used 0 of 1 attempt

i Answers are displayed within the problem

Question 11

0.0/1.0 point (graded)

The 95% confidence interval is given by (A, B) . What are the values of A and B ?

Please round to three decimal places. For instance, if your answer is .6789, please round to .679.

0.0/1.0 point (graded)

The 95% confidence interval is given by (A, B) . What are the values of A and B ?

Please round to three decimal places. For instance, if your answer is .6789, please round to .679.

A

Answer: .137

B

Answer: .257

Explanation

The 95% CI is given by:

$$(0.1969 - 1.96 * 0.03055, 0.1969 + 1.96 * 0.03055) = (0.137, 0.257)$$

[Show answer](#)

Submit

You have used 0 of 2 attempts

Now, imagine that you are considering a randomized experiment similar to the camera experiment. The exception is that you plan to give teachers lower incentives: half the monetary amount that was given in the previous experiment.

Question 12

0.0/1.0 point (graded)

Imagine that the relationship between incentives and the variable **open** is linear. What would be the expected ATE of this new intervention?

Please round your answer to the third decimal place, i.e. if it is 0.3414, please round to 0.341.

Answer: 0.098

Explanation

Since we are assuming a linear relationship then half of the incentives should have half of the effect. Then, our estimate would be $\frac{0.1969}{2} = 0.09845 \approx 0.098$

[Show answer](#)

Submit

You have used 0 of 2 attempts

Question 13

0.0/1.0 point (graded)

Assume that the value from Question 12 is the minimum ATE such that the intervention is cost-effective. What is the sample size required to have a power of at least 90% with the following properties?

- with a significance level of 5%
- an equal number of treated and control units
- σ^2 is the average of the variance of the control and the treatment group in the existing data

Hint: Recall that the formula for sample size is:

$$N = \frac{(\Phi^{-1}(1 - \beta) + \Phi^{-1}(1 - \frac{\alpha}{2}))^2}{\frac{\tau^2}{\sigma^2} \gamma (1 - \gamma)}$$

where β is the operating characteristic and $1 - \beta$ is the desired power.

Answer: 103

Explanation

We have that the formula for the sample size is given by:

$$N = \frac{(\Phi^{-1}(1-\beta) + \Phi^{-1}(1-\frac{\alpha}{2}))^2}{\frac{\tau^2}{\sigma^2} \gamma(1-\gamma)}$$

We are interested in $\alpha = 0.05$, $\beta = 0.1$. Since the number of treated and control units is the same, we will have that $\gamma = 0.5$. We will set $\tau = 0.09845$ as discussed above. Finally we will assume σ^2 is the average of s_c^2 and s_t^2 found above, $\hat{\sigma}^2 = \frac{49*0.126^2 + 51*0.1766^2}{100} = 0.0236$. You can get the values for the inverse normal distribution from various online calculators.

$$N_{\beta=0.1} = \frac{(1.282+1.96)^2}{\frac{0.09845^2}{0.0236} * 0.25} = 102.36$$

Since it is not possible to have 102.36 schools and this is the minimum sample size to be powered enough, the correct answer is 103 schools.

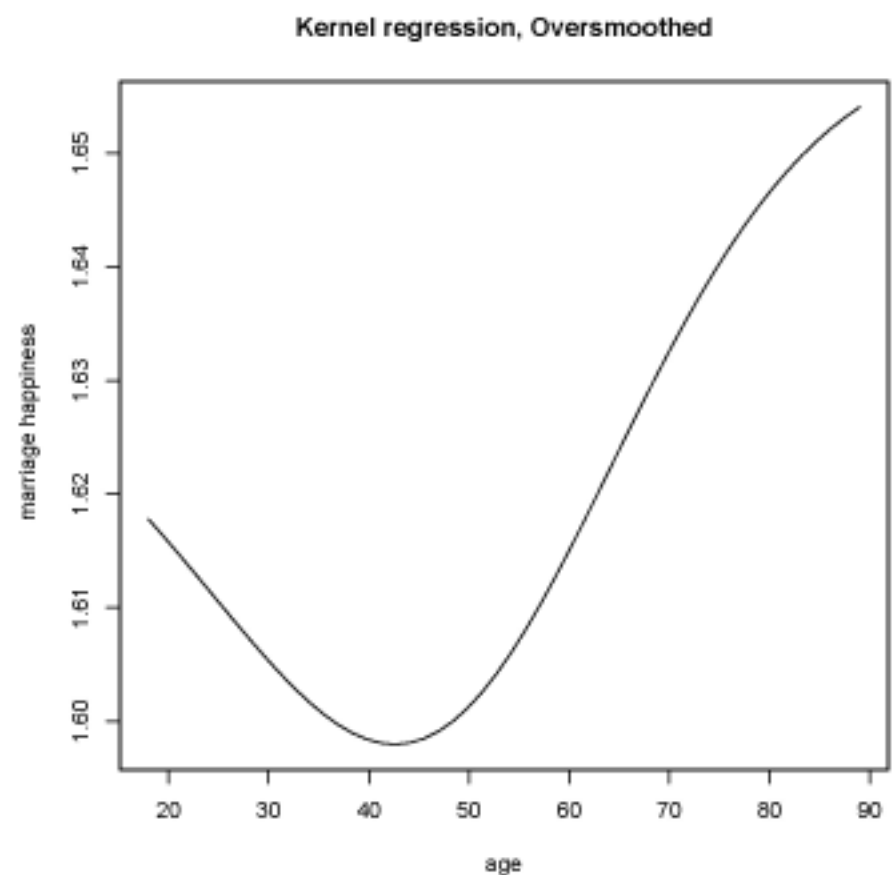
[Show answer](#)

Submit

You have used 0 of 2 attempts

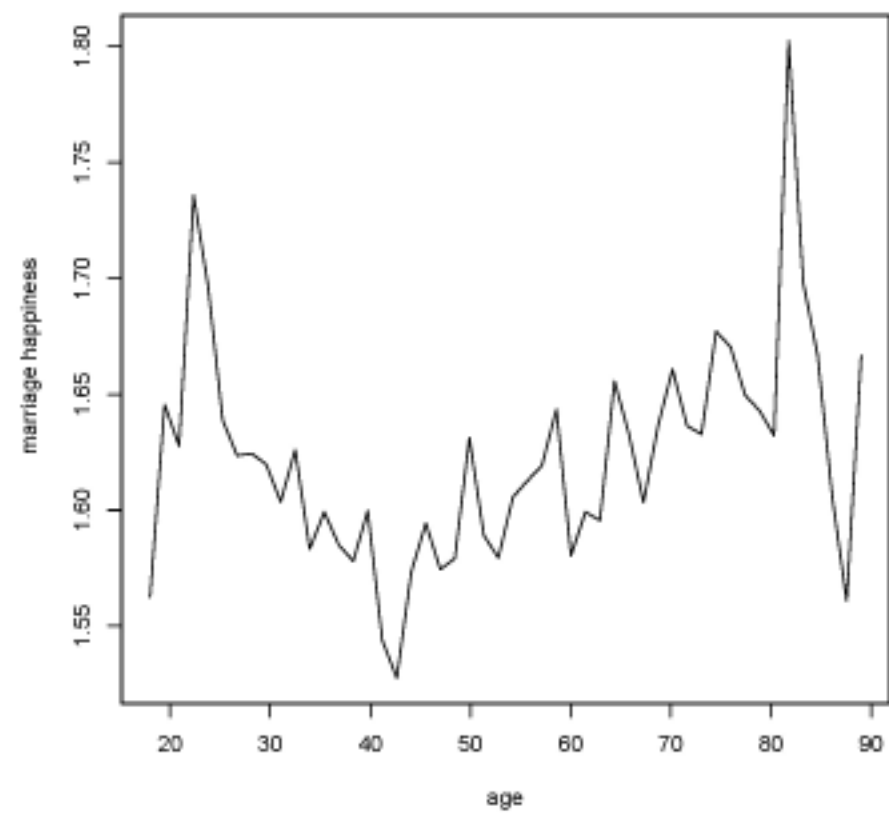
Now we are going to consider nonparametric regressions. The following plots show three different nonparametric regressions that relates the level of happiness in a marriage with age (where 2 corresponds to "very happy", 1 to "pretty happy", and 0 to "not too happy").

Plot A:

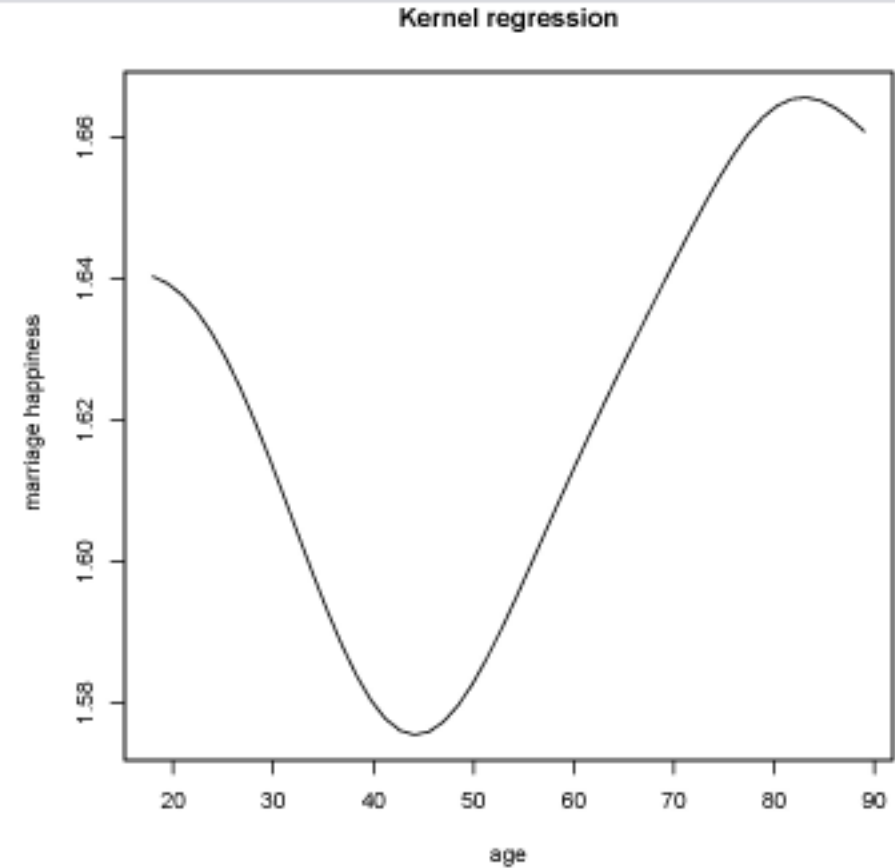


Plot B:

Kernel regression, Undersmoothed



Plot C:



Question 14

0.0/1.0 point (graded)

Rank the three plots from the one with the narrowest to the widest bandwidth.

☐ a, c, b

☐ b, a, c

☐ c, b, a

☒ b, c, a ✓

☐ c, a, b

☐ a, b, c

Explanation

The narrower the bandwidth the less smooth the plot will look like. The plot in (a) was done with a bandwidth of 15, the plot in (b) was done with a bandwidth of 0.1, and the plot in (c) was done with the optimal bandwidth in R that corresponds to a value of 6.36.

[Show answer](#)

Submit

You have used 0 of 2 attempts

i Answers are displayed within the problem

Going back to the data from `teachers_final.csv`, we are now going to focus on two variables:

- `pctpostwritten`, which denotes the mean student test scores after the intervention
- `open`

We want to see what the relationship between the fraction of days the school is open and student achievement. Use the code below (from lecture) to plot the kernel regression between these two variables using the R package `np`:

```
attach(schools)
plot <- npreg(xdat=XXX, ydat=XXX, bws=XXX, bandwidth.compute=FALSE)
plot(plot)
```

Question 15

1 point possible (graded)

Use your code to generate plots for the following bandwidths. Which of them seems most appropriate given the data?

☐ 20

☐ 0.001

☐ 1

☒ 0.04 ✓

Explanation

You can use the following code to generate these plots:

```
attach(schools)
```

```
bw_a <-npreg(xdat=pctpostwritten, ydat= open, bws=0.04,bandwidth.compute=FALSE) plot(bw_a)
```

```
bw_b <-npreg(xdat=pctpostwritten, ydat= open, bws=0.001,bandwidth.compute=FALSE) plot(bw_b)
```

```
bw_c <-npreg(xdat=pctpostwritten, ydat= open, bws=1,bandwidth.compute=FALSE) plot(bw_c)
```

```
bw_d <-npreg(xdat=pctpostwritten, ydat= open, bws=20,bandwidth.compute=FALSE) plot(bw_d)
```

Looking at the resulting plots, it's clear that the bandwidth of 0.001 is overfitting the data, whereas the plots using the bandwidths of 1 and 20 look like they are underfitting the data. Based on the plots, the bandwidth of 0.04 looks like the best options.

[Show answer](#)

Submit

You have used 0 of 2 attempts

i Answers are displayed within the problem

Question 16

0.0/1.0 point (graded)

Suppose we are interested in testing whether or not the **distribution** of the share of days a school is found to be open in the treatment group is statistically distinguishable from the **distribution** for that of the control group. Which of the following would be most useful for this purposes?

☐ Joint density plot

☐ Histogram of the variable by group

☒ Kolmogrov-Smirnov test ✓

☐ Kernel regression

☐ None of the above

Explanation

Recall, the Kolmogrov-Smirnov test (or KS test) is a test for the equality of distributions. So it is the ideal test given the question.

[Show answer](#)

Question 17

0.0/1.0 point (graded)

Let $i \in T, C$ index the cohort school i assigned. m_i denotes the sample mean of a variable (e.g. student scores) for group i , μ_i denotes the population mean of the variable, and F_i denotes the CDF for group i .

For each hypothesis test below, indicate which of the following methods is most useful for testing that hypothesis. **Enter N for using Neyman's method of inference, F for Fisher's exact test, and K for the KS test.**

A. $H_0 : \mu_T - \mu_C = 0$ vs. $H_1 : \mu_T - \mu_C \neq 0$

Answer: N

B. $H_0 : \mu_T - \mu_C > 0$ vs. $H_1 : \mu_T - \mu_C \leq 0$

Answer: N

C. $H_0 : m_T - m_C < 0$ vs. $H_1 : m_T - m_C \geq 0$

Answer: F

D. $H_0 : F_T = F_C$ vs. $H_1 : F_T \neq F_C$

Answer: K

E. $H_0 : F_T > G$ vs. $H_1 : F_T \leq G$ where $G \sim N(0, 1)$

Answer: K

Explanation

Neyman's Method of inference is used to make inferences about the underlying data based on the observed sample, since A and B both refer to population parameters, N is the correct answer. Fisher's exact test is useful for making inferences about the observed sample (C), whereas the KS test is used to compare any two distributions and is therefore appropriate for testing hypotheses D and E.

[Show answer](#)

Submit

You have used 0 of 2 attempts

i Answers are displayed within the problem

Question 18

0.0/1.0 point (graded)

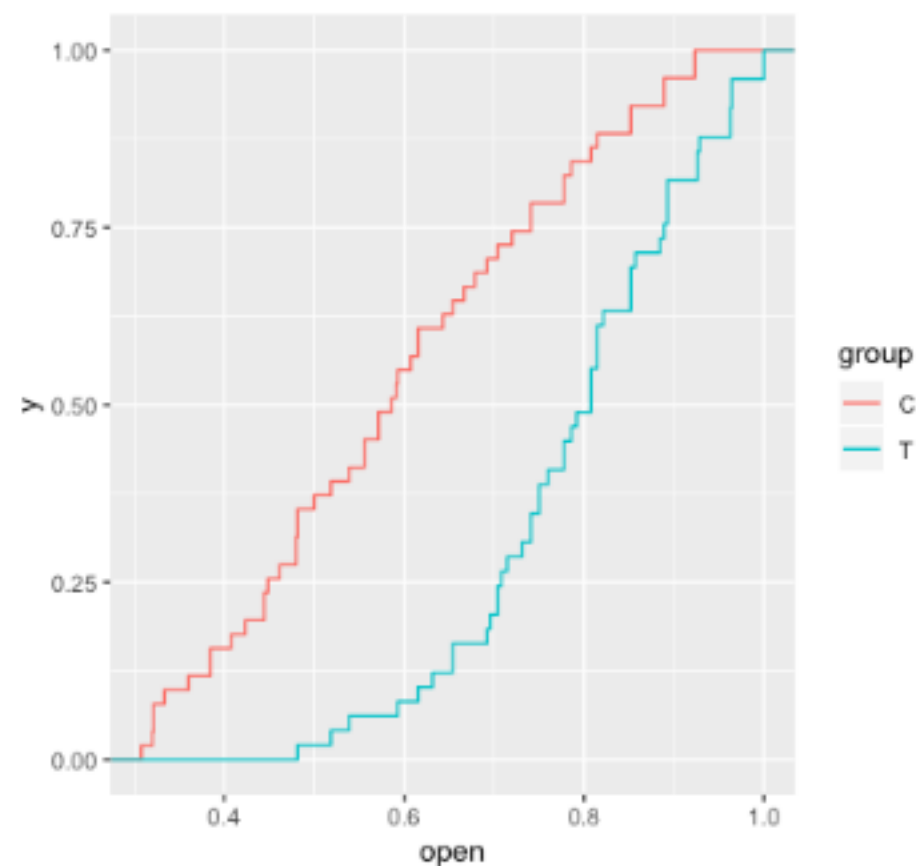
Use the R command `stat_ecdf()` to generate a plot of the CDFs for each cohort to see those results visually. Does the

distribution of `open` in the treatment group FOSD that of the control group?

☒ Yes ✓

☐ No

Explanation



As you can see from the plot, the distribution of `open` in the treatment group does FOSD that of the control group.

[Show answer](#)