

# NLP4CSS: Homework #3

Due 11:59pm EST on 03/08, 2024

*Instructor: Anjalie Field; Lead TA: Samuel Lefcourt; special thanks to Carlos Aguirre*

---

**Guidelines.** This assignment is to be completed **individually**. Be sure to comply with course policies on the course website.

**Starter Code.** Starter code is provided.

```
HW3
|- main.py
|- requirements.txt
```

**Submission.** This homework has written and coding components. For coding, you will complete the python file and submit it to gradescope. For the written part, you will write your answers in a PDF named `README.pdf` and also submit it to gradescope. Your PDF should contain answers to Problem 3 and Problem 5. Course Entry Code: YDPR48. Your final submission should have all the completed python file as well as your `README.pdf`.

## Introduction

In class you learned about multiple aspects of causal inference. Following is a summary of key concepts:

- **Causal Inference.** Process of establishing, and quantifying causal relationships empirically
- **Treatment and Outcome.** The variables that represents possible actions (treatment) and the observed result after treatment (outcome)
- **Confounders.** Factors that are common causes of both the treatment and the outcome

In this homework, you will estimate the treatment effect in a semi-synthetic dataset based on the 20 Newsgroup dataset. The 20 Newsgroup dataset is a collection of approximately 20K newsgroup documents corresponding to different topics; some of the topics are highly related and some are not, so we relabel the topics according to subject matter in the code.

The dataset is created by the following:

$$\mathbf{Y} = \text{const.} + U\_bias * \mathbf{U} + Z\_bias * \mathbf{Z} + \epsilon, \quad (1)$$

$$\mathbf{Z} = 1.0 * U + \epsilon > 0 \quad (2)$$

Where  $\mathbf{Z}$  is the treatment variable—for example, this could be the medium in which the documents were consumed: a mobile device or a desktop.  $\mathbf{Y}$  is the outcome variable; using the previous example, we could estimate the effect of the medium used on the outcome of the number of lines read. And  $\mathbf{U}$  is an unobserved confounder, in our example it could be the topic of the document, where we assume readers prefer to read certain topics on mobile vs desktops and the topic may also influence the number of lines in a document. Finally,  $\epsilon$  is random noise.

For our homework,  $\mathbf{U}$ , the unobserved confounder, is whether a document belongs to a specific topic (*religion* topic is the default). The data is constructed so that  $\mathbf{Z}$  is a binary variable that depends on the confounder, and  $\mathbf{Y}$  depends on both  $\mathbf{U}$  and  $\mathbf{Z}$ .

In the code, the default values for the bias values are: `const. = -0.5`, `U_bias = 5.0`, `Z_bias = 0.05`, and are hyperparameters to the `get_data()` function. To help your implementation, we will give you some of

the expected outputs with these parameters. We will grade submissions using different parameters. When debugging, feel free to change the values to obtain different datasets.

## Problem 1: Estimating the treatment effect by regressing $Y$ on $Z$ only.

(10 Points) **Complete the `regress_y_on_z` function.** First, investigate what would happen if you don't adjust for the confounder. Complete the function to estimate the effect of  $z$  on  $y$  using the following pseudo-code:

```
def regress_y_on_z(data, max_iter):
    # format data
    # create OLS model
    res = model.fit(method='pinv', maxiter=max_iter)
    print(res.summary(yname='Y', xname=['const', 'Z']))
```

You can expect the output to be similar to the following:

OLS Regression Results						
Dep. Variable:	Y	R-squared:	0.064			
Model:	OLS	Adj. R-squared:	0.064			
Method:	Least Squares	F-statistic:	673.1			
Date:	Thu, 29 Feb 2024	Prob (F-statistic):	1.35e-143			
Time:	15:27:51	Log-Likelihood:	-18833.			
No. Observations:	9816	AIC:	3.767e+04			
Df Residuals:	9814	BIC:	3.769e+04			
Df Model:	1					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-0.2841	0.025	-11.475	0.000	-0.333	-0.236
Z	0.8674	0.033	25.945	0.000	0.802	0.933
Omnibus:	3321.213	Durbin-Watson:	2.014			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	8233.646			
Skew:	1.941	Prob(JB):	0.00			
Kurtosis:	5.249	Cond. No.	2.75			

Where the *const.* and *Z\_bias* are wrongly estimated (see the `coef` column and the confidence intervals) due to confounding bias.

## Problem 2: Using Confounder to Estimate $Y$

(10 Points) **Complete the `regress_y_on_z_and_u` function.** Next, assume that you do observe the confounder. Based on the `regress_y_on_z` code from the previous section, estimate the treatment effect by regressing  $Y$  on  $Z$  and the confounder  $U$ . You should obtain values very close to the actual bias values, however, this is an unrealistic scenario as the confounder  $U$  is unobserved.

## Problem 3: Controlling for Confounders with Structured Text

In the previous section we had access to the confounder, but in real life scenarios this is unobserved. However, note that the confounder  $\mathbf{U}$  is based on whether a document belongs to a specific topic of documents. Could we capture the confounder from observed text?

**Part A (20 Points) Complete the `regress_y_on_z_and_topics` function.** Instead of regressing from word statistics from all of the vocabulary tokens in the documents (it would take a long time to converge), we will instead reduce the dimensionality of texts by employing topic models. As a starting point, use a matrix factorization (NMF from sklearn with TfidfVectorizer as inputs) topic model with 50 topics. Use the NMF model to extract topics from documents and use those lower-dimensional features as well as  $\mathbf{Z}$  to regress  $\mathbf{Y}$ . You should find that the  $Z\_bias$  estimate improves compared to Problem 1, but is less good than Problem 2.

**Part B (5 Points) Find a good control.** Play around with the model (e.g. try LDA instead on NMF), number of topics, other hyperparameters, or other methods (word embeddings?) to obtain a control that improves the estimate of the coefficient  $Z\_bias$ . Ideally, .05 would be within the confidence intervals, but you do not have to achieve this for full credit. In your written report, write a short paragraph describing what you tried and what the results were. Please submit your code so that the 50-topic NMF model described above runs by default (You can put your experimental code in a different function or use an optional parameter to control what model runs).

I tried out different number of topics for NMF and LDA, from 20 vs 50 vs 100. I found that consistently, it seemed like 50 topics was the golden number since it had a closer coefficient prediction than 20 topics and 100 topics respectively. I also attempted a word-embeddings method using pretrained word embeddings. I collected word embeddings over each word and then averaged their values throughout the entire document, the results of these embeddings would be used as a proxy of the confounder. I note that an increase of dimensions in glove-twitter seemed to yield better performance (glove-twitter-200 was quite close, and had the true value within its CIs), and glove-twitter-200 may be better than a comparable embedding like word2vec-google-news-300. Lastly, I thought of a weighted average word-embedding, where I use TF-IDF values of words and try to use those values as weights, but I didn't think it would computationally run in time for me.

Model	Dimensions/topics	Z	Lower CI	Upper CI
NMF	20	0.4657	0.418	0.513
NMF	50	0.4737	0.426	0.522
NMF	100	0.4662	0.418	0.514
LDA	20	0.5448	0.494	0.596
LDA	50	0.5118	0.462	0.561
LDA	100	0.5547	0.503	0.607
word2vec-google-news-300	300	0.4415	0.394	0.489
glove-twitter-200	200	0.4861	0.437	0.536
glove-twitter-50	50	0.5639	0.510	0.618

Table 1: Results of Regressions

## Problem 4: Controlling for Confounders with Inverse Probability Weighting

In class you learned about inverse probability weighting, which assigns different weights to subjects based on their propensity scores to account for confounders

**(20 Points) Complete the `reweigh_with_propensity_scores` function.** As in Problem 3, we will assume we do not directly know the confounder, and we will use the text as a proxy. Obtain propensity scores by training a logistic regression model using TF-IDF-weighted features (we suggest sklearn with `max_iter=2000` for the logistic regression model). To avoid overfitting, split the data into 2 equal-sized halves: *train* (used to train the propensity score model) and *test* (you can use `train_test_split` from sklearn). Use the propensity scores to reweigh **Y**, and return the both the unadjusted and the adjusted average treatment effect estimated over the test set. (Hint: Adjusted=0.580, Unadjusted=0.866)

## Problem 5: Draw Causal Graphs

For the following brief scenarios where text may be employed to help examine causal relations, identify the treatment and the outcome(s) and potential confounder(s) and briefly justify your answer. Speculate what may be the possible confounders if none are evident from the description. Then, draw the causal graph. Your answer should be in the form:

Treatment: *your answer*

Outcome(s): *your answer*

Confounder(s): *your answer*

*Your graph, using any software or hand draw-and-scan you prefer as long as we can read it*

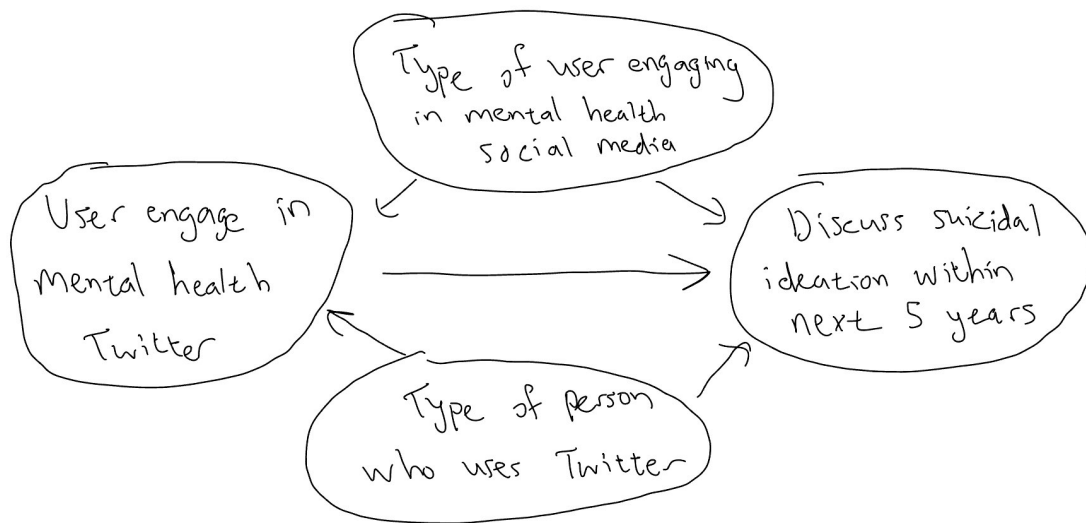
1. From the population of users tweeting about mental health, what linguistic structures or linguistic patterns differentiates those who proceed to discuss suicidal ideation in the future, from those who do not?

Treatment: Does user engage in mental health Twitter

Outcome(s): Does user discuss suicidal ideation within the next 5 years

Confounder(s): Users who engage in mental health social media, users of twitter

Users who engage in mental health social media already have a certain proclivity towards suicidal ideation since they decide to spend time engaged in this area (compared to people who don't). Users who engage in Twitter (and social media in general) have possibly different sharing behaviors, like they're more willing to share their internal thoughts so this affects the chance of them discussing/sharing suicidal ideation.



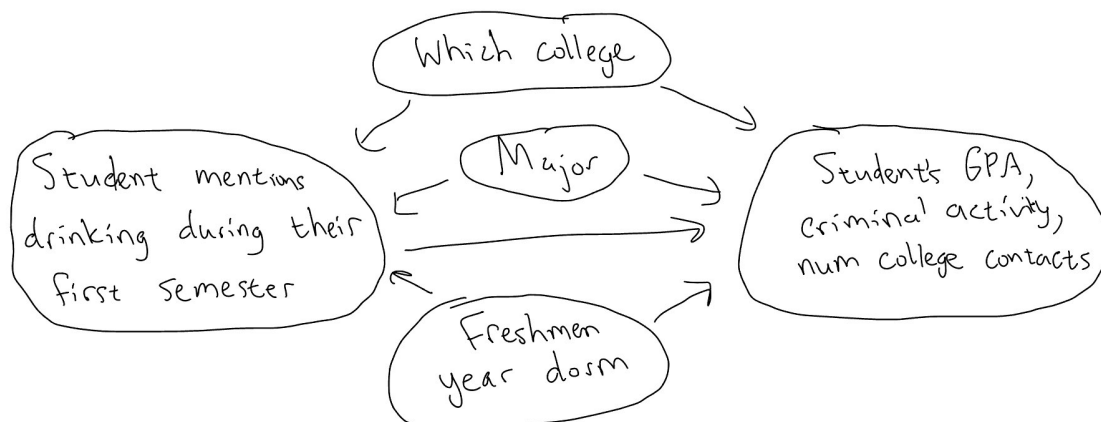
2. By collecting Reddit timelines from students entering college, we study what effect does drinking early in college have on college success, including habits, social relationships, and even criminal activity, of those who mention drinking during their first semester versus those that do not.

Treatment: Student mentions drinking during their first semester

Outcome(s): Student's GPA, criminal activity, and purported number of contacts from college

Confounder(s): Which college they attend, major, freshman year dorm living location

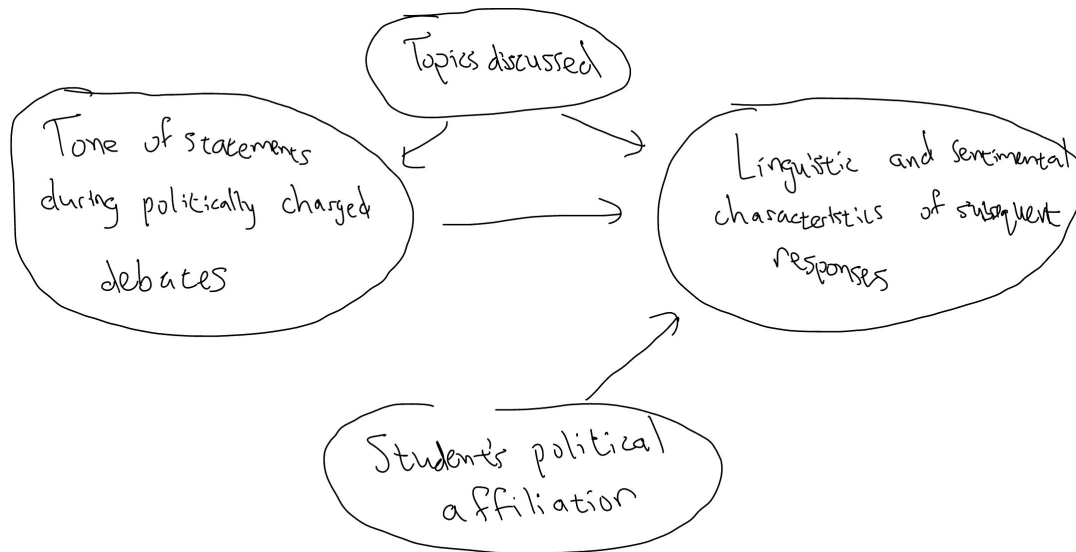
Each college has different general drinking behaviors trends and culture on campus and has a significant impact on student outcomes. Different majors have different inclinations towards students drinking and their outcomes. Where a student is housed for their freshman dorm has a huge impact on their drinking behavior while also having an impact on student outcomes due to numerous reasons like contacts/closeness/culture of that dorm.



3. We collect politically charged debates and investigate how the tones of the statements made during the debates changes the linguistic and sentiment characteristics in subsequent responses.

Treatment: The tone of statements made during politically charged debates  
Outcome(s): Linguistic and sentimental characteristics of subsequent responses  
Confounder(s): The topic discussed  
Covariate: Student's political affiliation

The topic discussed has a huge impact on the tone of the statement since some topics are more sensitive than others, even in a politically charged debate, in both the statement and the responses (outcome). And a student's political affiliation has an impact on the outcome since depending on if they agree or disagree with the statement they'll have different linguistic and sentimental characteristics. This is a big bonus due to the trend of students and younger people having more liberal beliefs.



4. If an AI article published under a woman's name were instead published in the same venue under the name of a man with the same scholarly credentials, would it be cited more?

Treatment: Gender of author's name associated with AI article  
Outcome(s): Number of citations  
Confounder(s): Niche in AI, institutional affiliation, co-authorship

Depending on which niche in AI the author talks about, this could influence both the gender of the author of who'd be interested in it, as well as the number of citations. Institutional affiliation, some institutions may be more male-dominated than others, and will also have their own relative brand of prestige or reach. Lastly, your co-authors has an impact on your gender (if your co-authors are mostly females/males) and your co-authors has a big impact on the number of citations you get.

