

CS 6603: AI, Ethics, and Society

Stats 101 Assignment

Cleo Zhang

yzhang3761@gatech.edu

Abstract—This assignment will take the Mental Health in Tech Survey dataset as an example and start exploring relationships in data with basic statistical measurements, including identifying dependent and independent variables and creating graphs to support the “fairness” hypothesis and “bias” hypothesis. Furthermore, this report will run the random sampling method using 50% of the original data.

1 OVERVIEW OF MENTAL-HEALTH-IN-TECH-SURVEY-2019.CSV

Dataset: Mental Health in Tech Survey

Number of Observations: 353

Number of Variables: 82

Regulated Domain in Law: HealthCare

Number of Protected Class Variables: 11

Table 1 — Variables associated with a protected class and the associated legal precedence

	Converted Variable Name	Protected Class	Law
Do you currently have a mental health disorder?	<i>have_mental_health_disorder</i>	Disability status	Rehabilitation Act of 1973; Americans with Disabilities Act
Have you ever been diagnosed with a mental health disorder?	<i>Diagnosed_with_mental_health</i>	Disability status	Rehabilitation Act of 1973; Americans with Disabilities Act
Have you had a mental health disorder in the past?	<i>had_mental_health_disorder</i>	Disability status	Rehabilitation Act of 1973; Americans with Disabilities Act

	Converted Variable Name	Protected Class	Law
What is your age?	<i>age</i>	Age	Age Discrimination in Employment Act of 1967
What is your gender?	<i>gender</i>	Sex	Equal Pay Act of 1963; Civil Rights Act of 1964, 1991
What is your race?	<i>race</i>	Race	Civil Rights Act of 1964, 1991

2 RELATIONSHIPS BETWEEN DEPENDENT AND INDEPENDENT VARIABLES

This section will compute the selected dependent variables (*treatment*, *mental_health_disclosure*, *mental_health_support*) as functions of the independent variables identified in *Table 1*.

- Frequency tables and histograms for *have_mental_health_disorder* VS. *treatment*, *mental_health_disclosure*, *mental_health_support*

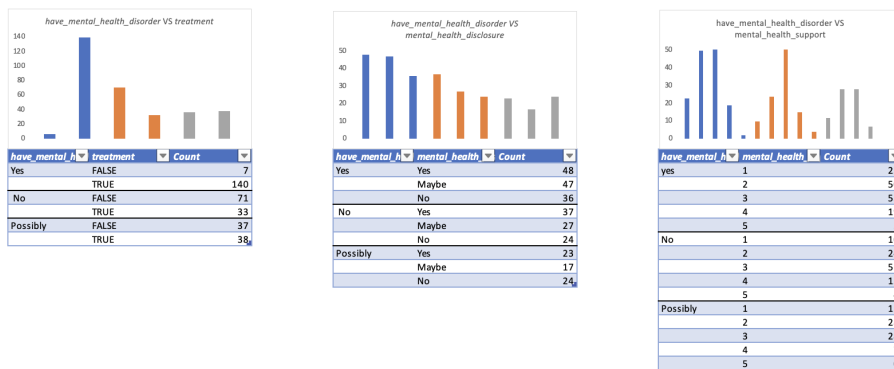


Figure 2.1 — *have_mental_health_disorder* VS. *treatment*, *mental_health_disclosure*, *mental_health_support*

- Frequency tables and histograms for *diagnosed_with_mental_health* VS. *treatment*, *mental_health_disclosure*, *mental_health_support*



Figure 2.2 — *diagnosed_with_mental_health* VS. *treatment*, *mental_health_disclosure*, *mental_health_support*

- Frequency tables and histograms for *had_mental_health_disorder* VS. *treatment*, *mental_health_disclosure*, *mental_health_support*

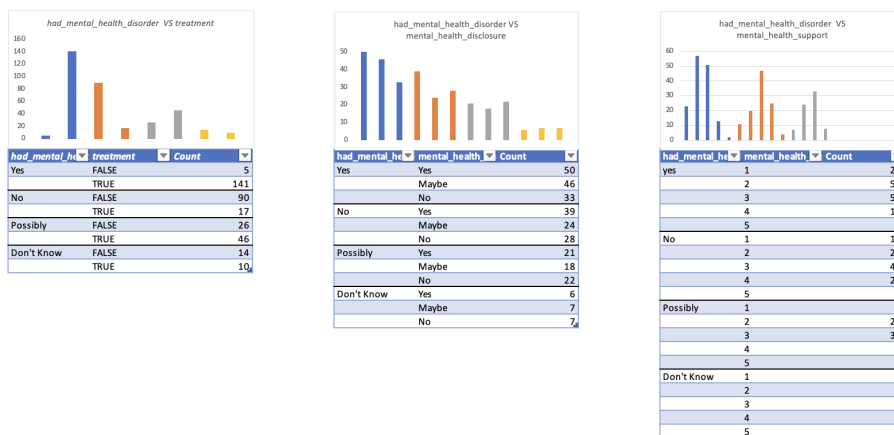


Figure 2.3 — *had_mental_health_disorder* VS. *treatment*, *mental_health_disclosure*, *mental_health_support*

- Frequency tables and histograms for *age VS. treatment, mental_health_disclosure, mental_health_support*

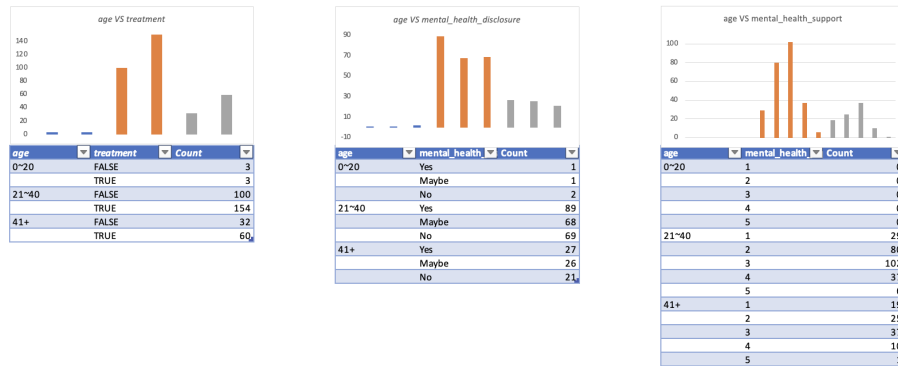


Figure 1 — age VS. treatment, mental_health_disclosure, mental_health_support

- Frequency tables and histograms for *gender VS. treatment, mental_health_disclosure, mental_health_support*

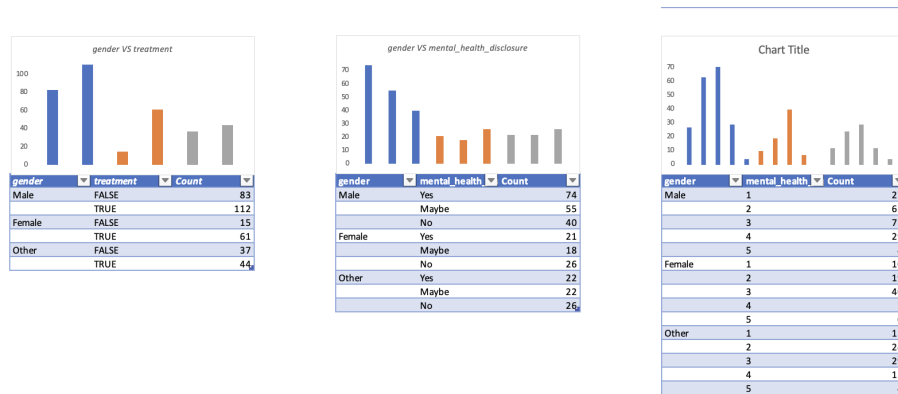


Figure 2 — gender VS. treatment, mental_health_disclosure, mental_health_support

- Frequency tables and histograms for *race* VS. *treatment*, *mental_health_disclosure*, *mental_health_support*

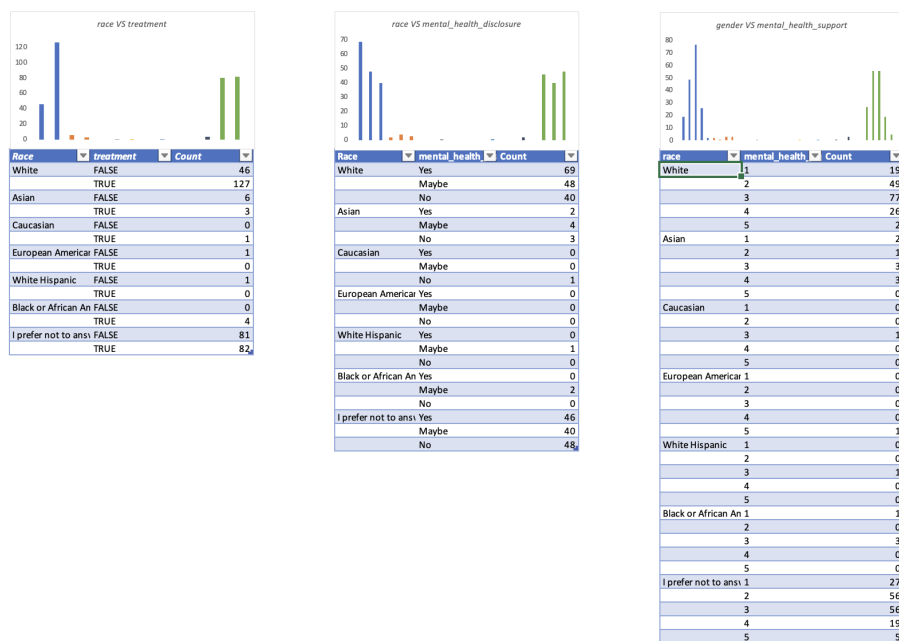


Figure 3 — *race* VS. *treatment*, *mental_health_disclosure*, *mental_health_support*

3 DATA MANIPULATION

This section will explain the manipulation of *gender* and *mental_health_disclosure* combination and create graphics to support Fair Hypothesis and Bias Hypothesis.

3.1 Create a graph to support the "fairness" hypothesis.

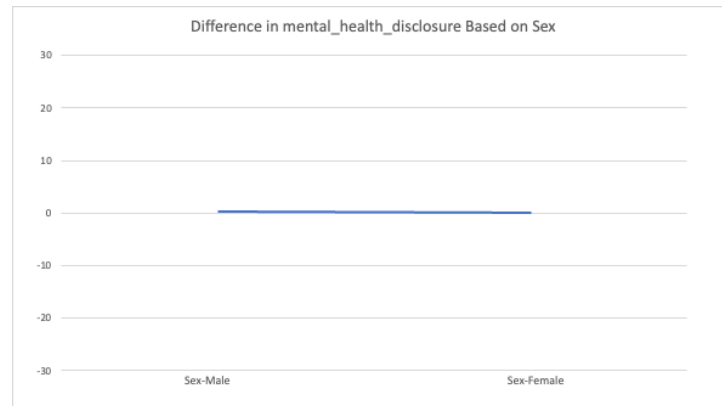


Figure 3.1 – Graphic to support “fairness” hypothesis

As seen from Figure 3.1, *mental_health_disclosure* decisions (Yes/No) are not dependent on the gender variable. [Manipulations: Used line graph; Increased the scale to +-30; Mapped the ratio of Yes decisions (74/304 versus 21/304)].

3.2 Create a graph to support the “bias” hypothesis.

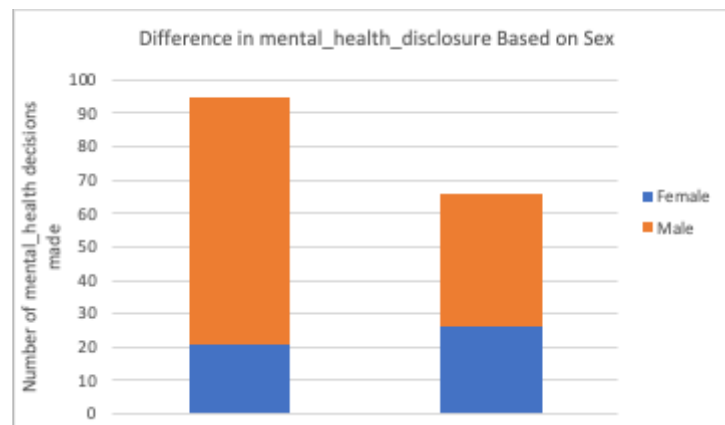


Figure 3.2 – Graphic to support “bias” hypothesis

As seen from Figure 3.2, *mental_health_disclosure* decisions (Yes/No) depend on the gender variable. [Manipulations: Used stacked bar graph; Used raw data from the frequency table in Section 2; Reworded labelled].

4 CALCULATE THE AVERAGE (USING MEAN, MEDIAN, AND MODE) OF THE PROTECTED CLASS GROUP.

Assign the answer Yes, Maybe, No to 1, 0, -1 for the *mental_health_disclosure* question; I can summarize my computation in Table 2.

Table 2 - Variables associated with a protected class and the associated legal precedence

Protected Class Variable (Sex)	Mean	Median	Mode
Original Data Set	o(Maybe)	o(Maybe)	1(Yes)
Reduced Data Set	o(Maybe)	o(Maybe)	1(Yes)
Difference	No Difference	No Difference	No Difference

5 EXPLAIN THE DIFFERENCES BETWEEN THE FREQUENCY TABLE AND HISTOGRAMS BETWEEN THE ORIGINAL AND REDUCED DATASET

Two of the most significant differences shown in *Figure 5* are:

- gender-male: the decisions of *Maybe* and *No* became the same, while the *Maybe* decision is much higher than *No* in the original dataset.
- gender-other: the *Maybe* decision (instead of *No* as in the original dataset) became the most popular one.

I used the Random Sampling Method to generate the reduced dataset, and I believe the members associated with the protected class variable are harmed by it. The reasons are:

- Even though the mean, median and mode values seemed the same, the data characteristics were not fully reflected in the reduced dataset, for example, standard deviation.
- The portion of each decision made by each gender became very different compared to the original data, which could be a result of small data samples.

gender	mental_health_disclosur	Count
Male	Yes	35
	Maybe	25
	No	25
Female	Yes	13
	Maybe	7
	No	12
Other	Yes	10
	Maybe	13
	No	12

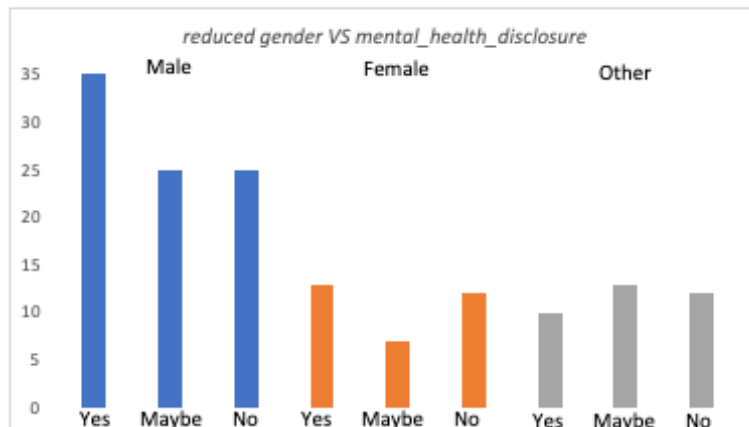


Figure 5 – Frequency table and histogram from the reduced dataset