

Quantitative Methods Assignment 1

Firstname Lastname

T-test Recap

One Sample T-test

In a one-sample test, we want to test whether the mean of a population is equal to a pre-specified value μ_0 based on a sample collected from the population. The hypothesis testing problem can be formulated as:

$$H_0 : \mu = \mu_0$$

$$H_a : \mu \neq \mu_0$$

Given sample mean \bar{x} and the standard error of the sample mean, $SE(\bar{x})$, we can construct a t-test:

$$T = \frac{\bar{x} - \mu_0}{SE(\bar{x})}$$

The null distribution of T is a t distribution with $n - 1$ degrees of freedom.

If we set a significance level of .05 (for a 2-tailed test) then: If $\text{abs}(T)$ is greater than the 97.5th quantile of a t distribution ($\text{df}=n-1$), we reject the null hypothesis in favor of the alternative that the population mean is different from μ_0 .

Note that for a sample size of n : $SE(\bar{x}) = \frac{s}{\sqrt{n}}$ where s is the sample standard deviation.

Two (Independent) Sample T-test

In a two sample t-test problem, we want to test whether the mean of population 1 is the same as the mean of population 2. The hypothesis testing problem can be formulated as:

$$H_0 : \mu_1 = \mu_2$$

$$H_a : \mu_1 \neq \mu_2$$

We can construct a t-test to test H_0 versus H_a ,

$$T = \frac{(\bar{x}_1 - \bar{x}_2)}{SE(\bar{x}_1 - \bar{x}_2)}$$

We can estimate the standard error of $\bar{x}_1 - \bar{x}_2$ using the following formula (derived in the math review packet) :

$$SE(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

In this formula, s_1 and s_2 represent the sample standard deviations for group 1 and group 2, respectively. If we believe that the two populations have unequal variances, then we can estimate s_1 and s_2 separately. If we believe that the two populations have equal variances, then we can say $s_1 = s_2 = s_{pooled}$ where:

$$s_{pooled} = \sqrt{\frac{s_1^2(n_1 - 1) + s_2^2(n_2 - 1)}{n_1 + n_2 - 2}}$$

Then, the null distribution of T is a t distribution with degrees of freedom as follows:

In general, we use the smaller of $n_1 - 1$ and $n_2 - 1$ as the degrees of freedom for the t distribution.

If assuming equal variances and using s_{pooled} , we can use $n_1 + n_2 - 2$ degrees of freedom.

Linear Regression with a Dummy Variable Predictor

First, let's read in some data. In order to read in the file, you will need to download height_sex.csv to your computer and include the full file path in the function call below. For example, I could put the file in my "Downloads" folder and call:

```
height <- read.csv("/Users/sophiesommer/Downloads/height_sex.csv")
```

```
#Read in data
```

```
height <- read.csv("height_sex.csv")
```

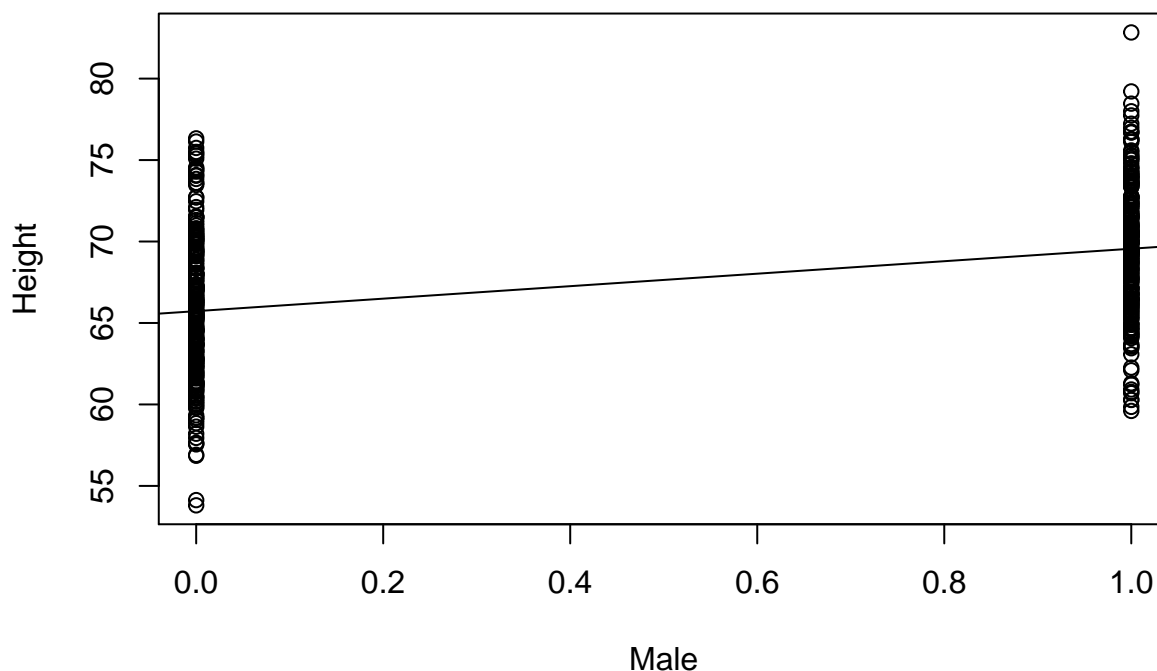
```
#Inspect the first few rows
```

```
height[1:5,]
```

```
##   X id Gender Birth.Order   Height   Weight Male
## 1 1 1      M           1 65.44441 166.14611     1
## 2 2 2      F           2 65.44836  76.74407     0
## 3 3 3      M           2 68.39149 140.14253     1
## 4 4 4      M           1 73.92210 165.30409     1
## 5 5 5      F           1 70.98392 117.61932     0
```

For this assignment, we are going to focus on the "Height" and "Male" variables. "Height" is given in inches and "Male" is an indicator variable (equal to 1 if the person is male and equal to 0 if the person is female). If we plot height against the gender indicator variable, we get:

```
plot(height$Male, height$Height, xlab="Male", ylab="Height")
abline(lm(height$Height~height$Male))
```



Suppose for a moment that you only have one set of heights, say, $\{60, 63, 63, 70, 72\}$ and you want to find some value, M , such that the sum of the squared differences between M and each value in your set of heights is minimized. It turns out that the value of M which satisfies this minimization problem is the mean of the dataset (you can prove this with some relatively simple calculus). Note: minimizing squared differences is akin to minimizing absolute distances (the mean is the solution to both minimization problems); however, we generally choose to minimize squared differences instead of absolute differences to make the calculus easier. The proof is left to the interested reader, but hopefully this claim already feels somewhat intuitive.

Now, suppose we try to fit an OLS regression line to this data. This would involve finding values of β_0

and β_1 such that the (squared) vertical distance between each point (in the above graph) and the line $Height = \beta_0 + \beta_1 * Male$ is minimized.

There are only two possible values of “Male” in this dataset: 0 and 1. Therefore, one way to think about minimizing the vertical distances between each data point and the line is this: Find the mean height of women (i.e., everyone with Male=0), find the mean height of men (i.e., everyone with Male=1), and connect the two points (0, mean womens height) and (1, mean mens height).

Now we can try to write the equation for this line in the form $Height = \beta_0 + \beta_1 * Male$. Using the strategy described above, we know that (0, mean womens height) is the y-intercept of the line. Therefore $\beta_0 = \text{mean womens height}$. We also know that β_1 is the slope of the line. Since we know two points on the line, we can calculate the slope as $\frac{y_2 - y_1}{x_2 - x_1} = \frac{(\text{mean mens height} - \text{mean womens height})}{(1 - 0)} = \text{mean mens height} - \text{mean womens height}$. So, $\beta_1 = \text{mean mens height} - \text{mean womens height}$.

A question to think about: is testing $\beta_1 = 0$ the same as testing the mean of men’s height and women’s height to be the same (this is the two sample t-test problem)?

Assignment

1. The mean height of men is calculated below and saved as MM. Fill in the rest of the code to a) save the mean height of women as MW and b) print the mean height of women

```
# Save mean height of men as MM
MM <- mean(height$Height[height$Male==1])

# Print mean height of men
MM
```

```
## [1] 69.55706
```

```
# Save mean height of women as MW

# Print mean height of women
```

2. Calculate the mean difference for the two groups and plot the regression line

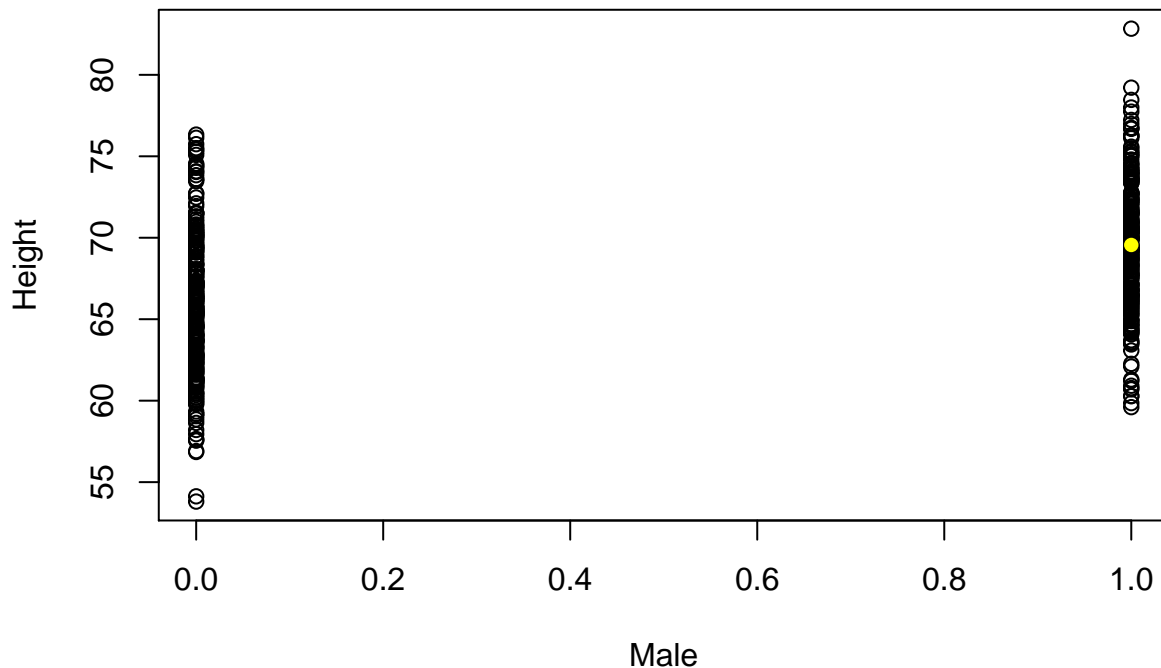
- a) Save and print: MeanDiff = MM - MW
- b) Add (0, MW) and (1,MM) to the plot from above (in yellow)
- c) Draw a line connecting (0, MW) and (1,MM) using the abline() function

```
# Save the difference in means (MM-MW) as the variable MeanDiff:

# Print MeanDiff

# Re-plot the data
plot(height$Male, height$Height, xlab="Male", ylab="Height")

# Add (0, MW) to the plot (in yellow) (note: (1,MM) has been added for you)
# Note: pch=16 makes the dot filled in, which makes it easier to see
points(1, MM, col=7, pch=16)
```



```
# Draw the line Height=b_0 + b_1*Male where b_0=MW and b_1=MeanDiff
# Use the abline() function which takes parameters a and b where a=y-intercept and b=slope
# Syntax is: abline(a,b)
```

3. Estimate the standard error for an independent sample T-test by hand (note: assume pooled standard deviation)

a) Calculate n_M and n_W , the number of men and women in the sample, respectively

b) Calculate s_M and s_W , the sample standard deviations for the men's heights and women's heights, respectively

c) Calculate s_{pooled} for the two groups using $s_{pooled} = \sqrt{\frac{s_1^2(n_1-1) + s_2^2(n_2-1)}{n_1+n_2-2}}$ and save the result as `spool`

d) Use the following equation: $SE(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ to estimate $SE(\bar{x}_1 - \bar{x}_2)$ where $s_1 = s_2 = s_{pooled}$ and $n_1 = n_W$, $n_2 = n_M$. Save the result as `SEMeanDiff` and print the result.

```
# Add code to calculate the number of women and save the result as nW
nM <- length(height$Height[height$Male==1])

# Add code to calculate sW, the sample standard deviation for women's heights
sM <- sd(height$Height[height$Male==1])

# Calculate pooled standard deviation and save as spool (note: sqrt() calculates square root)

# Print spool

# Estimate the standard error of the difference in means and save it as SEmeanDiff

# Print SEmeanDiff
```

4. Calculate a T-Statistic by hand to compare the null hypothesis that $MM - MW = 0$ versus the alternative hypothesis that $MM - MW \neq 0$

- a) Calculate $T = \frac{(\text{MeanDiff}) - 0}{\text{SEMeanDiff}}$ and save the result as Ttest
- b) Use the pt function to calculate a p-value for Ttest (degrees freedom= nW+nM-2). Hint: You will want to use the parameter lower.tail=FALSE because Ttest should be > 0 for this example, then multiply by 2 to get a 2-sided test

```
# Calculate a test statistic and save it as Ttest
```

```
# Print Ttest
```

```
# Use pt to calculate a p-value
```

5. Do the same thing using the t.test() function in R. Compare the value of t and p-value computed by the t.test() function to the values you calculated above. Note that setting var.equal=TRUE causes t.test() to use pooled variance. Do you get the same value (to two decimal places) for the T statistic?

```
# Use the t.test function to compare the men's heights and female's heights
```

```
# Note: the code has already been written for you (just uncomment the line below to run it)
```

```
# t.test(height$Height[height$Male==1],height$Height[height$Male==0], var.equal = TRUE)
```

6. Use the lm() function to implement ordinary least squares regression to estimate β_0 and β_1 in the equation $\text{Height} = \beta_0 + \beta_1 * \text{Male}$.

```
# Use lm() to estimate b_0 and b_1 in the equation Height=b_0+b_1*Male; save the result as linmod
```

```
# Note: this has been done for you
```

```
linmod <- lm(Height~Male, data=height)
```

```
#Inspect a summary of the results using the summary() function on linmod
```

7. Reprint the following values that you saved in previous problems: Mean women's height (MW), difference in mean height for men and women (MeanDiff), standard error of this mean difference (SEMeanDiff), and the T statistic for the difference in means (Ttest). Which of these values can you find in the lm() output above?

```
# Re-print MW
```

```
# Re-print MeanDiff
```

```
# Re-print SEMeandiff
```

```
# Re-print Ttest
```

8. Note that the T-test for the intercept in the lm() output above tests the null hypothesis: $\beta_0 = 0$ against the alternative hypothesis: $\beta_0 \neq 0$. In the space below, please show how you could calculate the Std. Error (.2576) and the t value (255.14) of the intercept without using lm() or t.test(). Hint: lm() uses s_{pooled} to estimate the standard deviation of heights for both groups.

```
# Estimate standard error of b_0
```

```
# (hint: remember that b_0=MW and you can estimate the SE of the mean using the central limit theorem)
```

```
# Use the standard error of b_0 to calculate the t value
```