



# Inferential Statistics Report

## 1. Introduction

As we see in daily life, there are a lot of differences between men and women. Some are identifiable, while others are more based on personalities. I have always seen tall men back in my home country, but coming to a diverse international space. I saw a lot of tall women as well. And mostly, I encountered women being taller than men. And I got confused if there is actually a difference between the height of men and women or is it just a perception, Hence the research question for my assignment - “Is there a difference between men’s and women’s height?”. To answer this question, we will use a difference of means test and a confidence interval will be created.

## 2. Dataset

The dataset used for this assignment was featured in Openintro Statistics Textbook (CDC,2019). It is a sample of a youth risk behavior surveillance system. The data set has 100 respondents and various questions were asked related to their physical strength, sleep, height, weight, gender, and driving habits.

However, for this paper, we will focus only on two variables namely Gender, which is considered a qualitative nominal variable. It does not have any order and is just based on categories - male or a female. The other variable is heights which is considered quantitative continuous as the heights are numbers with decimals and have an indefinite number of values between two numbers.<sup>1</sup>

## 3. Analysis

### Hypothesis

To answer our main question about whether there is a difference between the height of men and women, we will perform a difference of means significance test. We will consider 2 groups, men and women, which can be seen in appendix A. And by default, define our significance level = 0.05, which shows the probability of rejecting the null hypothesis when it is true and committing a type 1 error. So here, a significance level of 0.05 indicates a 5% risk of concluding that a difference exists when there is no actual difference.

Before we proceed, let us properly define our hypothesis.

---

<sup>1</sup> #variables: I distinguished between the two different type of variable which were important for my research question.

*Null hypothesis (H0)* : There is no difference of means between the height of female and males. ( $\bar{x}_2 - \bar{x}_1 = 0$ )

*Alternate hypothesis (HA)* : There is a difference of means between the heights of females and males. ( $\bar{x}_2 \neq \bar{x}_1$ )

Note: It is a two-tailed test because the difference of means can be either way. It can be that men are taller than women or vice versa.

## Summary Statistics

The summary statistics have been found using pandas in python. Before we go on further. We need to have descriptive statistics. Appendix A shows the relevant coding on the calculation of various statistics, and table 1 below indicates only relevant statistics for our significance test.

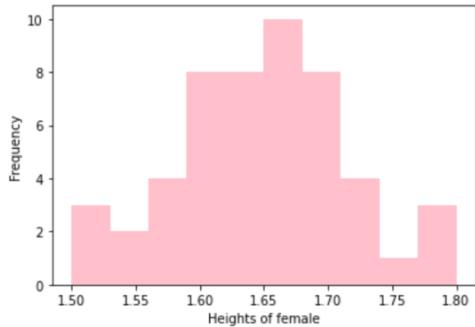
The histograms in Figures 1 and 2, which can be referred to in Appendix B, give us precise details about our statistics. We can see a clear difference between the mean values of heights of men and women. Further considering the median and mode are also quite different. However, the standard deviation and range are quite close, showing that the data does not vary much. Figure 3 shows a clear comparison graphically. And looking at the graphs, they do not look much skewed, except one outlier in male heights; other than that, they are more oriented towards a normal distribution.<sup>2</sup>

**Table 1: Summary statistics for the heights of the two sample groups: female and male.**

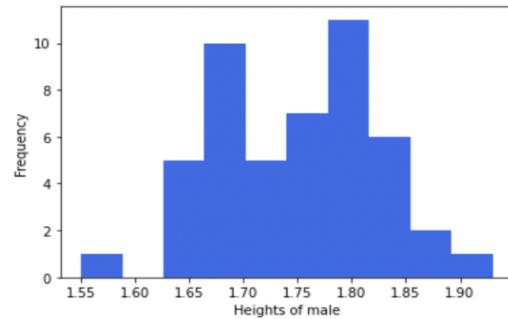
	<b>Female</b>	<b>Male</b>
<b>Count</b>	$n_1 = 51$	$n_2 = 48$
<b>Mean</b>	$\bar{x}_1 = 1.64$	$\bar{x}_2 = 1.74$
<b>Median</b>	1.65	1.75
<b>Mode</b>	1.60	1.78
<b>Standard Deviation</b>	$s_1 = 0.07$	$s_2 = 0.07$

<sup>2</sup> #descriptivestats: I calculated mean,median,mode, range to better interpret my data and used them later to conduct significance tests.

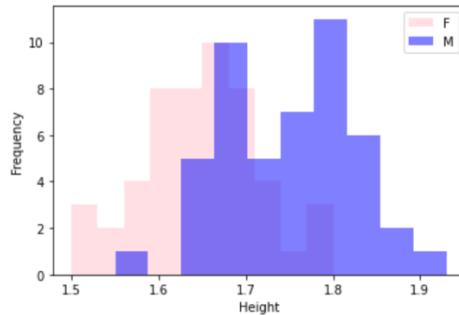
<b>Range</b>	0.30	0.37
--------------	------	------



**Fig 1.** Histogram showing the heights of female



**Fig 2.** Histogram showing the heights of male



**Fig 3.** Histogram showing comparison between heights of female and males<sup>3</sup>

### Conditions for inference:

Even though we have a sample size greater than 30 ( $n_1, n_2 > 30$ ) but we still don't know the population standard deviation and we are considering the standard error using sample standard deviation hence we will be using a ***t-distribution*** as our sampling distribution.

For this, we will consider 3 conditions:

- **Independent:**

Following the 10% rule and considering the whole population of more than a billion, we can clearly see that the sample data is less than 10% which makes our data to be independent.

---

<sup>3</sup> #dataviz: I generated a histogram to examine the distributions, and look at the skewness.

- **Normal Distribution:**

As the distribution has a sample size of greater than 30 ( $n>30$ ) and the values for mean, median and mode are quite close, we can conclude that the distribution is similar to a normal distribution.

- **Random Sample:**

As it is nowhere mentioned, if the dataset is random or not, we will consider this as a limitation and just assume the sample to be random.

As the conditions for inference are met, we will use the t distribution and apply central limit theorem.<sup>4</sup>

## Difference of Means Test

Now, we will conduct a difference of mean significance test and for this, first, we need to calculate the p-value. And before calculating the p-value, we need to find the

T-score which can be found using the formula  $T = (\bar{x}_2 - \bar{x}_1) \div \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ , where the

denominator is the standard error and the degrees of freedom is considered the minimum of one of the sample size subtracting 1.  $df = 47$ . See Appendix C for calculations in python.

We get a T score = 6.71 and the p-value as  $2.25 \times 10^{-8} < 0.05$ , and the value here is too small, which proves that there is convincing evidence to reject the null hypothesis and favor the alternate hypothesis.<sup>5</sup>

Let us also consider the practical significance, for which we will measure the hedges g because it helps us get more accurate results and eliminate the upward bias that we get using cohen's D. To find the value for hedges g, we first need to calculate pooled

standard deviation  $sp = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}}$ . Check Appendix C for calculations. We

get hedges g = 1.34 which means that there is a large effect size on the data means there is a big difference.<sup>6</sup>

## Confidence Intervals

---

<sup>4</sup> #distributions: I identified and justified a kind of distribution of my sampling distributions by checking all conditions carefully.

<sup>5</sup> #probability: I calculated the p value which states the conditional probability- that is it is the probability of getting the data observed or more extreme data if the null hypothesis is true.

<sup>6</sup> #significance: I calculated whether my dataset was statistically significant or not and then interpreted the results. Additionally, I also calculated the effect size for looking at the practical significance.

Now we will calculate the confidence interval for the mean of heights of men and women. We choose a 95% confidence interval as it will give us plausible range values. As we have discussed all the conditions for inference met. So we can use the t distribution and apply the Central Limit Theorem. Now let us calculate the interval using the formula  $CI = \bar{x} \pm t * (s/\sqrt{n})$ . Note that 95% confidence interval is a range of values that you can be 95% certain the interval contains the true mean of the population.

For calculations consider appendix D. Below are the confidence intervals for men and women.

- Female : [1.62, 1.65]
- Male: [1.72, 1.75]

We can clearly see that both the confidence intervals for the group do not overlap each other which states that there is a difference in their mean values and proves that we reject the null hypothesis and supports our difference of means test.<sup>7</sup>

## 4. Results and Conclusions

From all the previous calculations done above, We get a T score = 6.71, p-value =  $2.25 \times 10^{-8}$  and hedges g= 1.34 . All these values show that there is enough evidence to infer that there is a difference between the mean of heights of females and males. Hence, we can reject the null hypothesis. Also, the p-value shows that the difference is statistically significant, and as well as the hedges g shows the effect size which is large, conveying that the difference is practically significant as well. Also, the t score is large which makes the difference more certain.

The confidence interval for males [1.72, 1.75] and females [1.62, 1.65] are not overlapping, which further supports our inference that we can reject the null hypothesis which states that there is no difference between the means of females and males. And we can now answer the research question that yes, there is a difference between the heights of males and females. Also, looking at the means we can also conclude that males are taller than females.

Our inference was inductive because we took sample sizes and then discussed the population. So we went from specific to more general and the conclusion was not contained within the scope of the premises. And as we used confidence intervals, significance levels this implies that we are not 100% sure but strongly believe in our conclusions. Our inductive argument is strong and reliable because we justified all the

---

<sup>7</sup> #confidenceinterval:I calculated and constructed an appropriate confidence interval with clearly stated steps. And I interpreted exactly the meaning of the confidence interval and applied this into following tests to support my significance test.

three conditions for the t- distribution and applied the central limit theorem. The effect size was also large which also convinces that the inference was reliable.<sup>8</sup>

**Word Count:** 1498 words

## References

CDC.(2019). YRBSS Data and documentation.  
<https://www.cdc.gov/healthyyouth/data/yrbs/data.html>

## Reflection

I once got a 1 in distributions where I did not answer the question at all, we had to calculate the p value using the t distribution and I misinterpreted the question because of my unclear concepts related to the type of distribution but the feedback on my poll helped me out a lot. The feedback described a clear difference between a z distribution and t distribution and when to use which one. For example, a t distribution is used when we do not know the population standard deviation and the sample size is small and the feedback referred back to the book for revising the concept. I did that, which was super helpful for this assignment, now I clearly know when to use a t-distribution.

## Appendix

The full Jupyter notebook file and the data can be accessed in the zipped folder submitted as a secondary file.

### Appendix A: Dividing the dataset into subgroups and analyzing

---

<sup>8</sup> #induction: I identified why statistical inference are inductive in nature and discussed about the strength and reliability of the inference.

```
In [120]: 1 #dividing the dataset into subgroups - male and female
2 height_of_female = df['height'][df['gender']=='female']
3 height_of_male = df['height'][df['gender']=='male']
4 print(height_of_female.head(10))
5 print(height_of_male.head(10))

0    1.50
5    1.57
11   1.68
12   1.67
14   1.60
15   1.60
16   1.68
18   1.73
19   1.63
20   1.63
Name: height, dtype: float64
1    1.78
2    1.75
3    1.68
4    1.70
6    1.78
7    1.63
8    1.63
9    1.83
10   1.69
13   1.78
Name: height, dtype: float64
```

```
In [21]: 1 #print the summary statistics: mean,count,SD for height of female
2 height_of_female.describe()
```

```
Out[21]: count      51.000000
mean       1.648431
std        0.070011
min        1.500000
25%       1.600000
50%       1.650000
75%       1.700000
max        1.800000
Name: height, dtype: float64
```

```
In [25]: 1 #print the summary statistics: mean,count,SD for height of male
2 height_of_male.describe()
```

```
Out[25]: count      48.000000
mean       1.747917
std        0.076990
min        1.550000
25%       1.700000
50%       1.750000
75%       1.785000
max        1.930000
Name: height, dtype: float64
```

```
In [133]: 1 #import library
2 import statistics
3 print("statistics for female heights")
4 print("median=", statistics.median(height_of_female))
5 print("mode=", statistics.mode(height_of_female))
6 print("range=", (max(height_of_female)-min(height_of_female)))

statistics for female heights
median= 1.65
mode= 1.6
range= 0.3000000000000004
```

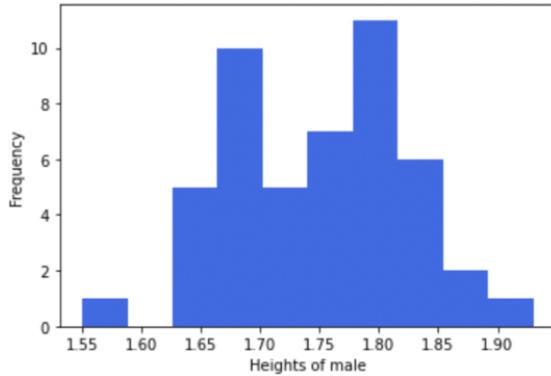
```
In [131]: 1 #using library to get relevant statistics
2 print("statistics for male heights")
3 print("median=",statistics.median(height_of_male))
4 print("mode=",statistics.mode(height_of_male))
5 print("range=", (max(height_of_male)-min(height_of_male)))
```

```
statistics for male heights
median= 1.75
mode= 1.78
range= 0.3799999999999999
```

## Appendix B: Visualizing Data

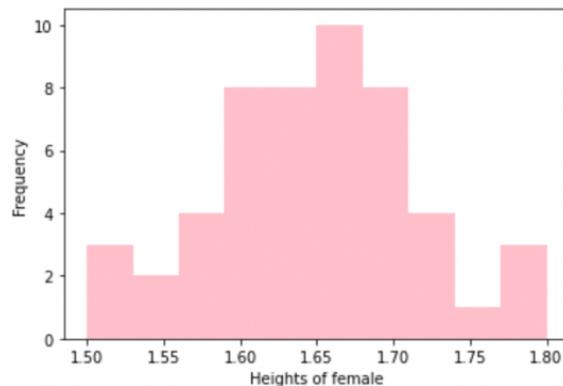
```
In [47]: 1 #plotting a histogram using the library matplotlib for males height
2 import matplotlib.pyplot as plt
3 plt.hist(height_of_male, color = 'royalblue')
4 plt.xlabel("Heights of male")
5 plt.ylabel("Frequency")
```

```
Out[47]: Text(0, 0.5, 'Frequency')
```



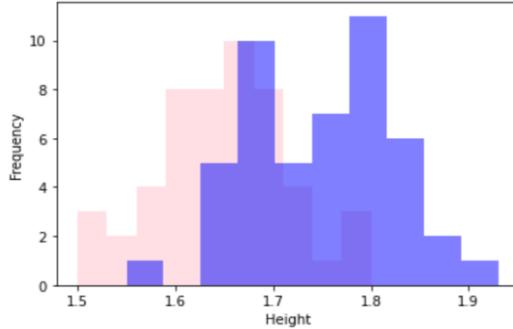
```
In [48]: 1 #plotting a histogram for females height
2 plt.hist(height_of_female, color = 'pink')
3 plt.xlabel("Heights of female")
4 plt.ylabel("Frequency")
```

```
Out[48]: Text(0, 0.5, 'Frequency')
```



```
In [111]: 1 #plotting a histogram with both variables to show comparison
2 plt.hist(height_of_female, color = 'pink', alpha= 0.5)
3 plt.hist(height_of_male, color = 'blue', alpha = 0.5)
4 plt.xlabel("Height")
5 plt.ylabel("Frequency")
```

```
Out[111]: Text(0, 0.5, 'Frequency')
```



## Appendix C: Difference of means test

```

1 #this code has been used from class Session 13.2 https://sle-collaboration.minervaproject.com/
2 #?url=https%3A//sle-authoring.minervaproject.com/api/v1/worksheets/c942ce85-aaca-4b51-b9de-d14
3 #b96d31616/&userId=11969&name=Sana+Mehta&avatar=https%3A//s3.amazonaws.com/picasso/fixtures/
4 #Sana_Mehta_11969_2021-08-30T23%3A44%3A59.323Z&noPresence=1&readOnly=1&isInstructor=0&signatu
5 #re=25b3277b489db7d96e36710ca965f652f7c937e4e4d4b95327d6db62d9c5df7b
6
7 import numpy as np
8 from scipy import stats
9
10 #defining a function
11 def difference_of_means_test(data1,data2,tails):
12
13     #finding the count of data
14     n1 = len(data1)
15     n2 = len(data2)
16
17     #finding the mean of the two groups
18     x1 = np.mean(data1)
19     x2 = np.mean(data2)
20
21     #finding the standard deviation of the dataset
22     s1 = np.std(data1,ddof=1) #bessels correction using n-1 as denominator
23     s2 = np.std(data2,ddof=1)
24
25     #finding the standard error
26     SE = np.sqrt((s1**2/n1 + s2**2/n2))
27
28
29     #finding the t score using standard error
30     Tscore = ((x2-x1)/SE)
31
32     #calculating degrees of freedom given the conservative estimate from openintro
33     df = min(n1,n2)-1
34
35     #calculating p value using stats library
36     pvalue = tails*stats.t.cdf(-Tscore,df)
37
38     #calculating pooled standard deviation using the formula
39     SDpooled = np.sqrt((s1**2*(n1-1) + s2**2*(n2-1))/(n1+n2-2))
40
41     #calculating cohens d using the formula
42     Cohensd = (x2-x1)/SDpooled
43
44     #calculating hedges g from cohens d
45     hedgesg = Cohensd*((4*(n1+n2)-9)-(3))/(4*(n1+n2)-9)
46
47     #printing all the values calculated using the function
48     print('t=',Tscore)
49     print('p =', pvalue)
50     print('d =',Cohensd)
51     print('g=',hedgesg)
52
53     #calling the function
54     difference_of_means_test(height_of_female,height_of_male,2)

t= 6.713450954353862
p = 2.22514300865667e-08
d = 1.3539963495503458
g= 1.3435002538173975

```

## Appendix D: Confidence Intervals

```
In [121]: 1 #importing relevant packages
2 #defining a function for calculating confidence interval
3 from scipy import stats
4 import numpy as np
5 def confidence_interval(mean,SD,n,level):
6     df = n-1
7     t = stats.t.ppf(1-(1-level)/2,df) #finding t score value where 1- (1-level)/2 shows
8     #sum of two tails and substrating it from 1 because to give the correct side of the area on the graph
9     interval=(mean-t*SE,mean+t*SE)
10    print("confidence_interval=",interval)
11
12 confidence_interval(1.64,0.7,51,0.95) #for heights of female
13 confidence_interval(1.74,0.7,48,0.95) #for heights of male

confidence_interval= (1.6203121827675093, 1.6596878172324905)
confidence_interval= (1.7202809988026393, 1.7597190011973607)
```