

Exercise 3 - Hypothesis Testing

Karlo Angelo C. Lazaro

2023-06-07

Problem

Refer to the Christmas spending dataset (see `christmas_spending.csv` in the class drive), which is fictional data on the total amount a random sample of customers spent in a department store during the Christmas season. It has been claimed that the customers of the department store spend a total of 5000 pesos on the average during the Christmas season. Test the claim at 0.05 level of significance.

Q1. What is your parameter of interest? What is your point estimator?

A1. The parameter of interest is the average spending of customers during Christmas season in the department store. The point estimator in this instance would be the sample mean (\bar{x}) computed from the provided data.

Q2. Load the data into R and compute the point estimate.

```
# Compute for the point estimate x-bar
s_mean <- mean(cspend$spending)
# Print the computed value of x-bar
print(s_mean)
```

A2.

```
## [1] 4658.872
```

The value of the sample mean (\bar{x}) is 4658.87 pesos.

Q3. State the problem as a two-sided hypothesis test. Give the null hypothesis and the alternative hypothesis.

A3. H_0 : The average Christmas spend of department store consumers is 5000 pesos.

H_a : The average Christmas spend of department store consumers is not 5000 pesos.

Q4. Why is the one-sample t-test appropriate for this problem?

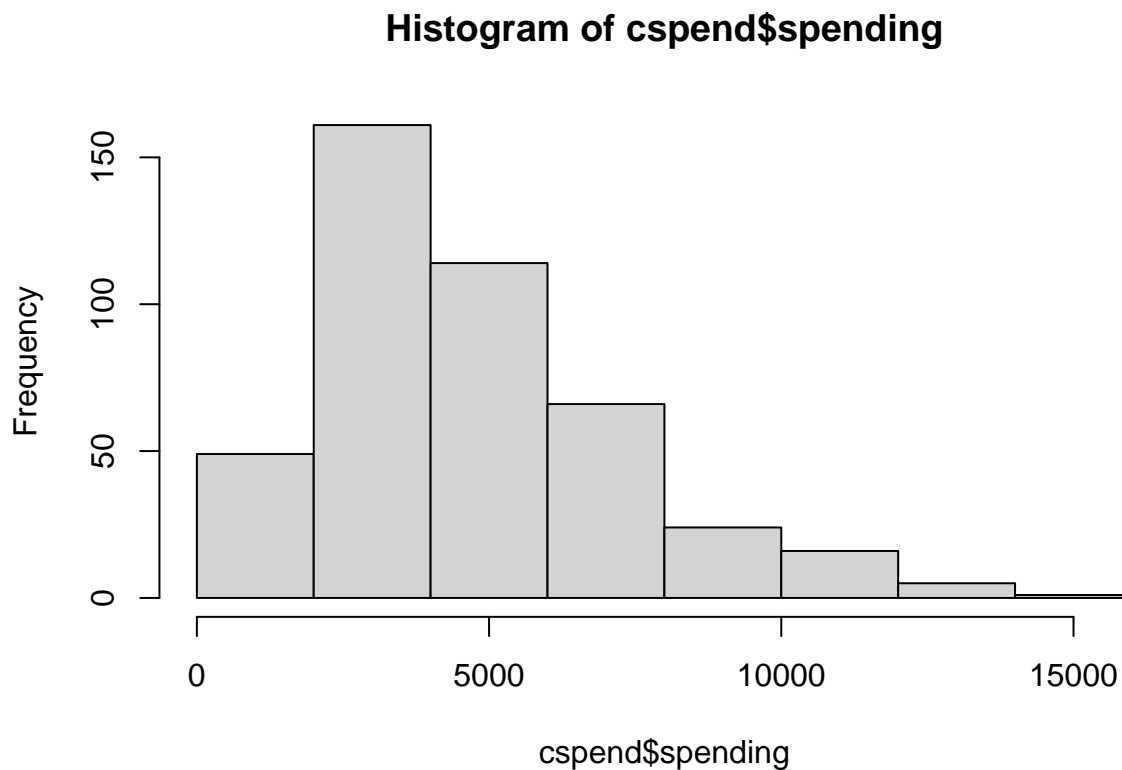
A4. Looking at the requirements for a one-sample t-test, there are four main assumptions to be met.

1. The data must be continuous.
2. It should be independent.
3. The distribution is approximately normal.
4. Homogeneity of the variance.

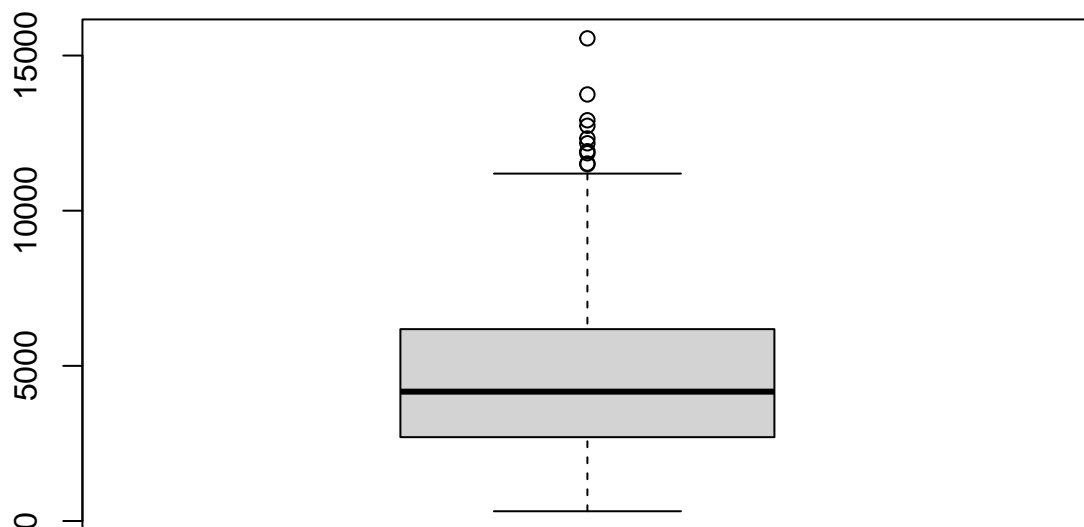
Now, looking at our given data, for (1) we have spending in pesos which is continuous. (2), one person's spending is independent of another person's so that also meets the criteria. (3), based on the histogram, our data appears to be skewed to the right. And (4), our boxplot is showing no egregious outliers.

Although our data does not meet the normality assumption, we can still go ahead with using one sample t-test as our sample size is robust enough and the violation is not that extreme. Also, the t-test is concerned with approximate normality and not strict normality. Real world data usually shows some degree of deviation from normality but as long as these deviations are acceptable and not severe then we can go ahead and use the one sample t-test.

```
# Visual test of normality of the data  
hist(cspend$spending)
```



```
# Visual test of data variability  
boxplot(cspend$spending)
```



Q5. Compute the observed value of your test statistic. The sample variance may be computed using the `var()` function in R.

A1. To compute for t we will use the following formula:

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

where,

\bar{x} is the sample mean.

μ_0 is the population mean.

s is the standard deviation, and

n is the sample size.

And right now we have the following values computed/given:

1. $\mu_0 = 5000$
2. $\bar{x} = 4658.872$
3. $n = 436$

```
# Given data from the problem
p_mean <- 5000
# Compute sample size
s_size <- NROW(cspend$spending)
```

To get the value of s , we use the following:

```

# Compute the sample variance
s_var <- var(cspend$spending)
# Compute the standard deviation
s_dev <- sqrt(s_var)
# Print the computed standard deviation
print(s_dev)

```

```
## [1] 2581.068
```

Now that we have all the variables we need to compute the one sample t-test, we can just plug them in to compute for our t statistic.

```

# Compute for the t-statistic
t <- (s_mean - p_mean) / (s_dev / sqrt(s_size))
# Print the computed t-statistic
print(t)

```

```
## [1] -2.759695
```

Q6. What is the null distribution? Don't forget to give the actual value of the parameter.

A6. The null distribution of our test statistic t follows a t-distribution with $n-1$ degrees of freedom.

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t_{n-1}$$

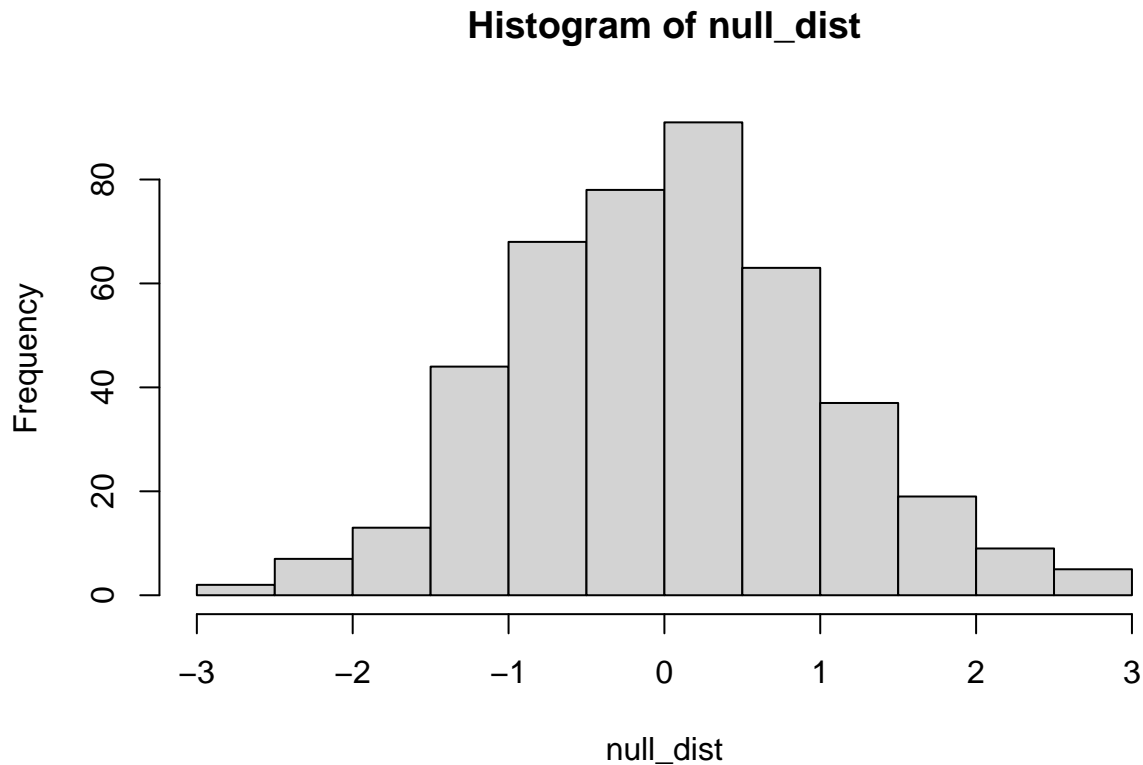
So, in this case, the null distribution follows a t-distribution with 435 degrees of freedom (t_{435}).

To visualize it further:

```

# Set seed so that the computed null distribution will be fixed
set.seed(31)
# Set the degrees of freedom
df <- 435
# Generate the null distribution
null_dist <- rt(s_size, df)
# Plot the null distribution
hist(null_dist)

```



Q7. Find the p -value corresponding to the observed value of your test statistic.

A7. For a two-sided t -test, the p -value is given by $p = P(|T| > t)$. Which is equal to the R code `2*(1-pt(abs(t), n-1))`.

Since we have all of these variables, we can compute the p -value by plugging in the different values we have.

```
# Plugging in the values
p_val <- 2*(1-pt(abs(t), df))
# Print the computed p-value
print(p_val)
```

```
## [1] 0.006029875
```

Q8. Based on your answer in (7), what is your conclusion?

A8.

9. Alternatively, what is the critical value corresponding to the significance level? You need to use the `qt()` function. Compare the critical value and the test statistic. What is your conclusion based on the critical value? Is your answer consistent with your answer in (8)?

Q10. Now perform the test again, but this time use the `t.test()` function in R. You should arrive at the same conclusion.

A10. Using the `t.test()` function, we get the following:

```
ttest <- t.test(cspend$spending, mu=p_mean, alternative = "two.sided")
print(ttest)
```

```
##
##  One Sample t-test
##
## data:  cspend$spending
## t = -2.7597, df = 435, p-value = 0.00603
## alternative hypothesis: true mean is not equal to 5000
## 95 percent confidence interval:
##  4415.924 4901.821
## sample estimates:
## mean of x
##  4658.872
```