

Part 1: Foundational Exercises

Lab 1: Hypothesis Testing - W203 Section 8

Team Anglerfish: Aruna Bisht, Jing Lange, Chase Madson, Maria Manna

Contents

1.1 Professional Magic	2
1.1.1 Type I Error Rate	2
1.1.2 Power	3
1.2 Wrong Test, Right Data	4
1.2.1 Violating t-test Assumptions	4
1.2.2 Remedial Measure	4
1.3 Test Assumptions	5
1.3.1 World Happiness	5
1.3.2 Legislators	7
1.3.3 Wine and Health	8
1.3.4 Attitudes Toward the Religious	9

1.1 Professional Magic

1.1.1 Type I Error Rate

Question: *What is the type I error rate of the test?*

- Let's denote the test statistic t with the given formula $t = X_1 + Y_1 + X_2 + Y_2 + X_3 + Y_3$.
- The type I error rate α is the probability that we incorrectly reject the null hypothesis $H_0 : p = \frac{1}{2}$. In our test, we reject the null hypothesis when the test statistic is either zero or six.

$$\begin{aligned}\alpha &= Pr(H_0 \text{ is rejected when it is true}) \\ &= Pr(t \in \{0, 6\} | p = \frac{1}{2})\end{aligned}$$

- We can rearrange $t = 0$ as 3 consecutive flips of $(X_i, Y_i) = (0, 0)$, and similarly $t = 1$ as $(X_i, Y_i) = (1, 1)$. We note that there is $\binom{6}{0} = 1$ way of observing $t = 0$ and $\binom{6}{6} = 1$ way of observing $t = 1$.

$$\begin{aligned}&= \binom{6}{0} Pr([(X_1, Y_1) = (0, 0)] \cap [(X_2, Y_2) = (0, 0)] \cap [(X_3, Y_3) = (0, 0)] | p = \frac{1}{2}) \\ &+ \binom{6}{6} Pr([(X_1, Y_1) = (1, 1)] \cap [(X_2, Y_2) = (1, 1)] \cap [(X_3, Y_3) = (1, 1)] | p = \frac{1}{2})\end{aligned}$$

- Given that each flip of the pair is independent of all other flips of the pair, we can restate as the product of probabilities.

$$\begin{aligned}&= Pr([(X_1, Y_1) = (0, 0)] | p = \frac{1}{2}) \times Pr([(X_2, Y_2) = (0, 0)] | p = \frac{1}{2}) \times Pr([(X_3, Y_3) = (0, 0)] | p = \frac{1}{2}) \\ &+ Pr([(X_1, Y_1) = (1, 1)] | p = \frac{1}{2}) \times Pr([(X_2, Y_2) = (1, 1)] | p = \frac{1}{2}) \times Pr([(X_3, Y_3) = (1, 1)] | p = \frac{1}{2})\end{aligned}$$

- We express this in terms of the joint distribution functions we are given.

$$= (f_{X,Y|p=\frac{1}{2}}(0, 0))^3 + (f_{X,Y|p=\frac{1}{2}}(1, 1))^3$$

- When $p = \frac{1}{2}$, then $f_{X_i, Y_i}(0, 0) = f_{X_i, Y_i}(1, 1) = \frac{1}{2} = \frac{1}{4}$, and thus we plug that probability into the above.

$$\begin{aligned}&= \left(\frac{1}{4}\right)^3 + \left(\frac{1}{4}\right)^3 \\ &= \frac{1}{32} = 0.03125\end{aligned}$$

1.1.2 Power

Question: What is the power of your test for the alternative hypothesis that $p = \frac{3}{4}$

- a. The power of this test is $1 - \beta$, where β is the type II error rate (i.e., the probability that we fail to reject H_0 when it is false). Specifically, we want to calculate the power under the alternative hypothesis $H_a : p = \frac{3}{4}$.

$$\begin{aligned}
 \text{Power} &= 1 - \beta \\
 &= 1 - \Pr(\text{Fail to reject } H_0 \text{ when } H_a \text{ is true instead}) \\
 &= 1 - \Pr(t \notin \{0, 6\} \mid H_a : p = \frac{3}{4}) \\
 &= 1 - [1 - \Pr(t \in \{0, 6\} \mid H_a : p = \frac{3}{4})] \\
 &= \Pr(t \in \{0, 6\} \mid H_a : p = \frac{3}{4})
 \end{aligned}$$

- b. We repeat our steps c, d, and e from our solution to 1.1.1 above to reach the joint probability density function when $p = \frac{3}{4}$.

$$= (f_{X,Y|p=\frac{3}{4}}(0,0))^3 + (f_{X,Y|p=\frac{3}{4}}(1,1))^3$$

- c. When $p = \frac{3}{4}$, then $f_{X_i,Y_i}(0,0) = f_{X_i,Y_i}(1,1) = \frac{\frac{3}{4}}{2} = \frac{3}{8}$, and thus we plug that probability into the above.

$$\begin{aligned}
 &= \left(\frac{3}{8}\right)^3 + \left(\frac{3}{8}\right)^3 \\
 &= \frac{27}{256} = \mathbf{0.10546875}
 \end{aligned}$$

1.2 Wrong Test, Right Data

Imagine that your organization surveys a set of customers to see how much they like your regular website, and how much they like your mobile website. Suppose that both of these preference statements are measured on 5-point Likert scales.

1.2.1 Violating t-test Assumptions

Question: *If you were to run a paired t-test using this data, what consequences would the violation of the metric scale assumption have for your interpretation of the test results?*

For the results of a t-test on some sample to be reliable, then each draw from the sample must come from an underlying random variable that is **metric** in scale, meaning they are continuous and numeric. As the preference statements in this survey are each measured on a 5-point Likert-scale, we know that this assumption is violated.

Running a paired t-test on Likert-scale data means calculating the test statistic $t = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$, which gives us the difference from the mean in units of standard deviation. These units lose their meaning when calculated using Likert-scale data; it doesn't make sense to subtract or take a mean of these values. This is because data from Likert-scale have **non-standard intervals**, are **categorical instead of continuous**, and have **incomparable magnitudes** as values are subjective to each customer in the sample (a “very good” from one person can mean something different than another person’s “very good”). Therefore, our interpretation of the test statistic would be invalid and so would be any conclusions we draw from it.

Optional Supplement to our Response

Assumptions: In order to run a paired t-test on some sample $\{X_i\}_{i=1}^n$ and receive reliable results, there are 3 assumptions that must be met.

1. Each X_i comes from an underlying random variable that is **metric** in scale, meaning they are continuous and numeric.
2. $\{X_i\}_{i=1}^n$ is an **independent and identically distributed** random sample.
3. $\{X_i\}_{i=1}^n$ is **normally distributed**, or $\{X_i\}_{i=1}^n \sim N(\mu_x, \sigma_x^2)$

1.2.2 Remedial Measure

Question: *What would you propose to do to remedy this problem?*

Since Likert-scale data contains paired values that are non-metric, it is better to use a **hypothesis of comparisons under the Wilcoxon Rank-Sum Test**. If we want to compare Likert-scale responses about customers' sentiment about our regular website (X) and our mobile website (Y), we would want to test the null hypothesis that $H_0 : P(X < Y) = P(X > Y)$, meaning equal likelihood. If the evidence suggests that we can reject this null hypothesis, then we can provide a useful conclusion to our organization about how our customers view these websites differently.

The assumptions required under this test are that the data must follow an ordinal scale (which Likert-scale does) and that each pair is drawn from the same distribution (which the organization can accomplish through random sampling).

Do we want to propose converting the data from Likert-scale to binary?

1.3 Test Assumptions

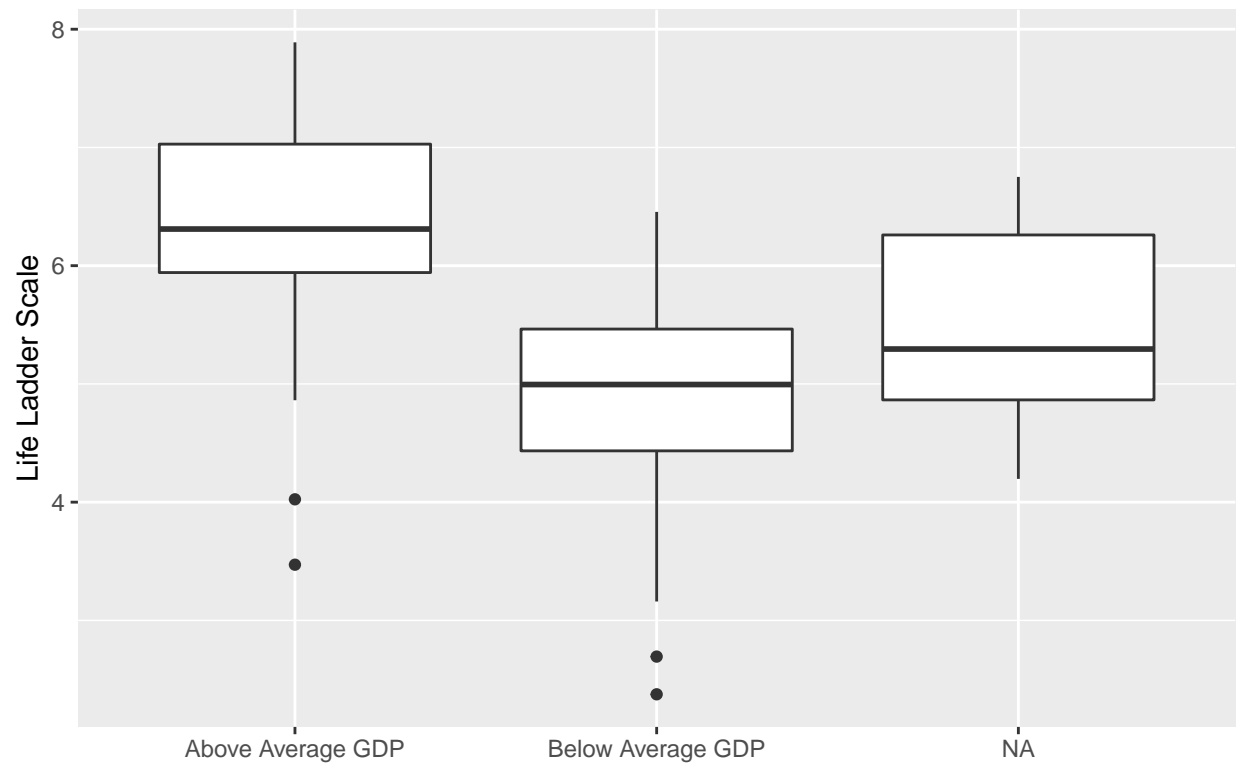
Evaluate the assumptions for each of the following tests.

1.3.1 World Happiness

Load the Data

```
happiness_WHR <-  
  'happiness_WHR.csv' %>%  
  read_csv(col_types = 'ciddddddd') %>%  
  clean_names()  
  
happiness_WHR  
  
## # A tibble: 239 x 11  
##   country_name year life_ladder log_gdp_per_capita social_support  
##   <chr>      <int>      <dbl>          <dbl>          <dbl>  
## 1 Afghanistan  2019         2.38          7.70          0.42  
## 2 Albania      2019         5.00          9.54          0.686  
## 3 Algeria      2019         4.74          9.34          0.803  
## 4 Argentina    2019         6.09          10           0.896  
## 5 Armenia      2019         5.49          9.52          0.782  
## 6 Australia    2019         7.23          10.8          0.943  
## 7 Austria      2019         7.20          10.9          0.964  
## 8 Azerbaijan   2019         5.17          9.57          0.887  
## 9 Bahrain      2019         7.10          10.7          0.878  
## 10 Bangladesh  2019         5.11          8.47          0.673  
## # ... with 229 more rows, and 6 more variables:  
## #   healthy_life_expectancy_at_birth <dbl>, freedom_to_make_life_choices <dbl>,  
## #   generosity <dbl>, perceptions_of_corruption <dbl>, positive_affect <dbl>,  
## #   negative_affect <dbl>  
  
happiness_WHR %>%  
  #filter(!is.na(log_gdp_per_capita)) %>%  
  mutate(above_average_gdp =  
    (log_gdp_per_capita >= mean(log_gdp_per_capita, na.rm = TRUE)) %>%  
    ifelse('Above Average GDP', 'Below Average GDP')) %>%  
  ggplot(mapping = aes(x = above_average_gdp, y = life_ladder)) +  
  geom_boxplot() +  
  labs(title = 'Box-Plots Comparing Above and Below Average GDP Countries',  
       x = '',  
       y = 'Life Ladder Scale')
```

Box-Plots Comparing Above and Below Average GDP Countries



1.3.2 Legislators

2 + 2

[1] 4

1.3.3 Wine and Health

```
2 + 2
```

```
## [1] 4
```


1.3.4 Attitudes Toward the Religious

2 + 2

[1] 4

Assumptions: In order to run a paired t-test on some sample $\{X_i\}_{i=1}^n$ and receive reliable results, there are 3 assumptions that must be met.

1. Each X_i comes from an underlying random variable that is **metric** in scale, meaning they are continuous and numeric.
2. $\{X_i\}_{i=1}^n$ is an **independent and identically distributed** random sample.
3. $\{X_i\}_{i=1}^n$ is **normally distributed**, or $\{X_i\}_{i=1}^n \sim N(\mu_x, \sigma_x^2)$
 - With sufficiently large n (e.g., surveying at least 30 customers), this assumption is met asymptotically since the Central Limit Theorem will cause a normal distribution of the mean.

Violations: