

TMFVC Assignment 12

Guangxin Zhao st194136

January 25, 2025

Exercise 12.1

(a)

Code:

```
1 # Exercise 12.1
2 # Load the data
3 X <- c(3.24, 4.39, 5.24, 3.83, 3.50,
4       3.75, 4.06, 3.83, 3.54, 3.20,
5       4.28, 3.65, 3.01, 4.69, 3.32)
6
7 mu_0 <- 3.5
8 sigma <- 0.76
9 n <- length(X)
10
11 mean_X <- mean(X)
12 var_X <- var(X)
13
14 # Calculate the z-statistic
15 z_stat <- (mean_X - mu_0) / (sigma / sqrt(n))
16
17 # Calculate the p-value
18 normal_pdf <- function(x) {
19   return(1 / sqrt(2 * pi) * exp(-x^2 / 2))
20 }
21
22 p_value <- integrate(normal_pdf, z_stat, Inf)$value
23
24 # Print the results
25 cat("Mean of X:", mean_X, "\n")
26 cat("Variance of X:", var_X, "\n")
27 cat("Z-statistic:", z_stat, "\n")
28 cat("P-value:", p_value, "\n")
```

Output:

```
1 Mean of X: 3.835333
2 Variance of X: 0.3718267
3 Z-statistic: 1.708869
4 P-value: 0.04373761
```

STANDARD NORMAL DISTRIBUTION: Table Values Represent AREA to the LEFT of the Z score.										
Z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.50000	.50399	.50798	.51197	.51595	.51994	.52392	.52790	.53188	.53586
0.1	.53983	.54380	.54776	.55172	.55567	.55962	.56356	.56749	.57142	.57535
0.2	.57926	.58317	.58706	.59095	.59483	.59871	.60257	.60642	.61026	.61409
0.3	.61791	.62172	.62552	.62930	.63307	.63683	.64058	.64431	.64803	.65173
0.4	.65542	.65910	.66276	.66640	.67003	.67364	.67724	.68082	.68439	.68793
0.5	.69146	.69497	.69847	.70194	.70540	.70884	.71226	.71566	.71904	.72240
0.6	.72575	.72907	.73237	.73565	.73891	.74215	.74537	.74857	.75175	.75490
0.7	.75804	.76115	.76424	.76730	.77035	.77337	.77637	.77935	.78230	.78524
0.8	.78814	.79103	.79389	.79673	.79955	.80234	.80511	.80785	.81057	.81327
0.9	.81594	.81859	.82121	.82381	.82639	.82894	.83147	.83398	.83646	.83891
1.0	.84134	.84375	.84614	.84849	.85083	.85314	.85543	.85769	.85993	.86214
1.1	.86433	.86650	.86864	.87076	.87286	.87493	.87698	.87900	.88100	.88298
1.2	.88493	.88686	.88877	.89065	.89251	.89435	.89617	.89796	.89973	.90147
1.3	.90320	.90490	.90658	.90824	.90988	.91149	.91309	.91466	.91621	.91774
1.4	.91924	.92073	.92220	.92364	.92507	.92647	.92785	.92922	.93056	.93189
1.5	.93319	.93448	.93574	.93699	.93822	.93943	.94062	.94179	.94295	.94408
1.6	.94520	.94630	.94738	.94845	.94950	.95053	.95154	.95254	.95352	.95449
1.7	.95543	.95637	.95728	.95818	.95907	.95994	.96080	.96164	.96246	.96327
1.8	.96407	.96485	.96562	.96638	.96712	.96784	.96856	.96926	.96995	.97062
1.9	.97128	.97193	.97257	.97320	.97381	.97441	.97500	.97558	.97615	.97670

Figure 1: Z-table from Z=0 to Z=1.99

According to Fig. 1, the Z-statistic $Z = 1.71$ corresponds to a cumulative probability $P(Z \leq 1.71) = 0.9564$. Using this value, the p-value for the right-tailed test is computed as $p = 1 - 0.9564 = 0.0436$. This matches the p-value obtained through numerical integration.

(b)

Since we have $p = 0.0437 < 0.05$, we reject the null-hypothesis.

Exercise 12.2

(a)

Multiple t-tests can only compare two group at a time, while ANOVA tests all group differences simultaneously. Meanwhile, if you perform multiple pairwise t-tests, the chance of a Type I error (false positive) increases because each t-test is conducted independently without adjusting for multiple comparisons.

(b)

- **Independence of observations:** The data in each group should be independent of data in other groups.
- **Normality:** The residuals of the data within each group should be normally distributed.
- **Homogeneity of variances:** The variance within each group should be the same across all groups.

(c)

Code:

```
1 # Exercise 12.2
2 library(ggplot2)
3 library(dplyr)
4 library(readxl)
5 library(car)
6
7 # Question (c)
8 # Load the data
9 data <- read_excel("hand-washing.xlsx")
10
11 summary <- data %>%
12   group_by(Method) %>%
13   summarise(mean=mean(Bacterial_Counts),
14             median=median(Bacterial_Counts),
15             sd=sd(Bacterial_Counts))
16
17 # Print the summary
18 print(summary)
19
20 # Visualize the data
21 # Box plot
22 ggplot(data, aes(x=Method, y=Bacterial_Counts, fill=Method)) +
23   geom_boxplot() +
24   geom_jitter(width=0.2, alpha=0.5) +
25   labs(title="Bacterial Counts by Hand-Washing Method (Box Plot)",
26        x="Hand-Washing Method",
27        y="Bacterial Counts") +
28   theme_minimal()
29
30 # Violin plot
31 ggplot(data, aes(x=Method, y=Bacterial_Counts, fill=Method)) +
32   geom_violin(trim=FALSE, alpha=0.7) +
33   labs(title="Bacterial Counts by Hand-Washing Method (Violin Plot)",
34        x="Hand-Washing Method",
35        y="Bacterial Counts") +
36   theme_minimal()
```

Output:

```
1 # A tibble: 4 × 4
2   Method      mean median    sd
3   <chr>      <dbl> <dbl> <dbl>
4 1 Alcohol Spray    37.5   34.5  26.6
5 2 Antibacterial Soap 92.5   91.5  42.0
6 3 Soap           106   105   47.0
7 4 Water           117  114.   31.1
```

The box plot is shown in Fig 2, and the violin plot is shown in Fig 3.

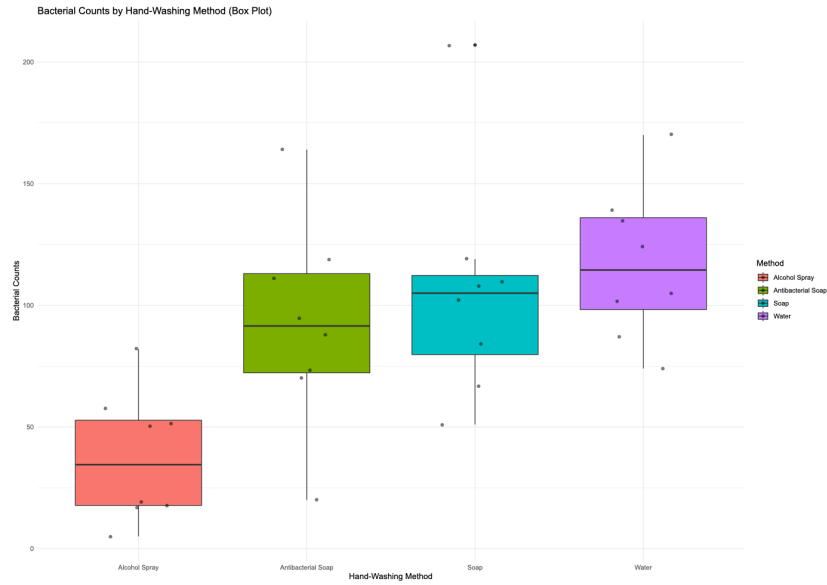


Figure 2: Box Plot of Bacterial Counts by Hand-Washing Method

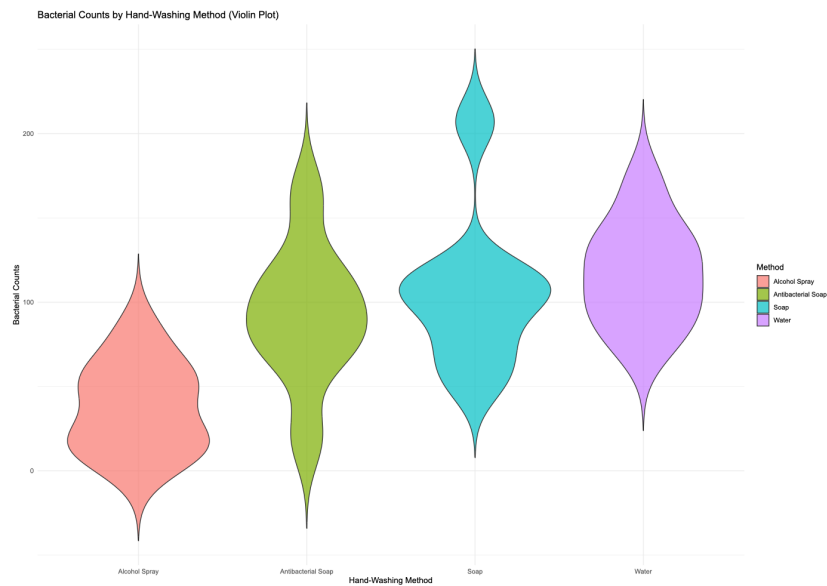


Figure 3: Violin Plot of Bacterial Counts by Hand-Washing Method

(d)

Code:

```
1 # Question (d)
2 # Test Homogeneity of Variance
3 data$Method <- as.factor(data$Method)
4 levene_test <- leveneTest(Bacterial_Counts ~ Method, data = data)
5
6 print(levene_test)
```

Output:

```
1 Levene's Test for Homogeneity of Variance (center = median)
2   Df F value Pr(>F)
3 group  3   0.1777 0.9106
4      28
```

According to Levene's Test, the p-value ($p = 0.9106$) is greater than the significance level ($\alpha = 0.05$). This means we fail to reject the null hypothesis that the variances across the groups are equal. Thus, the data meets the assumption of **homogeneity of variances**, which is a key requirement for performing ANOVA.

(e)

Code:

```
1 # Question (e)
2 # Perform ANOVA
3 anova_result <- aov(Bacterial_Counts ~ Method, data = data)
4
5 # Print the ANOVA table
6 print(summary(anova_result))
```

Output:

```
1           Df Sum Sq Mean Sq F value    Pr(>F)
2 Method      3  29882     9961   7.064 0.00111 **
3 Residuals   28  39484     1410
4 ---
5 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The ANOVA results indicate that there is a statistically significant difference in bacterial counts among the hand-washing methods ($p = 0.00111$, which is less than $\alpha = 0.05$). This means we reject the null hypothesis and conclude that at least one method differs significantly in its effectiveness at reducing bacterial counts.

(f)

Tukey's HSD (honestly significant difference) test is based on a formula very similar to that of the t-test. In fact, Tukey's test is essentially a t-test, except that it corrects for **family-wise error rate**.

(g)

Code:

```
1 # Question (g)
2 # Perform Tukey's HSD test
3 tukey_result <- TukeyHSD(anova_result)
4
5 # Print the Tukey's HSD test results
6 print(tukey_result)
7
8 # Convert Tukey HSD results to a data frame
9 tukey_data <- as.data.frame(tukey_result$Method)
10
11 # Add group comparison labels
12 tukey_data$Comparison <- rownames(tukey_data)
13
14 # View the processed data
15 print(tukey_data)
16
17 # Create the plot
18 ggplot(tukey_data, aes(x=Comparison, y=diff)) +
19   geom_point(size=4, color="blue") +
20   geom_errorbar(aes(ymin=lwr, ymax=upr), width=0.3,
21                 color="darkgray") +
22   labs(
23     title="Tukey HSD Confidence Intervals",
24     x="Group Comparisons",
25     y="Mean Difference"
26   ) +
27   theme_minimal() +
28   theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

Output:

```
1 Tukey multiple comparisons of means
2 95% family-wise confidence level
3
4 Fit: aov(formula = Bacterial_Counts ~ Method, data = data)
5 $Method
6
7 Antibacterial Soap-Alcohol Spray 55.0 3.735849 106.26415
8 Soap-Alcohol Spray 68.5 17.235849 119.76415
9 Water-Alcohol Spray 79.5 28.235849 130.76415
10 Soap-Antibacterial Soap 13.5 -37.764151 64.76415
11 Water-Antibacterial Soap 24.5 -26.764151 75.76415
12 Water-Soap 11.0 -40.264151 62.26415
13
14 p adj
15 Antibacterial Soap-Alcohol Spray 0.0319648
16 Soap-Alcohol Spray 0.0055672
17 Water-Alcohol Spray 0.0012122
18 Soap-Antibacterial Soap 0.8886944
19 Water-Antibacterial Soap 0.5675942
20 Water-Soap 0.9355196
```

The CI plot is shown in Fig 4

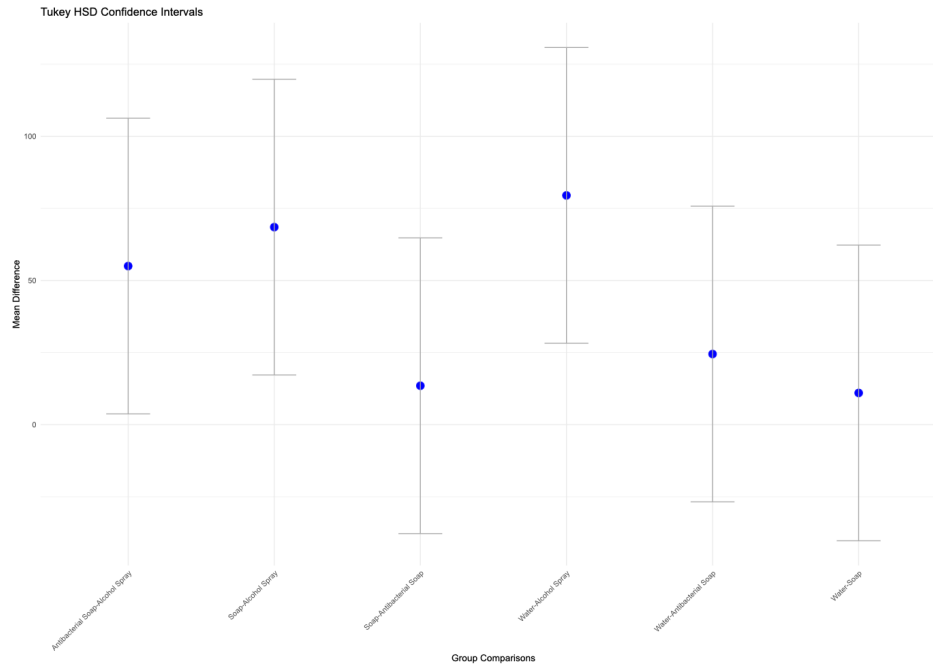


Figure 4: Confidence Interval Plot for Tukey HSD Pairwise Comparisons

(h)

1. **Increase Sample Size:** The current sample size is too small for a scientific journal. Increasing the sample size would improve statistical power and make the results more reliable.
2. **Randomize Experiment:** Conduct the experiment on random participants across different days to reduce potential biases and ensure generalizability of the results.
3. **Control Confounding Variables:** Standardize factors such as washing duration, water temperature, and handwashing technique to isolate the effect of the washing method on bacterial counts.