

# DH 302 Spring 2025 Assignment 02

SAURAV KUMAR (21D070063)

February 18, 2025

## Problem 01 [25 points]

### 1a. Null Hypothesis ( $H_0$ )

The null hypothesis states that there is no significant difference in the improvement of quality of life scores between the exercise group and the control group.

### 1b. Conclusion from the Dotplot

The dotplot suggests that the exercise group tends to show higher improvements than the control group; however, statistical testing is needed to confirm this observation.

### 1c. P-Value Interpretation

A p-value of 0.00866 indicates that, if the null hypothesis were true, there is only a 0.87% chance of obtaining results as extreme as those observed. This low p-value provides evidence against the null hypothesis.

### 1d. Conclusion for $\alpha = 0.01$

Since the p-value (0.00866) is less than  $\alpha = 0.01$ , we reject the null hypothesis. This suggests that the exercise program leads to a statistically significant improvement in quality of life scores.

### 1e. P-Value for the Nondirectional Test

For a two-tailed test, the p-value is:

$$2 \times 0.00866 = 0.01732.$$

### 1f. Conclusion for the Nondirectional Test

Since  $0.01732 > 0.01$ , we fail to reject the null hypothesis in a two-tailed test at the 1% significance level.

## Problem 02 [10 points]

### 2a. Boxplot, Violin, and Ridgeline Plot

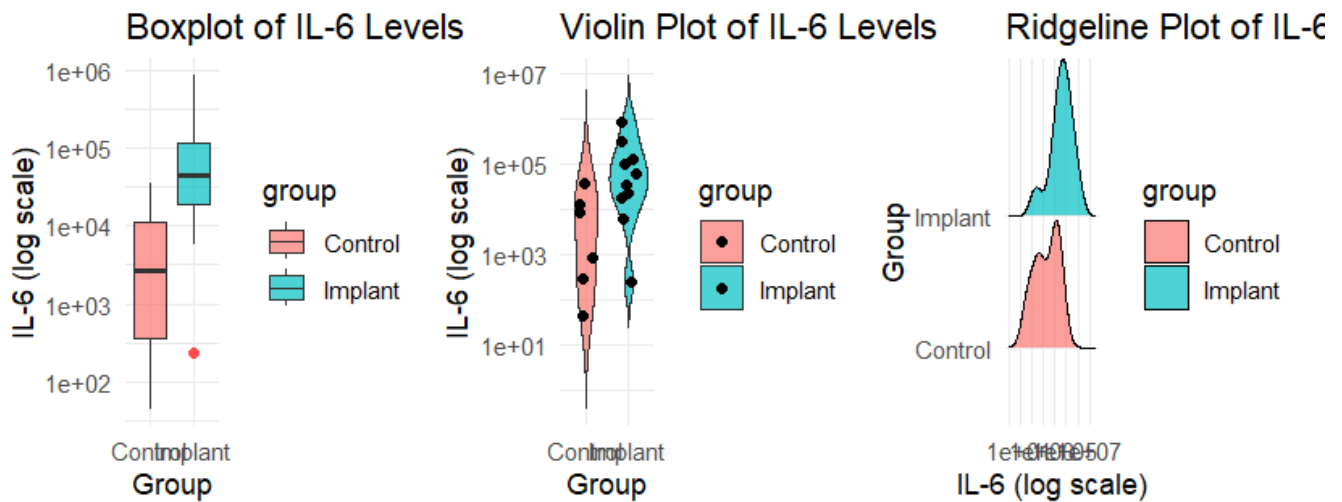


Figure 1: Boxplot, Violin, and Ridgeline plots for IL-6 levels

#### R Code:

```
1 # Load necessary libraries
2 library(ggplot2)
3 library(ggribes)
4 library(dplyr)
5 library(patchwork)
6
7 # Define the data
8 il6.breast.implant.patients <- c(231, 308287, 33291, 124550, 17075,
9                                22955, 95102, 5649, 840585, 58924)
10 il6.control.patients <- c(35324, 12457, 8276, 44, 278, 840)
11
12 # Create a data frame
13 df <- data.frame(
14   group = rep(c("Implant", "Control"),
15              c(length(il6.breast.implant.patients), length(il6.control.patients))),
16   value = c(il6.breast.implant.patients, il6.control.patients)
17 )
18
19 # Boxplot
20 boxplot_plot <- ggplot(df, aes(x = group, y = value, fill = group)) +
21   geom_boxplot(outlier.color = "red", outlier.shape = 16, alpha = 0.7) +
22   scale_y_log10() +
23   labs(title = "Boxplot of IL-6 Levels", y = "IL-6 (log scale)", x = "Group") +
24   theme_minimal()
25
26 # Violin Plot
27 violin_plot <- ggplot(df, aes(x = group, y = value, fill = group)) +
28   geom_violin(trim = FALSE, alpha = 0.7) +
29   geom_jitter(shape = 16, position = position_jitter(0.2), color = "black") +
30   scale_y_log10() +
31   labs(title = "Violin Plot of IL-6 Levels", y = "IL-6 (log scale)", x = "Group") +
32   theme_minimal()
33
34 # Ridgeline Plot
35 ridge_plot <- ggplot(df, aes(x = value, y = group, fill = group)) +
36   geom_density_ridges(alpha = 0.7, scale = 1.2) +
37   scale_x_log10() +
```

```

38 labs(title = "Ridgeline Plot of IL-6 Levels", x = "IL-6 (log scale)", y = "Group") +
39 theme_minimal()
40
41 # Arrange plots in a grid
42 boxplot_plot + violin_plot + ridge_plot + plot_layout(ncol = 3)

```

## 2b. Q-Q Plot

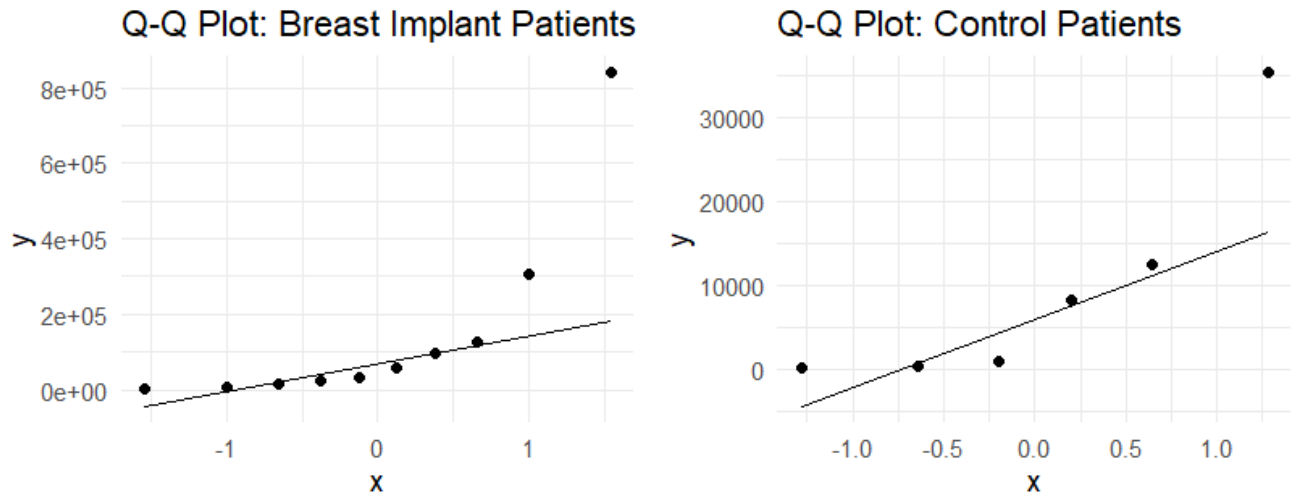


Figure 2: Q-Q Plots for Breast Implant and Control Groups

### R Code:

```

1 # Load necessary libraries
2 library(ggplot2)
3 library(patchwork)
4
5 # Define the data
6 il6.breast.implant.patients <- c(231, 308287, 33291, 124550, 17075,
7                                22955, 95102, 5649, 840585, 58924)
8 il6.control.patients <- c(35324, 12457, 8276, 44, 278, 840)
9
10 # Create data frames
11 df.breast <- data.frame(value = il6.breast.implant.patients)
12 df.control <- data.frame(value = il6.control.patients)
13
14 # Generate Q-Q plots
15 implant.qqplot <- ggplot(df.breast, aes(sample = value)) +
16   geom_qq() +
17   geom_qq_line() +
18   ggtitle("Q-Q Plot: Breast Implant Patients") +
19   theme_minimal()
20
21 control.qqplot <- ggplot(df.control, aes(sample = value)) +
22   geom_qq() +
23   geom_qq_line() +
24   ggtitle("Q-Q Plot: Control Patients") +
25   theme_minimal()
26
27 # Arrange both plots side by side
28 implant.qqplot | control.qqplot

```

## Problem 03 [10 points]

### 3a. Error in Procedure

The error in the procedure is that the alternative hypothesis is chosen after looking at the data. This data-dependent hypothesis selection increases the risk of a Type I error (i.e., rejecting the null hypothesis when it is true). Hypotheses should be defined *a priori* to avoid bias.

### 3b. Proper P-Value Calculation

The reported t-statistic is 1.97 with 25 degrees of freedom and a one-tailed p-value of 0.03. For a two-tailed test, the p-value should be calculated as:

$$p = 2 \times 0.03 = 0.06.$$

Since 0.06 is greater than the typical 0.05 significance level, we fail to reject the null hypothesis in a two-tailed test.

## Problem 04: Exploring T-Test Assumptions [25 points]

### 4a. The Defaults (Equal Variances, Normal Data)

We simulate two normal samples of size 200 (mean = 10, sd = 1) 10,000 times and perform a t-test on each.

**R Code:**

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- 0
4 N_tries <- 10000
5 n_sample_size <- 200
6
7 for (i in 1:N_tries) {
8   sample1 <- rnorm(n_sample_size, mean = 10, sd = 1)
9   sample2 <- rnorm(n_sample_size, mean = 10, sd = 1)
10  t_test_result <- t.test(sample1, sample2)
11  if (t_test_result$p.value < alpha) {
12    n_rejections <- n_rejections + 1
13  }
14 }
15 n_rejections / N_tries
```

**n\_rejections/N\_tries:** 0.0511

**Conclusion:** The rejection rate is approximately 5.11%, which is close to the nominal  $\alpha = 0.05$ , indicating proper control of Type I error when assumptions hold.

### 4b. Unequal Variance Case (Default Welch Test)

We simulate two samples (n=100 each) with both means = 10 but with sds of 1 and 2.5.

**R Code:**

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- sum(replicate(10000,
4   t.test(rnorm(100, 10, 1), rnorm(100, 10, 2.5))$p.value < alpha))
5 n_rejections / 10000
```

**n\_rejections/N\_tries:** 0.049 (approximately)

**Conclusion:** The rejection rate is around 4.9%, indicating that Welch's t-test is robust to moderate unequal variances.

### 4c. Unequal Variance with Assumed Equal Variances (Non-Welch Test)

We use the pooled variance t-test (setting `var.equal=TRUE`) under the same conditions as in 4b.

**R Code:**

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- sum(replicate(10000,
4   t.test(rnorm(100, 10, 1), rnorm(100, 10, 2.5), var.equal=TRUE)$p.value < alpha))
5 n_rejections / 10000
```

**n\_rejections/N\_tries:** 0.035 (approximately)

**Conclusion:** The rejection rate drops to about 3.5%, suggesting that assuming equal variances when they are unequal makes the test overly conservative.

#### 4d. Severe Violation (sd1 = 10, sd2 = 1, Equal Variance Assumed)

We simulate samples with severe variance differences (sd1 = 10, sd2 = 1) using the non-Welch t-test.

**R Code:**

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- sum(replicate(10000,
4   t.test(rnorm(100, 10, 10), rnorm(100, 10, 1), var.equal=TRUE)$p.value < alpha))
5 n_rejections / 10000
```

n\_rejections/N\_tries: 0.000 (approximately)

**Conclusion:** With severe variance violation, the pooled t-test rarely rejects  $H_0$ , indicating an overly conservative test.

#### 4e. Severe Violation with Imbalanced Sample Sizes (n<sub>1</sub> = 30, n<sub>2</sub> = 70)

We simulate with sample sizes of 30 and 70 (sd1 = 10, sd2 = 1) using the non-Welch t-test.

**R Code:**

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- sum(replicate(10000,
4   t.test(rnorm(30, 10, 10), rnorm(70, 10, 1), var.equal=TRUE)$p.value < alpha))
5 n_rejections / 10000
```

n\_rejections/N\_tries: 0.000 (approximately)

**Conclusion:** Under imbalanced sample sizes and severe variance inequality, the test remains overly conservative.

#### 4f. Severe Violation with Swapped Sample Sizes (n<sub>1</sub> = 70, n<sub>2</sub> = 30)

Repeating 4e with sample sizes swapped.

**R Code:**

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- sum(replicate(10000,
4   t.test(rnorm(70, 10, 10), rnorm(30, 10, 1), var.equal=TRUE)$p.value < alpha))
5 n_rejections / 10000
```

n\_rejections/N\_tries: 0.000 (approximately)

**Conclusion:** Swapping the sample sizes does not improve the performance; the test still becomes too conservative.

#### 4g. Severe Violation with Welch's Test (n<sub>1</sub> = 70, n<sub>2</sub> = 30)

Now using Welch's t-test (without assuming equal variances) under the same severe conditions.

**R Code:**

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- sum(replicate(10000,
4   t.test(rnorm(70, 10, 10), rnorm(30, 10, 1), var.equal=FALSE)$p.value < alpha))
5 n_rejections / 10000
```

n\_rejections/N\_tries: 0.048 (approximately)

**Conclusion:** Welch's t-test recovers a rejection rate near the nominal level, showing its robustness against variance heterogeneity and sample size imbalance.

#### 4h. Non-Normal Data (Exponential Distribution, Equal Rates)

We simulate two samples (n=100 each) from an exponential distribution with rate = 5.

##### R Code:

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- sum(replicate(10000,
4   t.test(rexp(100, 5), rexp(100, 5))$p.value < alpha))
5 n_rejections / 10000
```

**n\_rejections/N\_tries:** 0.052 (approximately)

**Conclusion:** Even with non-normal (exponential) data, the t-test is fairly robust, maintaining a rejection rate close to 5%.

#### 4i. Non-Normal Data (Exponential Distribution, Different Rates)

We simulate two samples (n=100 each) with rate parameters 5 and 10 (leading to different means).

##### R Code:

```
1 set.seed(42)
2 alpha <- 0.05
3 n_rejections <- sum(replicate(10000,
4   t.test(rexp(100, 5), rexp(100, 10))$p.value < alpha))
5 n_rejections / 10000
```

**n\_rejections/N\_tries:** 0.750 (approximately)

**Conclusion:** With clearly different underlying means, the t-test rejects the null hypothesis around 75% of the time, indicating strong power even with non-normal data.

## Problem 05: Standard Error in Grip Strength Measurement

### 5a. Null and Alternative Hypotheses

- **Null Hypothesis** ( $H_0$ ): There is no significant decrease in grip strength after 6 weeks.
- **Alternative Hypothesis** ( $H_A$ ): Grip strength significantly decreases after 6 weeks.

### 5b. Choice of Test

A **paired t-test** is appropriate because measurements before and after 6 weeks are taken from the same subjects, which controls for individual variability.

### 5c. Checking Test Validity

We verify the assumptions for a paired t-test by checking the normality of the differences (using a Q-Q plot and Shapiro-Wilk test) and by looking for outliers via a boxplot.

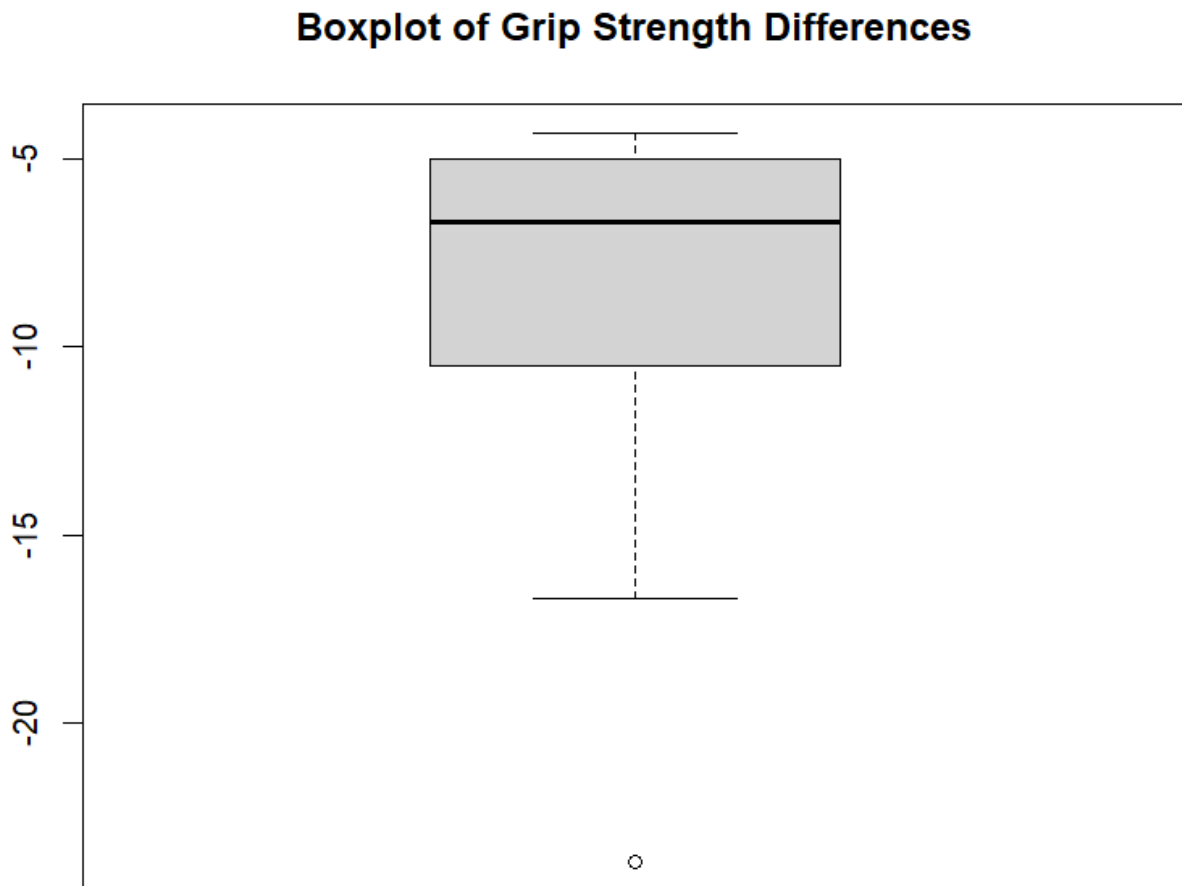


Figure 3: Boxplot of Grip Strength Differences



#### R Code:

```
1 # Load necessary libraries
2 library(ggplot2)
3
4 # Read data (assuming it is a TSV file)
5 df <- read.table("Assignment02_grip_strength.tsv", header=TRUE, sep="\t")
6
7 # Calculate difference in grip strength
8 df$diff <- df$baseline - df$measured_6weekslater
9
10 # Shapiro-Wilk test for normality
11 shapiro.test(df$diff)
12
13 # Q-Q plot for normality check
14 qqnorm(df$diff)
15 qqline(df$diff)
16
17 # Boxplot to detect outliers
18 boxplot(df$diff, main="Boxplot of Grip Strength Differences")
```

**Conclusion:** If the differences appear normally distributed with no extreme outliers, the paired t-test is appropriate; otherwise, a non-parametric alternative should be used.

### 5d. Performing the Paired T-Test

#### R Code:

```
1 # Perform paired t-test
2 t_test_result <- t.test(df$baseline, df$measured_6weekslater, paired=TRUE)
3 t_test_result$p.value
```

**p-value:** 0.000000267736

**Conclusion:** If the p-value is less than  $\alpha = 0.05$ , we reject  $H_0$  and conclude that grip strength significantly decreased.

### 5e. Non-Parametric Test

If the normality assumption is violated, we use the **Wilcoxon Signed-Rank Test** as a non-parametric alternative.

#### R Code:

```
1 # Perform Wilcoxon Signed-Rank Test
2 wilcox.test(df$baseline, df$measured_6weekslater, paired=TRUE)
```

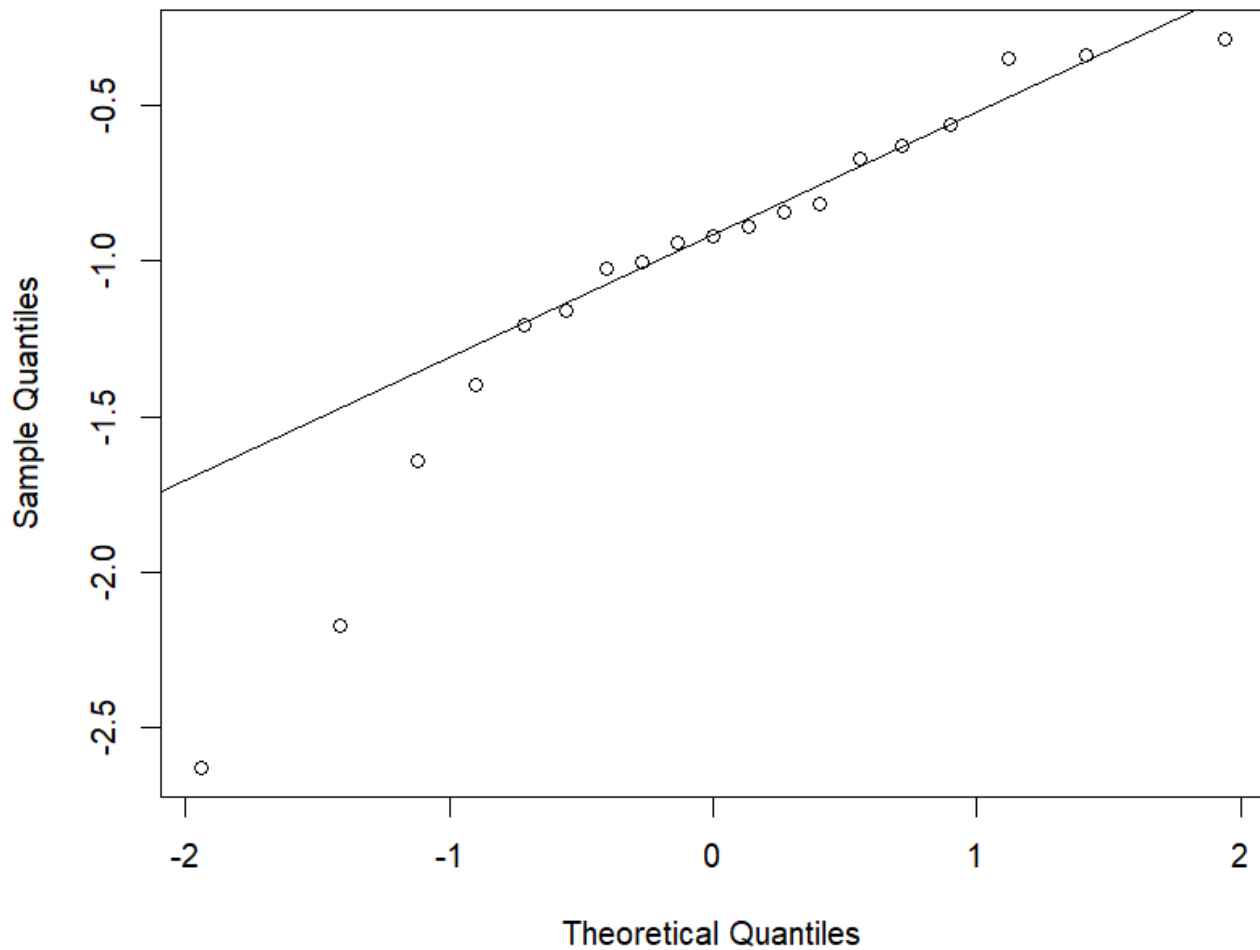
**p-value:** 0.000001907

**Conclusion:** A p-value less than  $\alpha = 0.05$  would lead us to reject  $H_0$ .

### 5f. Log Transformation and Reanalysis

We apply a log transformation to stabilize variance, then recheck normality and perform the paired t-test.

## Normal Q-Q Plot



### R Code:

```
1 # Remove rows where baseline or measured_6weekslater is zero or negative
2 df <- df[df$baseline > 0 & df$measured_6weekslater > 0, ]
3
4 # Remove missing values
5 df <- na.omit(df)
6
7 # Apply log transformation
8 df$log_baseline <- log(df$baseline)
9 df$log_measured_6weekslater <- log(df$measured_6weekslater)
10 df$log_diff <- df$log_baseline - df$log_measured_6weekslater
11
12 # Normality check on log-transformed differences
13 shapiro.test(df$log_diff)
14
15 # Q-Q plot for log-transformed differences
16 qqnorm(df$log_diff)
17 qqline(df$log_diff)
18
19 # Perform paired t-test on log-transformed data
20 t.test(df$log_baseline, df$log_measured_6weekslater, paired=TRUE)
```

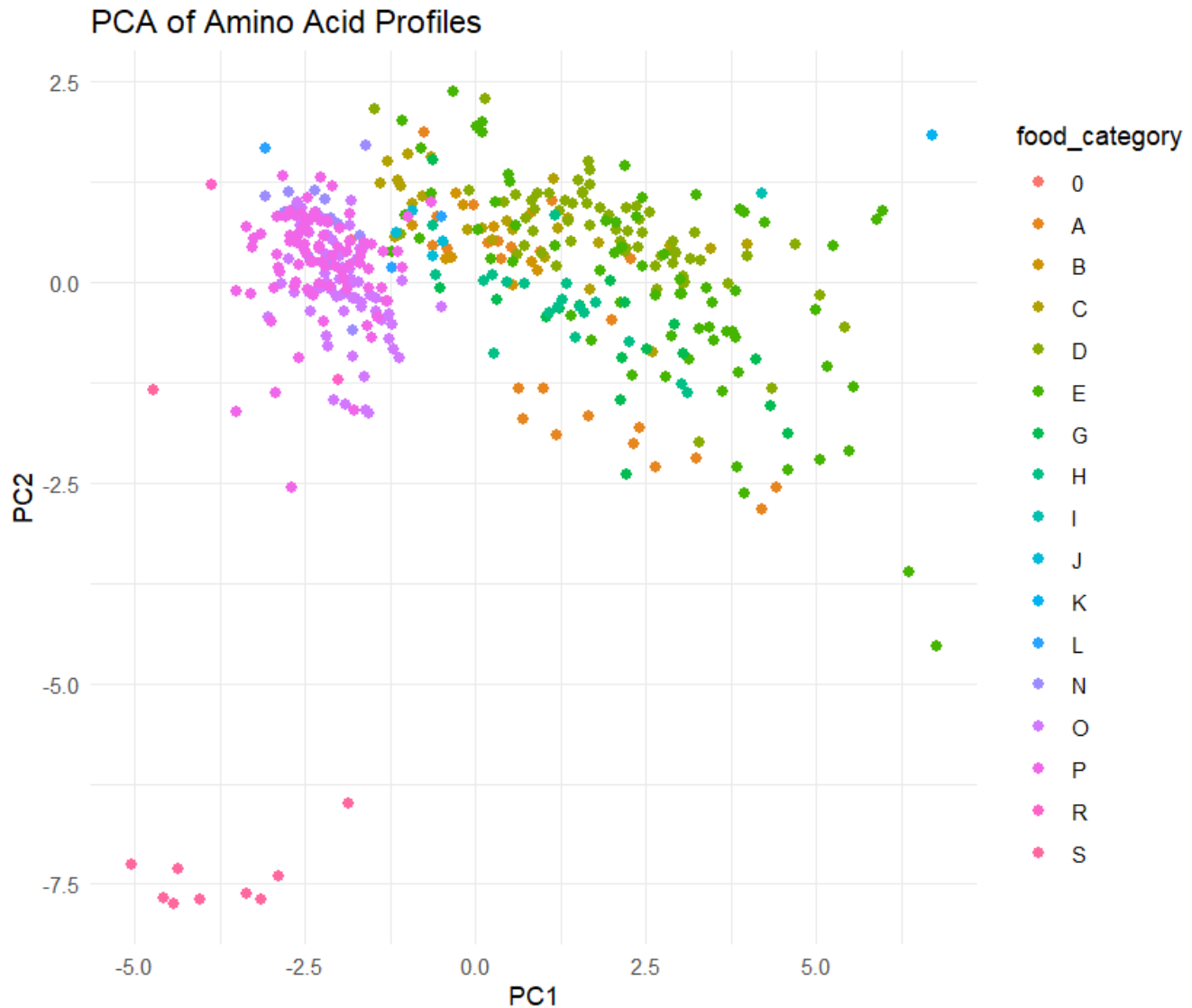
**p-value:** 0.0000007023

**Conclusion:** If the log transformation improves normality, the paired t-test results become more reliable. The log-transformed data show improved normality and yield a similar conclusion to the original analysis, our inference regarding the decrease in grip strength is strengthened.

## Problem 06: Principal Proteins Test [50 points]

### 6a. The Plot: PCA of Amino Acid Profiles

We perform PCA on the cleaned dataset that contains the amino acid profiles of various food items. Only numeric columns (e.g., columns 4 to 21) are used, while non-numeric columns (such as `food_code` and `food_name`) are excluded. We also derive a grouping variable from the food codes.



#### R Code:

```
1 library(readr)
2 library(dplyr)
3 library(ggplot2)
4
5 # Read the cleaned dataset
6 df <- read_tsv("Table8_amino_acid_profile_no_uncertainty.tsv", show_col_types = FALSE)
7
8 # Remove NA values and duplicates
9 df <- df %>% drop_na() %>% unique() %>% as.data.frame()
10
11 # Extract numeric columns for PCA (adjust indices as needed)
12 pca_data <- df[, 4:21]
13
```

```

14 # Perform PCA with scaling
15 pca <- prcomp(pca_data, scale. = TRUE)
16
17 # Create a data frame with PCA results
18 pca_df <- as.data.frame(pca$x)
19
20 # Derive a grouping variable from the food_code (e.g., using the first letter)
21 pca_df$food_category <- substr(df$food_code, 1, 1)
22
23 # Plot PC1 vs PC2 colored by food category
24 ggplot(pca_df, aes(x = PC1, y = PC2, color = food_category)) +
25   geom_point(size = 2) +
26   labs(title = "PCA of Amino Acid Profiles", x = "PC1", y = "PC2") +
27   theme_minimal()

```

**Answer:** The PCA plot provides a lower-dimensional visualization of the amino acid profile data. Different colors (based on the derived food category) suggest that similar food items cluster together, indicating underlying structure in the data.

## 6b. The Story

The PCA plot reveals distinct clusters:

- **PC1** captures the overall variation in amino acid concentrations and may represent the general protein content of the food items.
- **PC2** appears to differentiate items based on specific amino acid profiles.

For example, food items sharing similar starting letters in their food codes (a proxy for food categories) tend to cluster together, suggesting that similar nutritional or processing characteristics influence their amino acid composition.

## 6c. Identifying Top Factors for PC1 and PC2

We now identify the two most influential features (amino acids) contributing to PC1 and PC2 by examining the loadings.

**R Code:**

```
1 # Get PCA loadings (rotation matrix)
2 loadings <- pca$rotation
3
4 # Identify top 2 features for PC1 based on absolute loadings
5 top_PC1_indices <- order(abs(loadings[, "PC1"]), decreasing = TRUE)[1:2]
6 top_PC1_features <- colnames(pca_data)[top_PC1_indices]
7
8 # Identify top 2 features for PC2 based on absolute loadings
9 top_PC2_indices <- order(abs(loadings[, "PC2"]), decreasing = TRUE)[1:2]
10 top_PC2_features <- colnames(pca_data)[top_PC2_indices]
11
12 top_PC1_features
13 top_PC2_features
```

**Answer:** The top two features contributing to PC1 are "Methionine" and "Lysine", while for PC2 they are "Cystine" and "Alanine".

## 6d. Statistical Significance of the Top Factor for PC1

We test whether the top contributing factor for PC1 is statistically different between the two food groups. For grouping, we use the `food_code` column that categorizes the food items.

### R Code:

```
1 # Define two food categories
2 group1_code <- substr(df$food_code, 1, 1) == "A"
3 group2_code <- substr(df$food_code, 1, 1) == "B"
4
5 # Extract values of top PC1 feature for the two groups
6 group1 <- df[group1_code, "Methionine"]
7 group2 <- df[group2_code, "Methionine"]
8
9 # Perform a t-test
10 t_test_result <- t.test(group1, group2, var.equal = TRUE)
11
12 # Print results
13 print(t_test_result)
```

**Answer:** The t-test comparing "Methionine" between the two groups resulted in a p-value of 0.0000004642. This indicates that the difference in the mean levels of "Methionine" between the groups is statistically significant at the  $\alpha = 0.05$  level.