# Seminar 1
## Yufeng Deng      Yuanqing Wang

**Task1**

To determine the effectiveness of the treatment KTH001, it is essential to assess whether there is a statistically significant difference between the two samples. Typically, a t-test is employed for this purpose, under the assumption that the data are normally distributed. Therefore, as a preliminary step, it is necessary to evaluate the normality of both samples using the Shapiro-Wilk test.

Upon conducting the Shapiro-Wilk test on both samples, the results indicated that they are normally distributed, as the p-values for both samples exceed the 0.05 threshold.

```r
1. #import data
2. data <- read.csv("data_task1.csv")
3. placebo <- data$placebo
4. interve <- data$intervention
5. #Shapiro-Wilk test
6. shap_pla <- shapiro.test(placebo)
7. shap_int <- shapiro.test(interve)
8. print(shap_pla)
9. print(shap_int)
```

```
Output:
Shapiro-Wilk normality test
data: placebo
W = 0.95011, p-value = 0.1702
Shapiro-Wilk normality test
data: interve
W = 0.98567, p-value = 0.9481
```

However, it is important to note that the p-value is influenced by the sample size; smaller sample sizes are more likely to yield higher p-values. To bolster our confidence in the normality assumption, it is advisable to use graphical methods, such as the Quantile-Quantile (QQ) plot.

We generated QQ plots for both samples (Figure 1.1 & 1.2) and observed that the majority of points align closely with the reference line, indicating that the samples are approximately normally distributed. Although there are deviations at the extremes in the placebo sample, which may suggest some differences in the distribution's tails, these deviations are not substantial. Overall, we can reasonably assume that both samples are normally distributed.

```r
1. library(ggplot2)
2. ggplot(data, aes(sample = placebo)) +
3.   stat_qq() +
4.   stat_qq_line(color = "red") +
5.   labs(title = "QQ plot for placebo",
6.        x = "norm Quantiles",
7.        y = "Sample Quantiles") +
8.   theme_minimal()
9.
```

```
10. ggplot(data, aes(sample = interve)) +
11.   stat_qq() +
12.   stat_qq_line(color = "red") +
13.   labs(title = "QQ plot for intervention",
14.       x = "norm Quantiles",
15.       y = "Sample Quantiles") +
16.   theme_minimal()
```
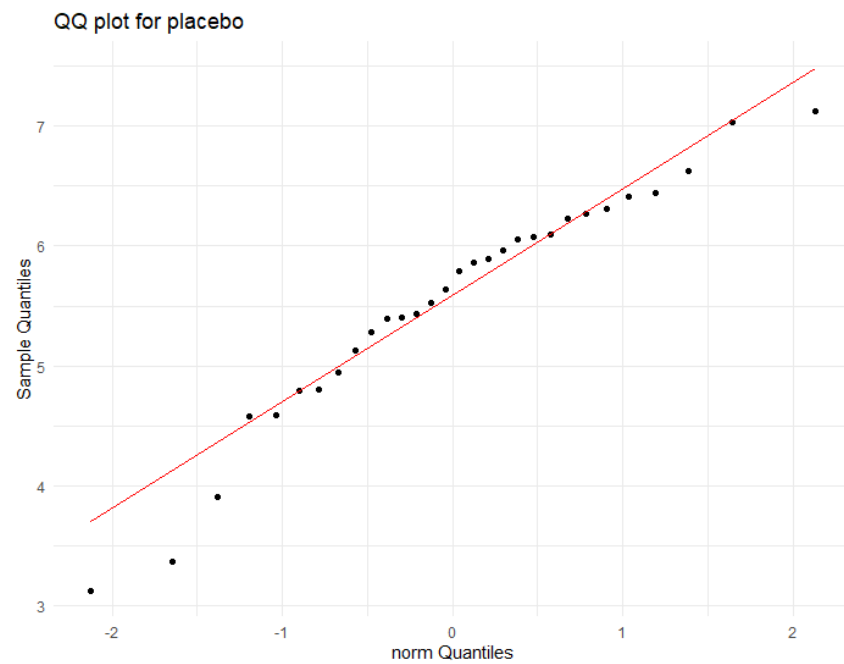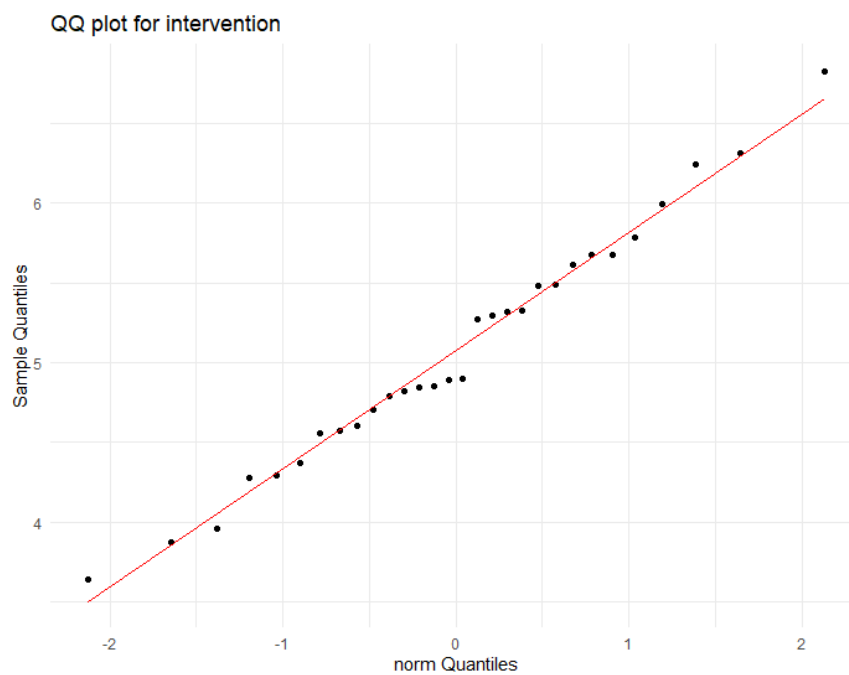


Figure1.1



Figure1.2

Given that the two samples are normally distributed, we aim to use the t-test to determine whether there is a significant difference between their means. Since the samples are independent and not paired, a paired t-test is inappropriate. Instead, we should employ an independent t-test.

Before proceeding with the t-test, we conducted an F-test to assess whether there is a significant difference between the variances of the two samples:

```
var.test(placebo, interve)
```

```
Output:
F test to compare two variances
data: placebo and interve
F = 1.6494, num df = 29, denom df = 29, p-value = 0.1839
alternative hypothesis: true ratio of variances is not equal to 1
```

Since the p-value exceeds 0.05, we conclude that there is no significant difference between the variances of the two samples. Hence, we used the independent t-test which assumes the two samples have the same variance:

```
t_test_result <- t.test(placebo, interve, var.equal = TRUE)
print(t_test_result)
```

```
Output:
Two Sample t-test
data: placebo and interve
t = 2.0512, df = 58, p-value = 0.04478
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 0.01110179 0.90989821
sample estimates:
mean of x mean of y
  5.5366   5.0761
```

Since the p-value is less than 0.05, this result indicates a statistically significant difference between the means of the placebo and intervention samples. Given that the mean of the intervention sample is lower, we can infer that the intervention may have had a measurable effect.

**Task2**

Given our assumption that the data in Task 1 are normally distributed, we can generate samples with any desired number of patients. Additionally, we can specify the effect size as needed. The means and standard deviations (SD) of each sample are presented in Table 2.1.

<div align="center">Table2.1</div>

| Sample | Mean | SD |
|---|---|---|
| Placebo | 5.5366 | 0.9702 |
| Intervention | 5.0761 | 0.7555 |

We generated new samples with varying sizes (ranging from 10 to 200, in increments of 10) while maintaining the original samples' means and standard deviations. Initially, we assumed an effect size (Cohen's d) of 0.5 and a significance level of 0.05. The resulting plot Figure 2.1 illustrates an increasing trend in statistical power as the sample size increases.

```r
1. library(pwr)
2. # Assume Cohen's d = 0.5, alpha = 0.05
3. effect_size <- 0.5
4. alpha <- 0.05
5. sample_sizes <- seq(10, 200, by = 10)
6. powers <- sapply(sample_sizes, function(n) {
7.   pwr.t.test(n = n, d = effect_size, sig.level = alpha, type = "two.sample", alternative = "two.sided")$power})
8. plot(sample_sizes, powers, type = "b", xlab = "Sample Size (per group)", ylab = "Power",
9.     main = "Power vs. Sample Size")
```
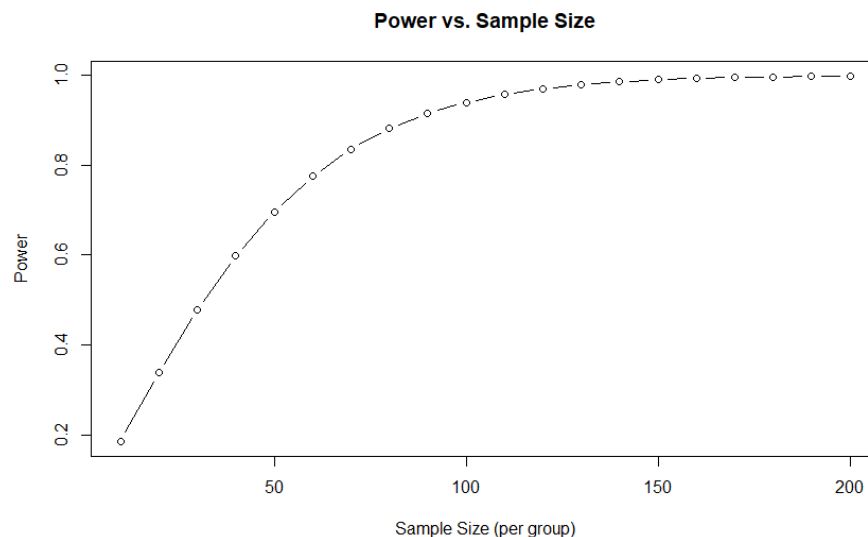


Figure2.1

Additionally, we generated curves for various effect sizes (Figure 2.2) and observed that the effect size significantly influences the power value. When the power does not approach 1, for a given sample size, an increase in effect size leads to a substantial rise in power, particularly when the effect size is small.

```r
1. sample_sizes <- seq(10, 200, by = 10)
2. effect_sizes <- c(0.1, 0.3, 0.5, 0.7, 0.9)
3. power_matrix <- outer(sample_sizes, effect_sizes, function(n, d) {
4.   pwr.t.test(n = n, d = d, sig.level = alpha, type = "two.sample", alternative = "two.sided")$power
5. })
6. power_df <- data.frame(Sample_Size = rep(sample_sizes, times = length(effect_sizes)),
7.                 Effect_Size = rep(effect_sizes, each = length(sample_sizes)),
8.                 Power = as.vector(power_matrix))
```

```
 9. library(ggplot2)
10. ggplot(power_df, aes(x = Sample_Size, y = Power, color = as.factor(Effect_Size))) +
11.   geom_line() +
12.   labs(x = "Sample Size (per group)", y = "Power", color = "Effect Size",
13.       title = "Power vs. Sample Size for Different Effect Sizes") +
14.   theme_minimal()
```
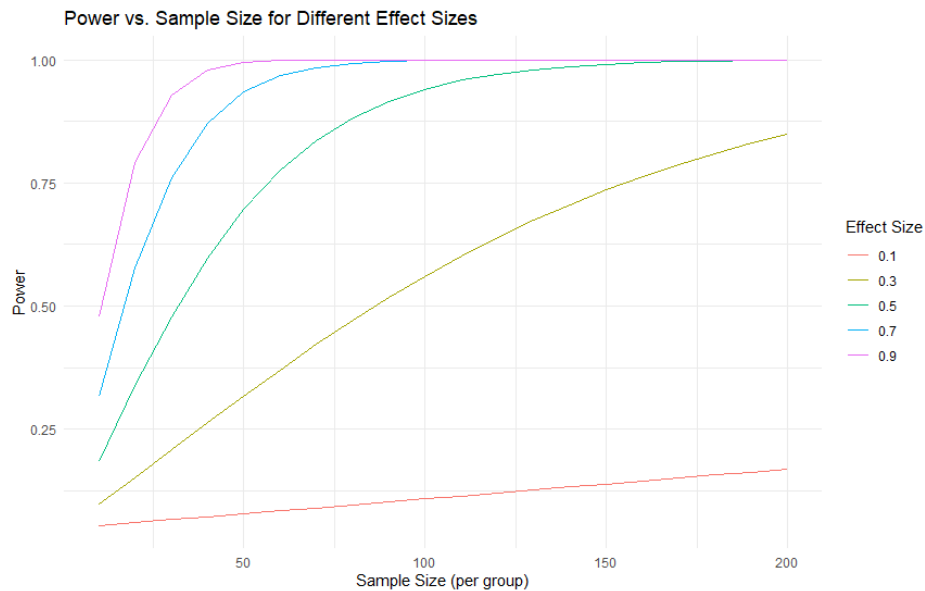


Figure2.2

## Task 3

Just like the Task1, we should first test whether these samples are normally distributed by using Shapiro-Wilk test.

```
 1. data <- read.csv("data_task3_crp.csv")
 2. plac <- data$crp_placebo
 3. int1 <- data$crp_intervention_1
 4. int2 <- data$crp_intervention_2
 5. #Shapiro-Wilk test
 6. shap_plac <- shapiro.test(plac)
 7. shap_int1 <- shapiro.test(int1)
 8. shap_int2 <- shapiro.test(int2)
 9. print(shap_plac)
10. print(shap_int1)
11. print(shap_int2)
```

```
Output:
> print(shap_plac)
 Shapiro-Wilk normality test
data: plac
W = 0.84466, p-value = 0.0004803
> print(shap_int1)
```

```
Shapiro-Wilk normality test
data: int1
W = 0.74996, p-value = 9.092e-06
> print(shap_int2)
Shapiro-Wilk normality test
data: int2
W = 0.71482, p-value = 2.583e-06
```

With the result of the Shapiro-Wilk test, we can say that none of the samples is normally distributed. So we cannot use the ANOVA test, which assumes the samples are normally distributed. Instead, we use the Kruskal-Wallis H test to detect whether there is a significant difference among the three samples.

```
1. # reform data
2. crp_data <- data.frame(
3.   group = factor(rep(c("Placebo", "Intervention_1", "Intervention_2"), each = 30)),
4.   crp = c(plac, int1, int2)
5. )
6.
7. # Kruskal-Wallis H test
8. kruskal.test(crp ~ group, data = crp_data)
```

```
Output:
Kruskal-Wallis rank sum test
data: crp by group
Kruskal-Wallis chi-squared = 9.9986, df = 2, p-value = 0.006743
```

Because the p-value is lower than 0.05, it indicates a significant difference among the three samples. Hence, we performed pairwise comparisons using the Mann-Whitney U test to identify which groups differ significantly.

```
1. pairwise.wilcox.test(crp_data$crp, crp_data$group, p.adjust.method = "bonferroni")
```

```
Output:
Pairwise comparisons using Wilcoxon rank sum exact test

data: crp_data$crp and crp_data$group

                 Intervention_1     Intervention_2
Intervention_2   0.7254                 -
Placebo          0.2718              0.0026

P value adjustment method: bonferroni
```

We analyzed the results and concluded that there is a significant difference between the Placebo and Intervention_2 groups (p = 0.0026). No significant differences were observed between the Placebo and Intervention_1 groups, nor between the Intervention_1 and Intervention_2 groups after applying

the Bonferroni correction. Furthermore, since the mean CRP value in the Intervention_2 group is lower than in the Placebo group, we conclude that the drug KTH002 has a measurable positive effect by reducing inflammation.

```
> mean(plac)
[1] 72.5312
> mean(int2)
[1] 24.84123
```

**Task 4**

First, we can draw the basic DV-TIME curve to observe the changes in drug concentration over time. In addition, we can draw a boxplot figure of different patients, as shown below.

```
1. # read file
2. data <- read.csv("conctimedata_reduced.csv")
3. # extract essential data
4. data_subset <- data[, c("ID", "TIME", "DV")]
5. data_subset$ID <- as.factor(data_subset$ID)
6. data_subset$ID <- factor(data_subset$ID, labels = paste0("Patient", 1:12))
7. patient_colors <- setNames(rainbow(12), paste0("Patient", 1:12))
8. # draw DV-TIME curve
9. p <- ggplot(data_subset, aes(x = TIME, y = DV, group = ID, color = ID)) +
10.   geom_line(size = 0.5) +
11.   labs(title = "Concentration-Time Profiles for 12 Patients",
12.        x = "Time (hours)", y = "Concentration (mg/L)") +
13.   scale_color_manual(values = patient_colors,
14.              labels = paste0("Patient", 1:12)) +
15.   theme_minimal() +
16.   theme(legend.title = element_blank(),
17.        plot.title = element_text(size = 12, face="bold", hjust = 0.5))
18. # save and adjust windows width & height
19. ggsave("concentration_profiles.jpg", plot = p,
20.        width = 1600/300, height = 1200/300, units = "in", dpi = 300)
21. print(p)
22. # draw boxplot
23. p_box <- ggplot(data_subset, aes(x = ID, y = DV, fill = ID)) +
24.   geom_boxplot() +
25.   labs(title = "Concentration Distribution by Patient",
26.        x = "Patient ID",
27.        y = "Concentration (mg/L)") +
28.   theme_minimal() +
29.   theme(legend.position = "none",
30.        plot.title = element_text(size = 12, face="bold", hjust = 0.5))
31. # save and adjust windows width & height
32. ggsave("concentration_boxplot.jpg", plot = p_box,
```

```
33.         width = 1600/300, height = 1200/300, units = "in", dpi = 300)
34. print(p_box)
```
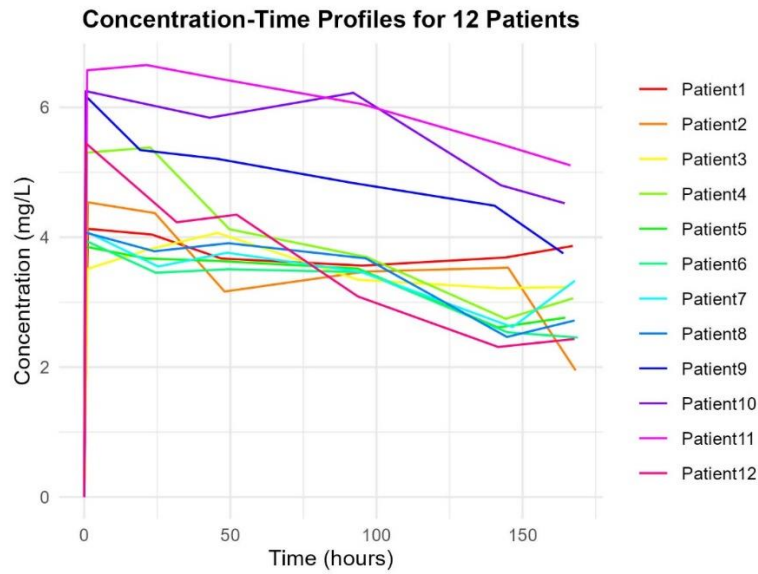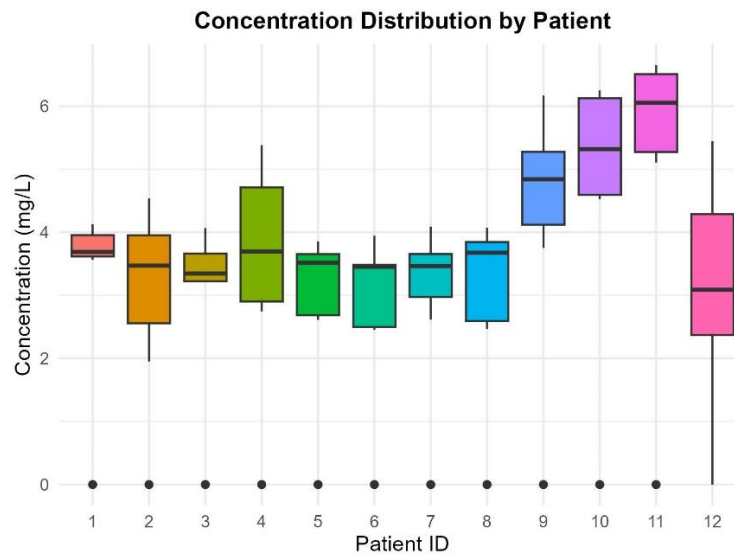


Figure4.1



Figure4.2

For the summary of the above curve, we need to introduce a concept, which is the Area Under Curve (AUC). It is the area under a known curve, i.e., integrating a particular interval. To achieve integration, we need first to make a cubic spline interpolation of the DV-TIME data points to smooth the curve. In the process, (0,0) points are discarded to reduce the sharp fluctuations of the curve. The cubic spline interpolation $S(x)$ satisfies,

$$S_i(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3$$

The condition at its nodes satisfies,

$$S_i(x_i) = y_i \quad S_i(x_{i+1}) = y_{i+1}$$
$$S_i(x_{i+1}) = S_{i+1}(x_{i+1})$$
$$S_i'(x_{i+1}) = S_{i+1}'(x_{i+1})$$

$$S_i''(x_{i+1}) = S_{i+1}''(x_{i+1})$$

```r
1. # Cubic Spline Interpolation
2. # create a new container to save data
3. interpolated_data <- data.frame(TIME = numeric(), DV = numeric(), ID = character())
4. # Cubic Spline Interpolation
5. for (id in levels(data_subset$ID)) {
6.    subset_data <- data_subset[data_subset$ID == id, ]
7.    subset_data <- subset_data[!(subset_data$TIME == 0 & subset_data$DV == 0), ]
8.    interpolated <- spline(subset_data$TIME, subset_data$DV, xout =
seq(min(subset_data$TIME), max(subset_data$TIME), length.out = 100))
9.    interpolated_data <- rbind(interpolated_data, data.frame(TIME = interpolated$x, DV =
interpolated$y, ID = id))
10. }
11. interpolated_data$ID <- factor(interpolated_data$ID, levels = levels(data_subset$ID))
12. # draw a new curve
13. p_interpolated <- ggplot(interpolated_data, aes(x = TIME, y = DV, group = ID, color =
factor(ID))) +
14.    geom_line(size = 0.5) +
15.    labs(title = "Interpolated Concentration-Time Profiles for 12 Patients",
16.        x = "Time (hours)", y = "Concentration (mg/L)") +
17.    scale_color_manual(values = patient_colors,
18.              labels = paste0("Patient", 1:12)) +
19.    theme_minimal() +
20.    theme(legend.title = element_blank(),
21.        plot.title = element_text(size = 12, face="bold", hjust = 0.5)) +
22.    coord_cartesian(ylim = c(0, NA))
23.
24. ggsave("interpolated_concentration_profiles.jpg", plot = p_interpolated,
25.        width = 1600/300, height = 1200/300, units = "in", dpi = 300)
26. print(p_interpolated)
```

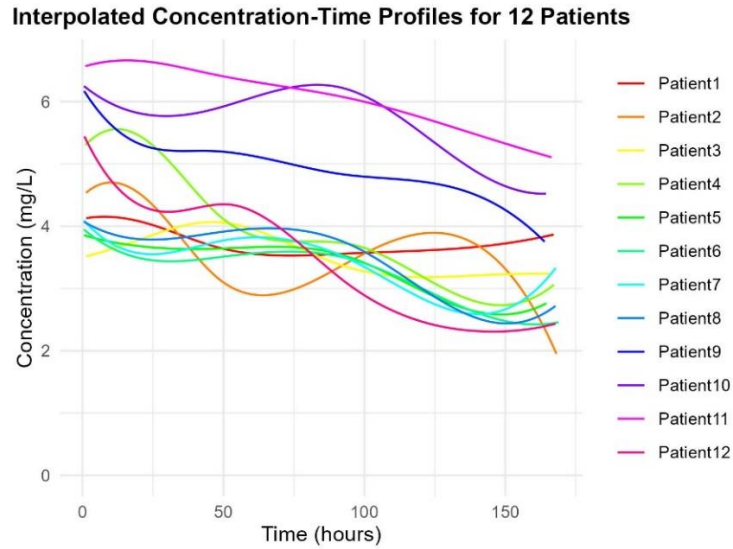**Interpolated Concentration-Time Profiles for 12 Patients**

Figure4.3

The curve in Figure 4.3 is smoother than in Figure 4.1. Then, perform integration over $x \in [1, 166]$ to obtain the result, and restructure to obtain a new dataset as shown below.

```
1. # initialize the container for storing integral results
2. integrated_results <- data.frame(ID = character(), Integral = numeric(), stringsAsFactors =
FALSE)
3. lower_bound <- 1
4. upper_bound <- 166
5. # Integral after interpolation
6. for (id in levels(interpolated_data$ID)) {
7.   subset_data <- interpolated_data[interpolated_data$ID == id, ]
8.   if (nrow(subset_data) > 1) {
9.     # create cubic spline interpolation function
10.    spline_function <- splinefun(subset_data$TIME, subset_data$DV)
11.    # perform the integral over the interval [1, 166]
12.    integral_value <- integrate(spline_function, lower = lower_bound, upper =
upper_bound)$value
13.    integrated_results <- rbind(integrated_results, data.frame(ID = id, Integral =
integral_value))
14.  } else {
15.    message(paste("ID", id, "error"))
16.  }
17. }
18. print(integrated_results)
```

Table4.1

| ID | AUC | WGT | BSA | AGE | HGT | DOSE | GFR |
|----|--------|-----|------|-----|-----|------|-----|
| 1 | 613.77 | 56 | 1.58 | 45 | 162 | 50 | 91 |
| 2 | 589.45 | 63 | 1.50 | 40 | 149 | 50 | 89 |
| 3 | 584.29 | 70 | 1.73 | 44 | 160 | 50 | 93 |

| 4 | 639.70 | 75 | 1.80 | 59 | 162 | 100 | 97 |
| 5 | 549.35 | 75 | 1.70 | 63 | 154 | 100 | 99 |
| 6 | 535.46 | 65 | 1.60 | 64 | 157 | 200 | 106 |
| 7 | 552.78 | 53 | 1.50 | 42 | 157 | 200 | 107 |
| 8 | 568.55 | 60 | 1.60 | 42 | 150 | 200 | 104 |
| 9 | 808.87 | 59 | 1.56 | 58 | 145 | 400 | 115 |
| 10 | 936.16 | 75 | 1.74 | 59 | 155 | 400 | 117 |
| 11 | 1000.14 | 63 | 1.65 | 34 | 153 | 400 | 111 |
| 12 | 567.22 | 67 | 1.70 | 46 | 160 | 100 | 100 |

Furthermore, we used multiple linear regression to verify the relationship between AUC and the six factors, and obtained the following results,

```
1. # read file
2. data <- read.csv("data_task4.csv", sep = ",")
3. # multiple linear regression
4. model <- lm(AUC ~ WGT + BSA + AGE + HGT + DOSE + GFR, data = data)
5. summary(model)
```

```
Output:
Residuals:
1    2    3    4    5    6    7    8    9    10   11   12
26.962 5.325 -16.329 -11.732 17.671 -58.328 10.616 -29.576 20.290 28.261 -20.128 26.968
Coefficients:
            Estimate Std.   Error   t value   Pr(>|t|)
 (Intercept) 1204.1046  541.9594   2.222 0.076958 .
WGT            5.9602    4.6517   1.281 0.256291
BSA          -34.0068  379.7916  -0.090 0.932128
AGE           -1.8867    1.7898  -1.054 0.340052
HGT            8.1983    4.1213   1.989 0.103347
DOSE           2.7067    0.3586   7.549 0.000646 ***
GFR          -25.0099    5.2831  -4.734 0.005178 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 40.51 on 5 degrees of freedomMultiple
R-squared: 0.9711,   Adjusted R-squared:  0.9363
F-statistic: 27.97 on 6 and 5 DF,  p-value: 0.001076
```

The regression analysis reveals that both DOSE and GFR are statistically significant predictors of AUC, with p-values of 0.000646 and 0.005178, respectively. This indicates a strong influence of both DOSE and GFR on AUC. Specifically, the positive coefficient of 2.71 for DOSE suggests that an increase in the dose leads to an increase in AUC. Conversely, the negative coefficient of -25.01 for GFR implies that an increase in GFR is associated with a reduction in AUC. Other variables, such as WGT, BSA, AGE, and HGT, did not show significant effects on AUC, as evidenced by their higher p-values. The model accounts for approximately 97.11% of the variance in AUC, with an adjusted R-squared value of 93.63%, indicating a robust overall fit. Thus, the analysis highlights that DOSE and GFR are the primary factors affecting AUC, whereas the other variables play a less

significant role.

To evaluate the appropriateness of the regression model, we conducted diagnostic analyses and visualized the results, as presented in Figure 4.4. The following sections provide a detailed examination of the diagnostic plots.
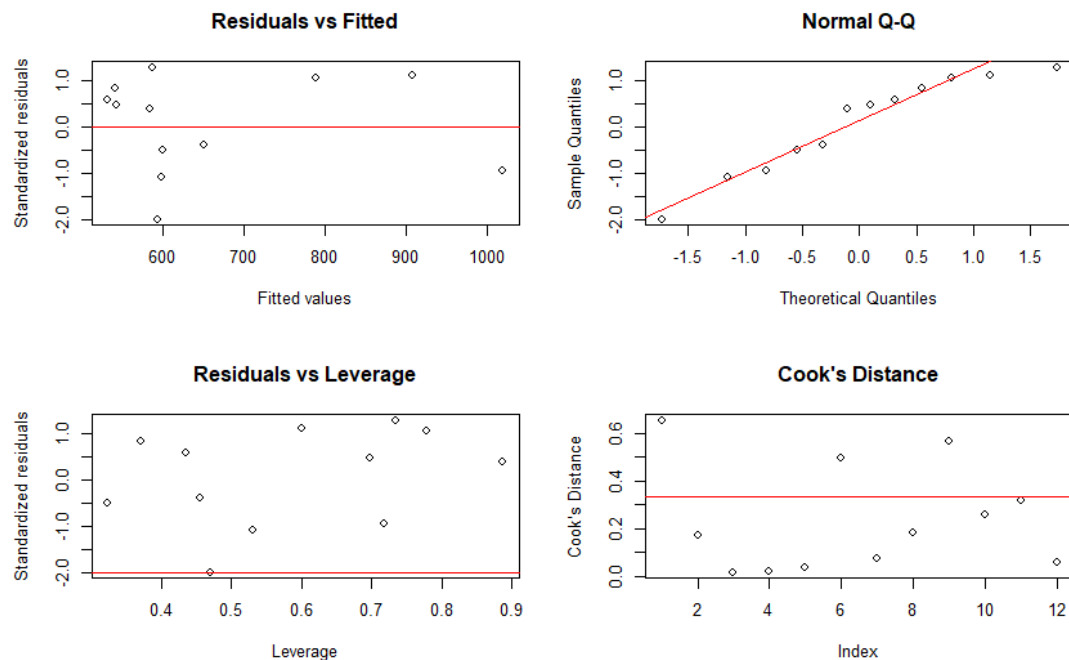


Figure 4.4

1. Residuals vs Fitted Plot

The plot does not display any significant curvature or recognizable patterns, suggesting that the assumption of linearity is reasonably upheld. The residuals are randomly scattered across the range of fitted values, indicating that the model sufficiently captures the relationship between the predictors and the response variable.

The residuals are scattered around the horizontal line at zero, which indicates that the model does not show significant signs of non-linearity. However, we can see a slight deviation at the higher fitted values (above 900), suggesting a potential issue. But the residuals range between approximately -1.5 and 1.2, meaning heteroscedasticity is not a major concern in this model.

2. Normal Q-Q Plot

A Q-Q plot (quantile-quantile plot) is a graphical tool used to compare the distribution of a dataset with a theoretical distribution, usually a normal distribution. It plots the quantiles of the sample data against the quantiles of the theoretical distribution. If the data follows the theoretical distribution closely, the points will approximately form a straight line along the 45-degree reference line. Deviations from this line indicate departures from the theoretical distribution, such as skewness or heavy tails.

The points mostly fall along the 45-degree reference line, which indicates that the residuals follow

an approximately normal distribution. There is a slight deviation at both ends, especially below -1.5 and above 1.5 on the theoretical quantile axis. These deviations suggest that some of the extreme values are slightly off from normality, but overall the residuals are close to normal.

3. Residuals vs Leverage Plot

Leverage[1] measures how much influence a data point has on the regression model, particularly in terms of how far the point is from the center of the predictor variables. Leverage values typically range between 0 and 1, high-leverage points have no neighboring points in $R^p$ space, where $p$ is the number of independent variables in a regression model. This makes the fitted model likely to pass close to a high leverage observation[2]. So we can consider that higher leverage means bigger influence.

The plot shows that most points exhibit relatively low leverage, and no points appear to be extreme outliers and most of the residuals fluctuate around 0. This suggests that no individual observation is disproportionately influencing the model's results. The leverage values range from approximately 0.45 to 0.9. While most residuals are between -1.5 and 1, a few points have high leverage values, especially those beyond 0.8. This suggests that these points have a higher influence on the model. However, the absence of a distinct pattern and the moderate residuals indicate that there are no extreme outliers in the data.

4. Cook's distance Plot

Cook's distance[3] is used to evaluate the influence of a data point when performing a least-squares regression analysis. One observation's Cook's distance is defined as the sum of all the changes in the model when this observation is removed from the model[4]. If a data point has a large Cook's distance, it indicates that removing the point would significantly alter the regression model's estimates, making it an influential point.

In this plot, just a few of the observations have Cook's distance values greater than the threshold, which is set at 4 divided by the sample size, implying that no highly influential points are present in the dataset. Consequently, the regression model is not unduly affected by outliers or influential data points. In addition, none of the points exceed the critical threshold of 0.5 (red line), suggesting no individual observation has an undue influence on the regression model. The Cook's distance values range from 0 to approximately 0.5, which is below the threshold of 1 typically used to flag influential observations.

Overall, the regression model appears to be well-specified, and the assumptions of linear regression are reasonably satisfied. Further adjustments to the model may be unnecessary unless future data indicates otherwise.

Additionally, we can use a correlation matrix to explore the linear relationships between pairs of variables. A correlation matrix is a symmetric matrix used to display the correlation coefficients between multiple variables. The correlation coefficient measures the linear relationship between two variables, typically ranging from −1 to 1. Each element in the correlation matrix represents the correlation between two variables, with the Pearson correlation coefficient being a commonly used

measure. The closer the absolute value is to 1, the stronger the linear relationship. As shown in the figure 4.4.

Perhaps we can use more flexible and intelligent regression methods, such as deep learning (neural networks). Nonlinear fitting with the sigmoid function might yield more accurate weights and prediction precision. However, due to the small sample size in this case, the deep learning packages provided by PyTorch would easily overfit, causing the fitting results to lose practical significance. In this task, we attempted this method, but we abandoned the results due to overfitting eventually.

**References**

[1] 'Leverage (statistics)', *Wikipedia*. Oct. 04, 2024. Accessed: Oct. 11, 2024. [Online]. Available: https://en.wikipedia.org/w/index.php?title=Leverage_(statistics)&oldid=1249350472#cite_note-1

[2] H. Joe, 'Everitt BS 2002: The Cambridge dictionary of statistics, Cambridge: Cambridge University Press. 420 pp.\pounds 30.00 (US $50.00)(HB). ISBN 05218 1099 X.' 2004.

[3] R. D. Cook, 'Detection of Influential Observation in Linear Regression', *Technometrics*, vol. 19, no. 1, pp. 15–18, Feb. 1977, doi: 10.1080/00401706.1977.10489493.

[4] 'Cook's Distance - MATLAB & Simulink - MathWorks 中国'. Accessed: Oct. 11, 2024. [Online]. Available: https://ww2.mathworks.cn/help/stats/cooks-distance.html