# STAT 6338
# Advanced Statistical Methods II
# Homework 2

**Name:** Harshul Shah
**NetID:** hxs230024

**Ans 1 16.5)**
(a) Representation of the ANOVA Model (Figure 16.2 Format)
The ANOVA model follows:

$$Y_{ij} = \mu_i + \epsilon_{ij}, \qquad \epsilon_{ij} \sim N(0, \sigma^2)$$

$\mu_i$ represents the mean length of hospital stay for income group i.
sigma^2 = 2.8 represents the common variance across groups.
The representation aligns with Figure 16.2, where different income groups (L1, L2, L3, L4) have different means (5.1, 6.3, 7.9, 9.5) but follow a normal distribution with overlapping variances.

(b) Expected Values of MSTR and MSE
Computed MSE (Mean Square Error): ≈ 2.87
Computed MSTR (Mean Square Treatment): ≈ 47.2
Comparison:
$E\{MSE\} \approx 2.87$ remains close to $\sigma^2 = 2.8$, as expected.
$E\{MSTR\} \approx 47.2$ is significantly larger than MSE.
Implication:
Since MSTR >> MSE, the variation between income groups is much larger than the variation within groups. This suggests that income groups significantly impact hospital stay length rather than individual variations alone.

(c) Effect of Changing $\mu_2$ and $\mu_3$ to 5.6 and 9.0
New MSTR = 63.1 (increased from 47.2)
MSE remains at 2.87, consistent with error variance.
Why did MSTR increase?
- The range of means (max-min) is still 4.4 days (same as before).
- However, more separation in the middle means increases overall variation, leading to a larger MSTR.
- This emphasizes that distribution shape impacts MSTR, not just range.

Final Conclusion
- ANOVA confirms that income group significantly affects hospital stay length.
- Since MSTR >> MSE, income groups explain a large portion of the variance.
- Adjusting middle group means increases MSTR more than expected, showing sensitivity to mean distribution, not just range.

**Ans 2)**
(a) Fit the ANOVA Model
ANOVA Table p-value:
The output shows Pr > F <.0001 for the Group effect.
Reject H0?:
 Yes, we reject the null hypothesis

H0:$\mu1=\mu2=\mu3$
at the 5% level because the p-value is less than 0.05. This indicates that there are significant differences in Duration among the different groups.

(b) Boxplot and Equal Variance Assumption
- Boxplot Observation: The boxplot visually represents the distribution of Duration for each group. Differences in the medians and interquartile ranges can be observed.
- Equal Variance Assumption: Without a specific test for equal variance (like Levene's test) included in this output, we need to visually inspect the boxplot. If the boxes (interquartile ranges) have roughly similar sizes across the groups, then the assumption of equal variance might be reasonable. If there are highly unequal box sizes, it's a concern.

(c) Qualitatively describe the relationship between physical fitness status and duration of required physical therapy.
Qualitative Description: The relationship between physical fitness status (represented by the groups) and the duration of required physical therapy can be qualitatively described based on the group means.
Intercept (baseline): 24.00
- Group 1: Intercept + 14.00 = 38.00
- Group 2: Intercept + 8.00 = 32.00
- Group 3: Intercept + 0.00 = 24.00
  - The mean duration for Group 1 is 38 days.
    The mean duration for Group 2 is 32 days.
    The mean duration for Group 3 is 24 days.
From these estimations, it can be said that Group 1 requires **more** physical therapy.

(d) Problem 17.10 parts b, c, d, and f.

b. Estimate with a 99 percent confidence interval the mean number of days required in therapy for persons of average physical fitness.
**Ans** Average physical fitness corresponds to the Intercept.
Intercept Estimate: 24.00
Standard Error of Intercept: 1.81702705
Degrees of Freedom (Error): 21
t-critical value for 99% CI (df=21): 2.831
Margin of Error: 2.831 * 1.81702705 ≈ 5.144
99% Confidence Interval: 24.00 ± 5.144 = [18.856, 29.144]

c. Obtain confidence intervals for D1 = $\mu1$ - $\mu3$ and D2 = $\mu2$ - $\mu3$: use the Bonferroni procedure with a 95 percent family confidence coefficient.
**Ans** We have two comparisons, so we adjust α to α/2 = 0.05/2 = 0.025. This makes it a 97.5% confidence interval, with alpha = 0.025.
t-critical value for 97.5% (df=21): 2.080
D1 = $\mu1$ - $\mu3$:
Estimate: 14.00 - 0.00 = 14.00
Standard Error: sqrt((2.4037)^2 + (0)^2) = 2.404
Margin of Error: 2.080 * 2.404 ≈ 4.999
95% Confidence Interval: 14.00 ± 4.999 = [9.001, 18.999]
D2 = $\mu2$ - $\mu3$:
Estimate: 8.00 - 0.00 = 8.00
Standard Error: sqrt((2.2983)^2 + (0)^2) = 2.2983
Margin of Error: 2.080 * 2.2983 ≈ 4.780
95% Confidence Interval: 8.00 ± 4.780 = [3.220, 12.780]

d. Would the Tukey procedure have been more efficient to use in part c above? Explain.
**Ans** Yes, Tukey would be more efficient. The output indicates Tukey's HSD test was conducted, and the confidence intervals directly address all pairwise comparisons. The Bonferroni correction is generally more conservative, especially with a small number of comparisons. The Tukey procedure gives you intervals specific to comparing all pairs of means, which is precisely what's asked for here.

f. Test for all pairs of factor level means whether or not they differ: use the Tukey procedure with $\alpha = 0.05$. Set up groups of factor levels whose means do not differ.
**Ans** The output provides the results of Tukey's Studentized Range (HSD) Test. All comparisons are marked with "***", indicating significant differences at the 0.05 level. Conclusion: All pairs of means differ significantly. There are no groups of factor levels whose means do not differ.

(e) Problem 17.15 parts a, b, c.

a. Estimate the contrast $L = (\mu_1 - \mu_2) - (\mu_2 - \mu_3)$ with a 99 percent confidence interval.
**Ans** Estimate for Contrast L: -22.0000000
Standard Error: 4.20279574
Degrees of Freedom: 21
t-critical value for 99% (df=21): 2.831
Margin of Error: 2.831 * 4.20279574 ≈ 11.90
99% Confidence Interval: $-22.00 \pm 11.90 = [-33.90, -10.10]$

b. Estimate the following comparisons using the Bonferroni procedure with a 95 percent family confidence coefficient:
D1 = $\mu_1 - \mu_2$
D2 = $\mu_1 - \mu_3$
D3 = $\mu_2 - \mu_3$
L1 = D1 - D3
**Ans** We have four comparisons, so we adjust $\alpha$ to $\alpha/4 = 0.05/4 = 0.0125$. This makes it a 98.75% confidence interval.
t-critical value for alpha = 0.0125 (df=21): 2.528
D1 = $\mu_1 - \mu_2$:
Estimate: 14.00 - 8.00 = 6.00
Standard Error = 3.32781196
Margin of Error: 2.528 * 3.32781196 ≈ 8.41
95% Confidence Interval: $6.00 \pm 8.41 = [-2.41, 14.41]$
D2 = $\mu_1 - \mu_3$:
Estimate: 14.00 - 0.00 = 14.00
Standard Error = 2.40370085
Margin of Error: 2.528 * 2.40370085 ≈ 6.077
95% Confidence Interval: $14.00 \pm 6.077 = [7.923, 20.077]$
D3 = $\mu_2 - \mu_3$:
Estimate: 8.00 - 0.00 = 8.00
Standard Error = 2.29837762
Margin of Error: 2.528 * 2.29837762 ≈ 5.809
95% Confidence Interval: $8.00 \pm 5.809 = [2.191, 13.809]$
L1 = D1 - D3:
Estimate: 6.00 - 8.00 = -2.00
Standard Error = 3.2756451
Margin of Error = 2.528 * 3.2756451 = 8.27
95% Confidence Interval = $-2.00 \pm 8.27$

c. Would the Scheffe procedure have been more efficient to use in part (b) than the Bonferroni procedure?

**Ans** It depends. With only 4 pre-planned comparisons, the Bonferroni method may provide tighter confidence intervals. However, if you were exploring other contrasts that weren't pre-planned, the Scheffé method would be more appropriate because it controls the family-wise error rate for all possible contrasts.


**Ans 3)**

(a) Fit the ANOVA Model

ANOVA Table p-value: The output shows $Pr > F < 0.0001$ for the Group effect, meaning that the p-value is less than 0.0001.

Reject H0?:

Yes, we reject the null hypothesis:

$$H0: \mu1 = \mu2 = \mu3 = \mu4 = \mu5$$

at the 5% level (and even at the 0.01% level) because the p-value ($<0.0001$) is much less than 0.05. This indicates that there are significant differences in Time Lapse among the different groups.

(b) Box-plot and Equal Variance Assumption

Box-plot Observation: While the visual box-plot isn't available, the group means and standard deviations are:
- Group 1: Mean = 24.55, Std Dev = 2.48
- Group 2: Mean = 22.55, Std Dev = 2.98
- Group 3: Mean = 11.75, Std Dev = 2.57
- Group 4: Mean = 14.80, Std Dev = 2.53
- Group 5: Mean = 30.10, Std Dev = 3.09

Levene's Test: The test for homogeneity of variance gives $p = 0.7836$. Since $0.7836 > 0.05$, we do not reject the null hypothesis that variances are equal across groups. Thus, the assumption of equal variance appears reasonable.

(c) Relationship Between Agent (Group) and Time Lapse

Qualitative Description:
- Group 5 has the highest mean Time Lapse (30.1).
- Group 3 has the lowest mean (11.75).
- Groups 1 and 2 have similar means (24.55 and 22.55).
- Group 4 has a mean (14.80) that falls between Groups 3 and (1&2).
- This confirms a clear relationship between agent groups and Time Lapse.

(d) Tukey's Test and Confidence Intervals

Tukey's Test Output: The test identifies significant differences between groups with $\alpha = 0.1$. The Minimum Significant Difference = 2.1654.

Pairwise Comparisons:

|Group 1 - Group 2| = |24.55 - 22.55| = 2.00 → Not significantly different

|Group 1 - Group 3| = |24.55 - 11.75| = 12.80 → Significantly different

Confidence Interval for Mean of Group 1:

Mean = 24.55

90% Confidence Interval: [23.59, 25.51]

Confidence Interval for $D = \mu2 - \mu1$:

Parameter Estimate = -2.00

Standard Error = 0.8673

t Value = -2.31, p-value = 0.0233
90% Confidence Interval: [-3.44, -0.56]
Since the interval does not include 0, this suggests that Group 2 has a significantly lower mean than Group 1.

(e) Contrast Estimation
Estimate $L=(\mu1+\mu2)/2-(\mu3+\mu4)/2$
Parameter Estimate = 10.275
Standard Error = 0.6133
90% Confidence Interval: [9.26, 11.30]
Since this interval does not contain 0, it confirms that the average of means for Groups 1 and 2 is significantly greater than that for Groups 3 and 4.

Scheffe Procedure for Multiple Comparisons:
$D1 = \mu1-\mu2$     (Estimate = 2.00)
$D2 = \mu3-\mu4$ (Estimate = -3.05)
$L1 = (\mu1+\mu2)/2-\mu5$ (Estimate = -6.55)
$L2 = (\mu3+\mu4)/2-\mu5$ (Estimate = -16.825)
$L3 = (\mu1+\mu2)/2-(\mu3+\mu4)/2$ (Estimate = 10.275)
Using F(alpha=0.1, df1=4, df2=95) = 2.07, the Scheffe confidence intervals can be calculated accordingly.

(f) Residual Analysis
Constant Variance:
Levene's test
$p = 0.7836 > 0.05$, indicating homoscedasticity (equal variance assumption is met).
Outliers:
Normality:
Histogram of Residuals: Ideally should be bell-shaped for normality assumption to hold.