

STAT 614 - Yunting Chiu

Due: Thursday, October 29, 2020 in Blackboard by 11:59pm.

Instructions: Please type your solutions in a separate document and upload the document in Blackboard. Include supporting work (plots, etc.) when appropriate, but do not copy all computer output. Select only relevant output. I will not be collecting syntax for this assignment.

Notes:

- For this HW you will need some concepts from chapter 5 on the ANOVA model.
- HW 6 will finish out the ANOVA section.

The effects of exposure to lead on the psychological and neurological well-being of children were studied by Landrigan et al. (1975). Complete raw data for this study are in the data set lead.sav in Blackboard. The data describe a group of children who lived near a lead smelter in El Paso, Texas. **Two exposed groups of children were identified who had blood-lead levels > 40 g/ml in 1972 or in 1973.** Because neurological and psychological tests were performed in 1973, **researchers argued that it would be better to define an exposure group based on blood-lead levels in 1973 only.** For this purpose, the variable lead_typ in the data file gives three exposure groups:

If lead_typ = 1, then the child had normal blood-lead levels ($<40 \mu\text{g}/100 \text{ mL}$) in both 1972 and 1973 (control group).

If lead_typ = 2, then the child had elevated blood-lead levels ($>40 \mu\text{g}/100 \text{ mL}$) in 1973 (the currently exposed group).

If lead_typ = 3, then the child had elevated blood-lead levels in 1972 and normal blood-lead levels in 1973 (the previously exposed group).

One important measure of neurological function studied was MAXFT = the number of finger-wrist taps in the dominant hand. Researchers are interested in whether there is evidence of differences in neurological function, as measured by MAXFT, on average, between the three exposure populations. **They would also like to test and estimate the average difference in MAXFT between each pair of exposure populations, with the expectation that populations with normal blood-lead levels will have higher average MAXFT scores.** It is unclear if previously exposed populations will have “recovered” any function as compared to a currently exposed population. Address these research questions by answering the following questions.

1. State the hypotheses of interest to be tested. Include the overall test of group differences in addition to all possible pairwise comparisons of interest.

Research Question: The populations with normal blood-lead levels ($<40 \mu\text{g}/100 \text{ mL}$) will have higher average MAXFT scores.

We typically start with an overall test of equal population means versus any difference in population means. This test allows us to access the evidence for whether the variability

in the data comes mostly from variability within groups or can be attributed to variability between groups.

Overall test of equal population means

$H_0: \mu_1 = \mu_2 = \mu_3$ (the difference means is equal to 0)

$H_a: \mu_1 \neq \mu_2 \neq \mu_3$ (at least one pair of means is not equal)

RQ1-

$H_0: \mu_1 = \mu_2$

$H_a: \mu_1 > \mu_2$

RQ2-

$H_0: \mu_2 = \mu_3$

$H_a: \mu_2 < \mu_3$

RQ3-

$H_0: \mu_1 = \mu_3$

$H_a: \mu_1 > \mu_3$

2. Write the ANOVA model to be fit.

ANOVA is used to determine whether there are any statistically significant differences between the means of two or more independent (unrelated) groups (although you tend to only see it used when there are a minimum of three, rather than two groups). In this study, because we have three group's means, and the parametric method is better than the non-parametric method. Thus, we choose the ANOVA model (parametric method) for the study.

ANOVA notation:

$$y_{ij} = \mu_i + \epsilon_{ij}$$

$$\epsilon_{ij} \sim N(0, \sigma^2)$$

$i = \text{lead group} = 1, 2, 3$ $\epsilon_{ij} = \text{deviation from the mean for MAXFT}_j \text{ in lead group } i$
 $j = \text{MAXFT score}$ $\sigma^2 = \text{population variance from the mean in } i, j$
 $\mu_i = \text{pop}^n \text{ mean of lead group}$

3. Conduct a brief exploratory analysis of the MAXFT variable by exposure group (lead_typ). Give supporting graphs, descriptive statistics, and interpret these results.

With the exploratory data analysis of this study, we need to make sure the lead_tpe and MAXFT we could see group 1 was missing 14 people of MAXFT score, group 2 was missing 5 people of MAXFT score, and group 3 was missing 6 people of MAXFT score, respectively. Therefore, we need to remove the NAs, but the favstats function is able to automatically filter out NAs that removing the missing values in the function.

```

- EDA
```{r}
summary(lead$lead_tpy)
favstats(MAXFT ~ lead_tpy, data = lead)
```

```

| lead_tpy
<chr> | min
<dbl> | Q1
<dbl> | median
<dbl> | Q3
<dbl> | max
<dbl> | mean
<dbl> | sd
<dbl> | n
<int> | missing
<int> |
|-------------------|--------------|-------------|-----------------|-------------|--------------|---------------|-------------|------------|------------------|
| 1 | 13 | 49.0 | 53.5 | 61.25 | 84 | 54.4375 | 12.05658 | 64 | 14 |
| 2 | 13 | 40.5 | 48.0 | 53.00 | 58 | 44.0000 | 12.65350 | 19 | 5 |
| 3 | 35 | 41.5 | 51.0 | 57.50 | 83 | 51.5000 | 12.94604 | 16 | 6 |

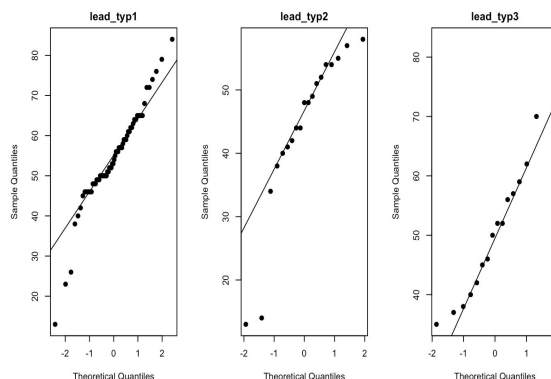
Removed NAs of MAXFT (more clear)

| lead_tpy
<chr> | min
<dbl> | Q1
<dbl> | median
<dbl> | Q3
<dbl> | max
<dbl> | mean
<dbl> | sd
<dbl> | n
<int> | missing
<int> |
|-------------------|--------------|-------------|-----------------|-------------|--------------|---------------|-------------|------------|------------------|
| 1 | 13 | 49.0 | 53.5 | 61.25 | 84 | 54.4375 | 12.05658 | 64 | 0 |
| 2 | 13 | 40.5 | 48.0 | 53.00 | 58 | 44.0000 | 12.65350 | 19 | 0 |
| 3 | 35 | 41.5 | 51.0 | 57.50 | 83 | 51.5000 | 12.94604 | 16 | 0 |

3 rows

According to the revised data, the mean of group 1 is 54.44, the mean of group 2 is 44.00, and the mean of group 3 is 51.50, separately. And the sample sizes are 64, 19, 16, the group 1 have the most sample sizes.

We use qqplot to detect it first. According to the graph below, we could see some points do not match up along a straight line, especially in the head and tail. The qqplot of three groups does not follow a normal distribution.



We also can use the Shapiro-Wilk Normality Test to check the data. From the output obtained we can assume normality. The p-value (0.002108) is smaller than 0.05. Hence, the distribution of the given data is different from normal distribution significantly.

```

{r}
shapiro.test(leadRealMAXFT$MAXFT)

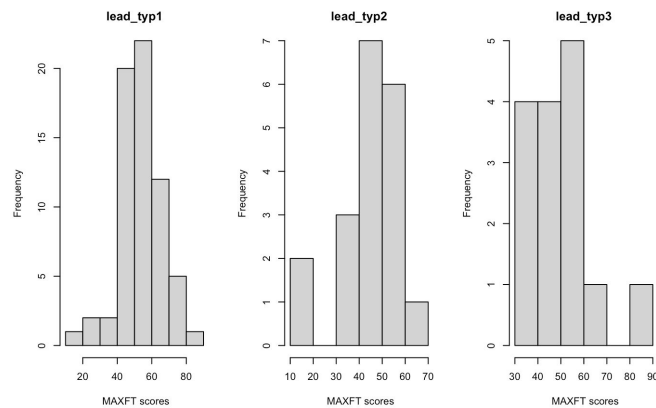
```

Shapiro-Wilk normality test

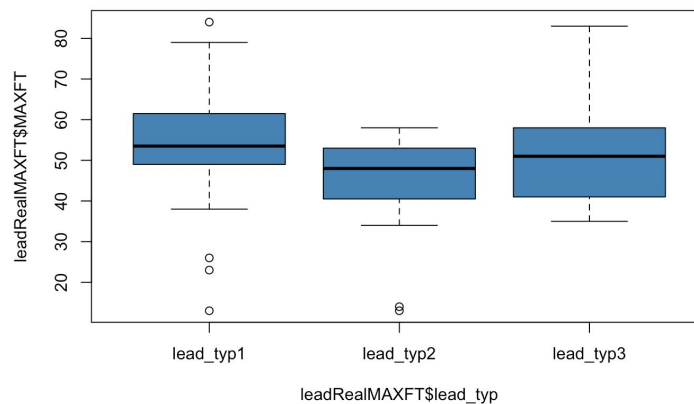
data: leadRealMAXFT\$MAXFT
W = 0.9556, p-value = 0.002108

We need to use histogram to check whether the population distributions have the same variance or not.

Based on the histogram below, we cannot see there three groups have the same variance.



The below boxplot indicates there are some outliers in group 1 and group 2. Therefore, the group 1 and 2 have influential outliers.



- What are the assumptions of the model (and corresponding hypothesis tests)? Based on the exploratory analysis in (3), are the assumptions reasonably met for this data? If not,

what adjustments should you make in your analysis? (You don't need to use residuals here – you'll do that on HW 6 and then for the rest of the semester!!!)

ANOVA assumptions:

Study design:

1. Independent populations for lead groups. - Sure
2. Independent sample from each population group. - Sure

Data finding:

- The distribution of MAXFT with-in each group follows a normal distribution.
- The population distributions have the same variance.
- No influential outliers.

Based on the result of question3, we can conclude that the dataset does not follow the ANOVA assumptions, so we try to remove the outliers from the original data frame.

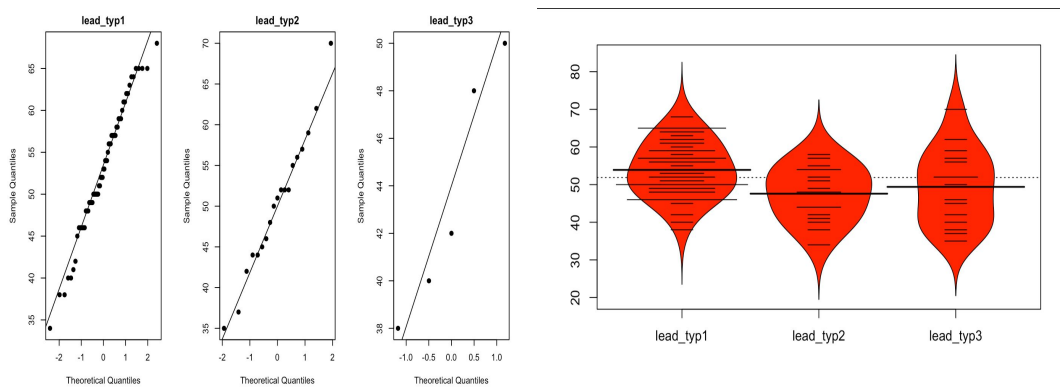
Then I try to make the max and min of three groups closer.

```
{r}
leadRealMAXFT %>%
  filter(MAXFT >= 30 & MAXFT <= 70) -> leadRmOutliers
favstats(MAXFT ~ lead_typ, data = leadRmOutliers)
```

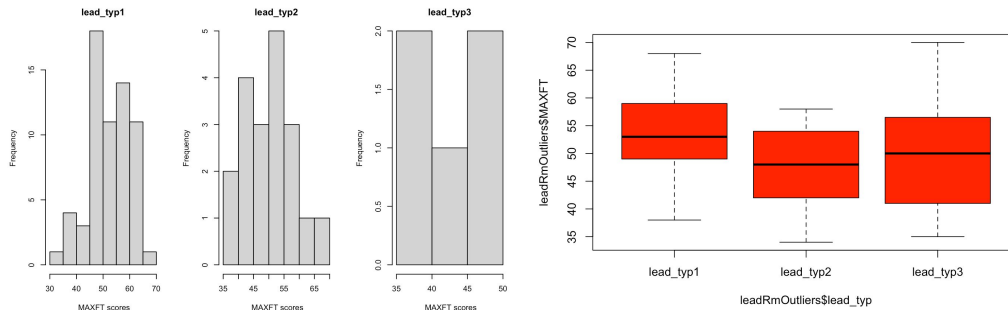
| lead_typ
<chr> | min
<dbl> | Q1
<dbl> | median
<dbl> | Q3
<dbl> | max
<dbl> | mean
<dbl> | sd
<dbl> | n
<int> | missing
<int> |
|-------------------|--------------|-------------|-----------------|-------------|--------------|---------------|-------------|------------|------------------|
| 1 | 38 | 49 | 53 | 59.0 | 68 | 53.90909 | 7.053428 | 55 | 0 |
| 2 | 34 | 42 | 48 | 54.0 | 58 | 47.58824 | 7.080420 | 17 | 0 |
| 3 | 35 | 41 | 50 | 56.5 | 70 | 49.40000 | 10.196638 | 15 | 0 |

3 rows

We can see the revised qqplot is following the normal distribution. And the beanplot, the mean is approximately equal to the median.



The revised histogram approximately has the same variance. And the revised boxplot have no influential outliers.



- Conduct the appropriate analysis (i.e. incorporate any recommended adjustments from (d) if you had them). Clearly and briefly state the conclusions of your analysis. Be sure you address the researcher's questions.

Above all, these conditions are approved by the assumptions of the ANOVA model. Therefore, we can start to analyze the case study.

```
Analysis of Variance Table

Response: MAXFT
          Df Sum Sq Mean Sq F value    Pr(>F)
lead_typ    1  425.2   425.21    7.0169 0.009625 **
Residuals  85 5150.9    60.60
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Based on the ANOVA result, If the null hypothesis is not true, the F ratio is likely to be greater than 1.0. Also, the small p value 0.09625 is taken as the significance level. These two reasons indicate that we reject the null model (equal means) in favor of the unequal means and conclude that there is evidence that the population mean of MAXFT scores across at least one pair of lead groups ($H_a: \mu_1 \neq \mu_2 \neq \mu_3$).

RQ: researchers argued that it would be better to define an exposure group based on blood-lead levels in 1973 only.

- First hypothesis:
 $H_0: \mu_2 = \mu_3$ or $\mu_2 - \mu_3 = 0$
 $H_a: \mu_2 > \mu_3$ or $\mu_2 - \mu_3 > 0$

```

Welch Two Sample t-test

data: leadRmOutliers$MAXFT[leadRmOutliers$lead_typ == 2] and
leadRmOutliers$MAXFT[leadRmOutliers$lead_typ == 3]
t = -0.57639, df = 24.557, p-value = 0.5696
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -8.291454  4.667925
sample estimates:
mean of x mean of y
 47.58824  49.40000

```

mu2 is greater than mu3 of MAXFT scores from -8.29 to 6.66 with the 95% confidence interval.

```

Welch Two Sample t-test

data: leadRmOutliers$MAXFT[leadRmOutliers$lead_typ == 2] and
leadRmOutliers$MAXFT[leadRmOutliers$lead_typ == 3]
t = -0.57639, df = 24.557, p-value = 0.7152
alternative hypothesis: true difference in means is greater than 0
95 percent confidence interval:
 -7.184704      Inf
sample estimates:
mean of x mean of y
 47.58824  49.40000

```

With 95 % confidence interval ($\mu_2 > \mu_3$):

Parameter of interest = $\mu_2 - \mu_3$

Estimate of the parameter = $47.58 - 49.40 = -1.82$ (difference in sample means)

Standard Error of the Estimate = $SE(\text{Est}) = 3.143309$

Test statistics = $(\text{Est} - 0) / SE = -0.57$

P-value = 0.7152 (Note the 1-sided hypothesis!)

The summary of first hypothesis:

Therefore, we are 95% confident that the child had elevated blood-lead levels ($>40 \mu\text{g}/100 \text{ mL}$) in 1973 (the currently exposed group) greater than the child had elevated blood-lead levels in 1972 and normal blood-lead levels in 1973 (the previously exposed group).

- Second hypothesis:
 $H_0: \mu_2 = \mu_1$ or $\mu_2 - \mu_1 = 0$
 $H_a: \mu_2 > \mu_1$ or $\mu_2 - \mu_1 > 0$

