

Assignment 6: GLMs week 1 (t-test and ANOVA)

Jared Wang

OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics on t-tests and ANOVAs.

Directions

1. Change “Student Name” on line 3 (above) with your name.
2. Work through the steps, **creating code and output** that fulfill each instruction.
3. Be sure to **answer the questions** in this assignment document.
4. When you have completed the assignment, **Knit** the text and code into a single PDF file.
5. After Knitting, submit the completed exercise (PDF file) to the dropbox in Sakai. Add your last name into the file name (e.g., “Salk_A06_GLMs_Week1.Rmd”) prior to submission.

The completed exercise is due on Tuesday, February 18 at 1:00 pm.

Set up your session

1. Check your working directory, load the **tidyverse**, **cowplot**, and **agricolae** packages, and import the NTL-LTER_Lake_Nutrients_PeterPaul_Processed.csv dataset.
2. Change the date column to a date format. Call up **head** of this column to verify.

```
#1
```

```
getwd()
```

```
## [1] "C:/Users/wangc/Box/Home Folder cw369/Private/Duke/Course/Spring 2020/Environmental_Data_Analyti
```

```
library(tidyverse)
```

```
library(cowplot)
```

```
library(agricolae)
```

```
library(FSA)
```

```
library(dunn.test)
```

```
library(ggthemes)
```

```
df.lk <-
```

```
  read.csv("../Data/Processed/NTL-LTER_Lake_Nutrients_PeterPaul_Processed.csv")
```

```
#2
```

```
df.lk <- mutate(df.lk, sampledate = as.Date(sampledate, "%Y-%m-%d"))
```

```
class(df.lk$sampledate)
```

```
## [1] "Date"
```

```
head(df.lk$sampledate)
```

```
## [1] "1991-05-20" "1991-05-20" "1991-05-20" "1991-05-20" "1991-05-20"
```

```
## [6] "1991-05-20"
```

Wrangle your data

3. Wrangle your dataset so that it contains only surface depths and only the years 1993-1996, inclusive. Set month as a factor.

```
df.lk0 <- df.lk %>%
  filter(depth == 0 & year4 %in% c(1993:1996)) %>%
  mutate(month = as.factor(month))
```

Analysis

Peter Lake was manipulated with additions of nitrogen and phosphorus over the years 1993-1996 in an effort to assess the impacts of eutrophication in lakes. You are tasked with finding out if nutrients are significantly higher in Peter Lake than Paul Lake, and if these potential differences in nutrients vary seasonally (use month as a factor to represent seasonality). Run two separate tests for TN and TP.

4. Which application of the GLM will you use (t-test, one-way ANOVA, two-way ANOVA with main effects, or two-way ANOVA with interaction effects)? Justify your choice.

Answer: Because influence of seasonality on water column nutrient levels is likely to be influenced by other factors of a lake (e.g. size & depth), I expect that the effects of seasonality and lake name on surface nutrient level are dependent on each other. Therefore, would use two-way ANOVA with interaction effects.

5. Run your test for TN. Include examination of groupings and consider interaction effects, if relevant.

6. Run your test for TP. Include examination of groupings and consider interaction effects, if relevant.

```
#5 - ANOVA for TP
#format as aov
aov.tp <- aov(data = df.lk0, tp_ug ~ month * lakename)
summary(aov.tp)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## month         4      671      168   1.623 0.1730
## lakename       1    10370    10370 100.283 <2e-16 ***
## month:lakename 4      1014       254   2.452 0.0496 *
## Residuals    119    12305       103
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 1 observation deleted due to missingness
```

```
#format as lm
lm.tp <- lm(data = df.lk0, tp_ug ~ month * lakename)
#post-hoc test for pairwise differences
TukeyHSD(aov.tp)
```

```
##      Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = tp_ug ~ month * lakename, data = df.lk0)
##
## $month
##              diff              lwr              upr              p adj
## 6-5  5.9146220  -3.234390  15.063634  0.3837749
## 7-5  7.9267363  -1.222276  17.075748  0.1224572
## 8-5  4.3748753  -4.706921  13.456671  0.6703911
## 9-5  3.8207521  -8.393804  16.035308  0.9085595
## 7-6  2.0121143  -4.721376   8.745605  0.9215444
## 8-6 -1.5397467  -8.181621   5.102128  0.9677800
## 9-6 -2.0938698 -12.621493   8.433754  0.9816312
## 8-7 -3.5518610 -10.193735   3.090013  0.5765788
```

```

## 9-7 -4.1059841 -14.633608 6.421639 0.8162959
## 9-8 -0.5541231 -11.023385 9.915139 0.9998946
##
## $lakename
##               diff      lwr      upr p adj
## Peter Lake-Paul Lake 17.91381 14.36807 21.45955 0
##
## $`month:lakename`
##               diff      lwr      upr      p adj
## 6:Paul Lake-5:Paul Lake -0.9178824 -16.4886641 14.652899 1.0000000
## 7:Paul Lake-5:Paul Lake -1.7271111 -17.1846493 13.730427 0.9999981
## 8:Paul Lake-5:Paul Lake -2.0872222 -17.5447604 13.370316 0.9999902
## 9:Paul Lake-5:Paul Lake -0.7380000 -20.5935673 19.117567 1.0000000
## 5:Peter Lake-5:Paul Lake 4.3135714 -13.9293175 22.556460 0.9989515
## 6:Peter Lake-5:Paul Lake 16.8838889 1.4263507 32.341427 0.0206973
## 7:Peter Lake-5:Paul Lake 22.9304706 7.3596889 38.501252 0.0002415
## 8:Peter Lake-5:Paul Lake 15.0200000 -0.3355071 30.375507 0.0607728
## 9:Peter Lake-5:Paul Lake 14.7452500 -6.4208558 35.911356 0.4316694
## 7:Paul Lake-6:Paul Lake -0.8092288 -11.8989312 10.280474 1.0000000
## 8:Paul Lake-6:Paul Lake -1.1693399 -12.2590423 9.920363 0.9999989
## 9:Paul Lake-6:Paul Lake 0.1798824 -16.5021309 16.861896 1.0000000
## 5:Peter Lake-6:Paul Lake 5.2314538 -9.4943403 19.957248 0.9787107
## 6:Peter Lake-6:Paul Lake 17.8017712 6.7120688 28.891474 0.0000401
## 7:Peter Lake-6:Paul Lake 23.8483529 12.6013419 35.095364 0.0000000
## 8:Peter Lake-6:Paul Lake 15.9378824 4.9908457 26.884919 0.0003006
## 9:Peter Lake-6:Paul Lake 15.6631324 -2.5591082 33.885373 0.1584032
## 8:Paul Lake-7:Paul Lake -0.3601111 -11.2902412 10.570019 1.0000000
## 9:Paul Lake-7:Paul Lake 0.9891111 -15.5872518 17.565474 1.0000000
## 5:Peter Lake-7:Paul Lake 6.0406825 -8.5653181 20.646683 0.9437275
## 6:Peter Lake-7:Paul Lake 18.6110000 7.6808700 29.541130 0.0000101
## 7:Peter Lake-7:Paul Lake 24.6575817 13.5678793 35.747284 0.0000000
## 8:Peter Lake-7:Paul Lake 16.7471111 5.9617574 27.532465 0.0000827
## 9:Peter Lake-7:Paul Lake 16.4723611 -1.6532090 34.597931 0.1087387
## 9:Paul Lake-8:Paul Lake 1.3492222 -15.2271407 17.925585 0.9999999
## 5:Peter Lake-8:Paul Lake 6.4007937 -8.2052070 21.006794 0.9208652
## 6:Peter Lake-8:Paul Lake 18.9711111 8.0409811 29.901241 0.0000062
## 7:Peter Lake-8:Paul Lake 25.0176928 13.9279904 36.107395 0.0000000
## 8:Peter Lake-8:Paul Lake 17.1072222 6.3218685 27.892576 0.0000523
## 9:Peter Lake-8:Paul Lake 16.8324722 -1.2930979 34.958042 0.0926020
## 5:Peter Lake-9:Paul Lake 5.0515714 -14.1485150 24.251658 0.9975850
## 6:Peter Lake-9:Paul Lake 17.6218889 1.0455259 34.198252 0.0276305
## 7:Peter Lake-9:Paul Lake 23.6684706 6.9864574 40.350484 0.0004851
## 8:Peter Lake-9:Paul Lake 15.7580000 -0.7232597 32.239260 0.0735733
## 9:Peter Lake-9:Paul Lake 15.4832500 -6.5132124 37.479712 0.4163366
## 6:Peter Lake-5:Peter Lake 12.5703175 -2.0356832 27.176318 0.1571717
## 7:Peter Lake-5:Peter Lake 18.6168992 3.8911050 33.342693 0.0032014
## 8:Peter Lake-5:Peter Lake 10.7064286 -3.7915495 25.204407 0.3464892
## 9:Peter Lake-5:Peter Lake 10.4316786 -10.1207861 30.984143 0.8273658
## 7:Peter Lake-6:Peter Lake 6.0465817 -5.0431207 17.136284 0.7595330
## 8:Peter Lake-6:Peter Lake -1.8638889 -12.6492426 8.921465 0.9999197
## 9:Peter Lake-6:Peter Lake -2.1386389 -20.2642090 15.986931 0.9999970
## 8:Peter Lake-7:Peter Lake -7.9104706 -18.8575073 3.036566 0.3778093
## 9:Peter Lake-7:Peter Lake -8.1852206 -26.4074611 10.037020 0.9089776
## 9:Peter Lake-8:Peter Lake -0.2747500 -18.3133864 17.763886 1.0000000

```

```
tp.int <- with(df.lk0, interaction(month, lakename))
aov.tp.int <- aov(data = df.lk0, tp_ug ~ tp.int)
tp.groups <- HSD.test(aov.tp.int, trt = "tp.int", group = TRUE)
tp.groups
```

```
## $statistics
##      MSerror Df      Mean      CV
##    103.4055 119 19.07347 53.3141
##
## $parameters
##      test name.t ntr StudentizedRange alpha
##    Tukey tp.int  10          4.560262  0.05
##
## $means
##              tp_ug      std  r   Min   Max   Q25   Q50   Q75
## 5.Paul Lake  11.474000  3.928545  6  7.001 17.090  8.1395 11.8885 13.53675
## 5.Peter Lake 15.787571  2.719954  7 10.887 18.922 14.8915 15.5730 17.67400
## 6.Paul Lake  10.556118  4.416821 17  1.222 16.697  7.4430 10.6050 13.94600
## 6.Peter Lake 28.357889 15.588507 18 10.974 53.388 14.7790 24.6840 41.13000
## 7.Paul Lake   9.746889  3.525120 18  4.501 21.763  7.8065  9.1555 10.65700
## 7.Peter Lake 34.404471 18.285568 17 19.149 66.893 21.6640 24.2070 50.54900
## 8.Paul Lake   9.386778  1.478062 18  5.879 11.542  8.4495  9.6090 10.45050
## 8.Peter Lake 26.494000  9.829596 19 14.551 49.757 21.2425 23.2250 27.99350
## 9.Paul Lake  10.736000  3.615978  5  6.592 16.281  8.9440 10.1920 11.67100
## 9.Peter Lake 26.219250 10.814803  4 16.281 41.145 19.6845 23.7255 30.26025
##
## $comparison
## NULL
##
## $groups
##              tp_ug groups
## 7.Peter Lake 34.404471    a
## 6.Peter Lake 28.357889   ab
## 8.Peter Lake 26.494000  abc
## 9.Peter Lake 26.219250 abcd
## 5.Peter Lake 15.787571  bcd
## 5.Paul Lake  11.474000   cd
## 9.Paul Lake  10.736000   cd
## 6.Paul Lake  10.556118    d
## 7.Paul Lake   9.746889    d
## 8.Paul Lake   9.386778    d
##
## attr(,"class")
## [1] "group"
```

```
#6 - ANOVA for TN
aov.tn <- aov(data = df.lk0, tn_ug ~ month * lakename)
summary(aov.tn)
```

```
##              Df  Sum Sq Mean Sq F value    Pr(>F)
## month          4  429686  107421    1.585    0.185
## lakename        1 2498451 2498451   36.855 2.47e-08 ***
## month:lakename  4  288272   72068    1.063    0.379
## Residuals      97 6575834   67792
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 23 observations deleted due to missingness

#format as lm
lm.tn <- lm(data = df.lk0, tn_ug ~ month * lakename)
summary(lm.tn)

##
## Call:
## lm(formula = tn_ug ~ month * lakename, data = df.lk0)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -357.88 -118.10  -10.41   50.58 1353.86
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      300.51     106.30   2.827  0.0057 **
## month6           23.61     123.64   0.191  0.8489
## month7           53.12     127.05   0.418  0.6768
## month8           36.00     127.05   0.283  0.7775
## month9          105.82     184.11   0.575  0.5668
## lakenamePeter Lake    84.43     144.86   0.583  0.5614
## month6:lakenamePeter Lake 200.49     170.90   1.173  0.2436
## month7:lakenamePeter Lake 271.82     176.18   1.543  0.1261
## month8:lakenamePeter Lake 325.05     174.20   1.866  0.0651 .
## month9:lakenamePeter Lake  59.70     278.35   0.214  0.8306
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 260.4 on 97 degrees of freedom
## (23 observations deleted due to missingness)
## Multiple R-squared:  0.3285, Adjusted R-squared:  0.2662
## F-statistic: 5.272 on 9 and 97 DF,  p-value: 7.729e-06

#ignore interactions
aov.tn.maineff <- aov(data = df.lk0, tn_ug ~ month + lakename)
summary(aov.tn.maineff)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## month          4  429686  107421    1.581    0.185
## lakename        1 2498451 2498451   36.763 2.33e-08 ***
## Residuals     101 6864107   67961
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 23 observations deleted due to missingness

#post-hoc test for pairwise differences
TukeyHSD(aov.tn.maineff)

##      Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = tn_ug ~ month + lakename, data = df.lk0)
##
```

```

## $month
##      diff      lwr      upr      p adj
## 6-5 116.294252 -120.82276 353.4113 0.6529963
## 7-5 179.189801 -65.25955 423.6392 0.2565155
## 8-5 202.333236 -39.36107 444.0275 0.1454295
## 9-5 118.016415 -263.04241 499.0752 0.9106171
## 7-6 62.895549 -125.01416 250.8053 0.8847262
## 8-6 86.038984 -98.27247 270.3504 0.6938688
## 9-6 1.722164 -345.78381 349.2281 1.0000000
## 8-7 23.143436 -170.51017 216.7970 0.9973397
## 9-7 -61.173385 -413.72325 291.3765 0.9888483
## 9-8 -84.316821 -434.96203 266.3284 0.9628103
##
## $lakename
##      diff      lwr      upr p adj
## Peter Lake-Paul Lake 305.0995 205.1061 405.0929 0
tn.int <- with(df.lk0, interaction(month, lakename))
aov.tn.int <- aov(data = df.lk0, tn_ug ~ tn.int)
tn.groups <- HSD.test(aov.tn.int, "tn.int", group = TRUE)
tn.groups

## $statistics
##      MSerror Df      Mean      CV
## 67792.1 97 487.4077 53.41917
##
## $parameters
##      test name.t ntr StudentizedRange alpha
##      Tukey tn.int 10      4.579991 0.05
##
## $means
##      tn_ug      std r      Min      Max      Q25      Q50
## 5.Paul Lake 300.5115 67.85647 6 244.870 417.345 251.0738 275.0400
## 5.Peter Lake 384.9389 62.65797 7 312.133 460.791 333.7260 373.0810
## 6.Paul Lake 324.1245 117.32193 17 45.670 439.984 307.8120 342.8260
## 6.Peter Lake 609.0427 379.99046 16 379.781 1962.902 462.9225 497.8530
## 7.Paul Lake 353.6341 40.78474 14 281.421 412.669 328.0188 351.6630
## 7.Peter Lake 709.8848 422.31321 13 352.001 2048.151 571.0920 590.7920
## 8.Paul Lake 336.5081 118.22435 14 163.148 499.251 233.8633 356.6185
## 8.Peter Lake 745.9833 349.34126 15 448.049 1924.631 579.3500 688.5110
## 9.Paul Lake 406.3360 169.15898 3 223.799 557.812 330.5980 437.3970
## 9.Peter Lake 550.4680 183.97504 2 420.378 680.558 485.4230 550.4680
##      Q75
## 5.Paul Lake 329.5267
## 5.Peter Lake 440.5575
## 6.Paul Lake 422.2600
## 6.Peter Lake 606.3447
## 7.Paul Lake 385.5945
## 7.Peter Lake 707.7710
## 8.Paul Lake 423.1365
## 8.Peter Lake 781.0950
## 9.Paul Lake 497.6045
## 9.Peter Lake 615.5130
##
## $comparison

```

```
## NULL
##
## $groups
##          tn_ug groups
## 8.Peter Lake 745.9833    a
## 7.Peter Lake 709.8848    a
## 6.Peter Lake 609.0427   ab
## 9.Peter Lake 550.4680   ab
## 9.Paul Lake  406.3360   ab
## 5.Peter Lake 384.9389   ab
## 7.Paul Lake  353.6341    b
## 8.Paul Lake  336.5081    b
## 6.Paul Lake  324.1245    b
## 5.Paul Lake  300.5115    b
##
## attr(,"class")
## [1] "group"
```

7. Create two plots, with TN (plot 1) or TP (plot 2) as the response variable and month and lake as the predictor variables. Hint: you may use some of the code you used for your visualization assignment. Assign groupings with letters, as determined from your tests. Adjust your axes, aesthetics, and color palettes in accordance with best data visualization practices.
8. Combine your plots with cowplot, with a common legend at the top and the two graphs stacked vertically. Your x axes should be formatted with the same breaks, such that you can remove the title and text of the top legend and retain just the bottom legend.

```
#7
#set the theme
theme.hc01.legendtop <- theme_hc() +
  theme(axis.title = element_text(family = "serif", size = (10)),
        axis.text = element_text(family = "serif", size = (8), color = "black"),
        legend.title = element_text(size = 10, family = "serif"),
        legend.text = element_text(size = 10, family = "serif"),
        legend.key = element_rect(color = NA, fill = NA),
        legend.background = element_rect(color = NA, fill = NA),
        legend.position = c(0.2, 0.9))

theme.hc01.nolegend <- theme_hc() +
  theme(axis.title = element_text(family = "serif", size = (10)),
        axis.text = element_text(family = "serif", size = (8), color = "black"),
        legend.title = element_text(size = 10, family = "serif"),
        legend.text = element_text(size = 10, family = "serif"),
        legend.key = element_rect(color = NA, fill = NA),
        legend.background = element_rect(color = NA, fill = NA),
        legend.position = "none")

#plot TP
box.tp <- ggplot(df.lk0, aes(x = month, y = tp_ug,
                           color = lakename, fill = lakename)) +
  geom_boxplot(alpha = 0.8) +
  labs(x = "Month",
       y = expression(paste("Total Phosphorous (", mu, "g)")),
       color = "", fill = "") +
  scale_color_manual(values = c("cadetblue4", "cornsilk3")) +
```

```

scale_fill_manual(values = c("cadetblue4", "cornsilk3")) +
stat_summary(geom = "text", fun.y = max, vjust = -1, size = 4,
              label = c("bcd", "cd", "ab", "d", "a", "d",
                        "abc", "d", "abcd", "cd"),
              position = position_dodge(width = 0.75)) +
theme.hc01.legendtop
#plot TN
box.tn <- ggplot(df.lk0, aes(x = month, y = tn_ug,
                           color = lakename, fill = lakename)) +
  geom_boxplot(alpha = 0.8) +
  labs(x = "Month",
       y = expression(paste("Total Nitrogen (", mu, "g)")),
       color = "", fill = "") +
  scale_y_continuous(limits = c(0, 1000)) +
  scale_color_manual(values = c("cadetblue4", "cornsilk3")) +
  scale_fill_manual(values = c("cadetblue4", "cornsilk3")) +
  stat_summary(geom = "text", fun.y = max, vjust = -1, size = 4,
              label = c("b", "ab", "ab", "ab", "b",
                        "a", "b", "a", "ab", "ab"),
              position = position_dodge(width = 0.75)) +
theme.hc01.nolegend

#8
plot_grid(box.tp, box.tn, nrow = 2)

```


