

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

Alicia Zhao

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] "/Users/mac/Desktop/Data Analytics/Environmental_Data_Analytics_2020"

library(tidyverse)
library(cowplot)
library(agricolae)

#2
PeterPaul.chem.nutrients <-
  read.csv("./Data/Processed/NTL-LTER_Lake_Chemistry_Nutrients_PeterPaul_Processed.csv")

PeterPaul.chem.nutrients$sampldate <- as.Date(
  PeterPaul.chem.nutrients$sampldate, format = "%Y-%m-%d")

head(PeterPaul.chem.nutrients$sampldate)

## [1] "1984-05-27" "1984-05-27" "1984-05-27" "1984-05-27" "1984-05-27"
## [6] "1984-05-27"
```

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.

```
PeterPaul.surface <- PeterPaul.chem.nutrients %>%
  filter(depth == 0 & year4 %in% c(1993, 1994, 1995, 1996))

PeterPaul.surface$month = as.factor(PeterPaul.surface$month)
class(PeterPaul.surface$month)
```

```
## [1] "factor"
```

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: To test if nutrients are significantly different between Peter Lake and Paul Lake and across seasons, I would use a two-way ANOVA, as I would be testing a continuous response variable (TP, TN) against two categorical explanatory variables (month, lake).

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
TN.anova.2way <- aov(data = PeterPaul.surface, tn_ug ~ lakename + month)
summary (TN.anova.2way)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## lakename    1 2468595 2468595    36.32 2.75e-08 ***
## month       4  459542  114885     1.69   0.158
## Residuals 101 6864107   67961
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 27 observations deleted due to missingness
```

```
TukeyHSD(TN.anova.2way)
```

```
##      Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = tn_ug ~ lakename + month, data = PeterPaul.surface)
##
## $lakename
##              diff          lwr          upr p adj
## Peter Lake-Paul Lake 303.796 203.8026 403.7894    0
##
## $month
##              diff          lwr          upr    p adj
## 6-5 132.58168 -104.53533 369.6987 0.5307817
## 7-5 196.50011  -47.94924 440.9495 0.1761663
## 8-5 208.77984  -32.91447 450.4741 0.1238871
## 9-5 160.08048 -220.97835 541.1393 0.7701126
## 7-6  63.91843 -123.99128 251.8281 0.8785969
## 8-6  76.19815 -108.11330 260.5096 0.7803543
```

```
## 9-6 27.49879 -320.00718 375.0048 0.9994732
## 8-7 12.27972 -181.37388 205.9333 0.9997809
## 9-7 -36.41964 -388.96950 316.1302 0.9984948
## 9-8 -48.69936 -399.34457 301.9458 0.9952369
```

```
TN.anova.2way.int <- aov(data = PeterPaul.surface, tn_ug ~ lakename * month)
summary(TN.anova.2way.int)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## lakename      1 2468595 2468595   36.414 2.91e-08 ***
## month         4  459542  114885    1.695   0.157
## lakename:month 4  288272   72068    1.063   0.379
## Residuals    97 6575834   67792
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 27 observations deleted due to missingness
```

#Interaction is not significant, so sticking with the 2-way without interaction effect.

```
TN.2way <- with(PeterPaul.surface, lakename)
TN.anova.2way.int.2 <- aov(data = PeterPaul.surface, tn_ug ~ TN.2way)
TN.groups <- HSD.test(TN.anova.2way.int.2, "TN.2way", group = TRUE)
TN.groups
```

```
## $statistics
##      MSerror Df      Mean      CV
## 69749.03 105 487.4077 54.1847
##
## $parameters
##      test name.t ntr StudentizedRange alpha
##      Tukey TN.2way 2          2.804124 0.05
##
## $means
##              tn_ug      std r      Min      Max      Q25      Q50      Q75
## Paul Lake 336.9293 100.2745 54  45.670  557.812 284.0107 344.243 411.5165
## Peter Lake 640.7253 361.3738 53 312.133 2048.151 448.0490 571.092 692.4860
##
## $comparison
## NULL
##
## $groups
##              tn_ug groups
## Peter Lake 640.7253      a
## Paul Lake 336.9293      b
##
## attr("class")
## [1] "group"
```

```
#6
TP.anova.2way <- aov(data = PeterPaul.surface, tp_ug ~ lakename + month)
summary(TP.anova.2way)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## lakename      1  10228   10228  94.453 <2e-16 ***
## month         4    813     203   1.876   0.119
## Residuals   123  13320     108
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 5 observations deleted due to missingness
TukeyHSD(TP.anova.2way)

## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tp_ug ~ lakename + month, data = PeterPaul.surface)
##
## $lakename
##           diff          lwr          upr p adj
## Peter Lake-Paul Lake 17.80939 14.18208 21.43669    0
##
## $month
##           diff          lwr          upr          p adj
## 6-5  6.3451786  -3.012727 15.703084 0.3350273
## 7-5  8.8661326  -0.491773 18.224038 0.0723646
## 8-5  4.8191843  -4.469970 14.108339 0.6055077
## 9-5  5.4951391  -6.998304 17.988582 0.7410806
## 7-6  2.5209540  -4.366278  9.408186 0.8487741
## 8-6 -1.5259943  -8.319518  5.267530 0.9713266
## 9-6 -0.8500395 -11.618033  9.917954 0.9994865
## 8-7 -4.0469483 -10.840472  2.746576 0.4691480
## 9-7 -3.3709935 -14.138987  7.397000 0.9084852
## 9-8  0.6759548 -10.032345 11.384255 0.9997883

TP.anova.2way.int <- aov(data = PeterPaul.surface, tp_ug ~ lakename * month)
summary(TP.anova.2way.int)

##           Df Sum Sq Mean Sq F value Pr(>F)
## lakename    1  10228    10228  98.914 <2e-16 ***
## month        4     813       203   1.965 0.1043
## lakename:month 4     1014       254   2.452 0.0496 *
## Residuals   119  12305       103
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 5 observations deleted due to missingness
# Interaction effect is significant!
TukeyHSD(TP.anova.2way.int)

## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tp_ug ~ lakename * month, data = PeterPaul.surface)
##
## $lakename
##           diff          lwr          upr p adj
## Peter Lake-Paul Lake 17.80939 14.26365 21.35513    0
##
## $month
##           diff          lwr          upr          p adj
## 6-5  6.3451786  -2.8038335 15.494191 0.3119085
## 7-5  8.8661326  -0.2828796 18.015145 0.0622967
```

```

## 8-5 4.8191843 -4.2626118 13.900980 0.5839528
## 9-5 5.4951391 -6.7194172 17.709695 0.7243206
## 7-6 2.5209540 -4.2125367 9.254445 0.8376355
## 8-6 -1.5259943 -8.1678685 5.115880 0.9688094
## 9-6 -0.8500395 -11.3776631 9.677584 0.9994372
## 8-7 -4.0469483 -10.6888225 2.594926 0.4453729
## 9-7 -3.3709935 -13.8986170 7.156630 0.9012092
## 9-8 0.6759548 -9.7933076 11.145217 0.9997679
##
## $\lakenamename:month`
##
## diff lwr upr p adj
## Peter Lake:5-Paul Lake:5 4.3135714 -13.9293175 22.5564604 0.9989515
## Paul Lake:6-Paul Lake:5 -0.9178824 -16.4886641 14.6528993 1.0000000
## Peter Lake:6-Paul Lake:5 16.8838889 1.4263507 32.3414270 0.0206973
## Paul Lake:7-Paul Lake:5 -1.7271111 -17.1846493 13.7304270 0.9999981
## Peter Lake:7-Paul Lake:5 22.9304706 7.3596889 38.5012523 0.0002415
## Paul Lake:8-Paul Lake:5 -2.0872222 -17.5447604 13.3703159 0.9999902
## Peter Lake:8-Paul Lake:5 15.0200000 -0.3355071 30.3755071 0.0607728
## Paul Lake:9-Paul Lake:5 -0.7380000 -20.5935673 19.1175673 1.0000000
## Peter Lake:9-Paul Lake:5 14.7452500 -6.4208558 35.9113558 0.4316694
## Paul Lake:6-Peter Lake:5 -5.2314538 -19.9572479 9.4943403 0.9787107
## Peter Lake:6-Peter Lake:5 12.5703175 -2.0356832 27.1763181 0.1571717
## Paul Lake:7-Peter Lake:5 -6.0406825 -20.6466832 8.5653181 0.9437275
## Peter Lake:7-Peter Lake:5 18.6168992 3.8911050 33.3426933 0.0032014
## Paul Lake:8-Peter Lake:5 -6.4007937 -21.0067943 8.2052070 0.9208652
## Peter Lake:8-Peter Lake:5 10.7064286 -3.7915495 25.2044066 0.3464892
## Paul Lake:9-Peter Lake:5 -5.0515714 -24.2516579 14.1485150 0.9975850
## Peter Lake:9-Peter Lake:5 10.4316786 -10.1207861 30.9841433 0.8273658
## Peter Lake:6-Paul Lake:6 17.8017712 6.7120688 28.8914737 0.0000401
## Paul Lake:7-Paul Lake:6 -0.8092288 -11.8989312 10.2804737 1.0000000
## Peter Lake:7-Paul Lake:6 23.8483529 12.6013419 35.0953640 0.0000000
## Paul Lake:8-Paul Lake:6 -1.1693399 -12.2590423 9.9203626 0.9999989
## Peter Lake:8-Paul Lake:6 15.9378824 4.9908457 26.8849190 0.0003006
## Paul Lake:9-Paul Lake:6 0.1798824 -16.5021309 16.8618956 1.0000000
## Peter Lake:9-Paul Lake:6 15.6631324 -2.5591082 33.8853729 0.1584032
## Paul Lake:7-Peter Lake:6 -18.6110000 -29.5411300 -7.6808700 0.0000101
## Peter Lake:7-Peter Lake:6 6.0465817 -5.0431207 17.1362841 0.7595330
## Paul Lake:8-Peter Lake:6 -18.9711111 -29.9012412 -8.0409811 0.0000062
## Peter Lake:8-Peter Lake:6 -1.8638889 -12.6492426 8.9214648 0.9999197
## Paul Lake:9-Peter Lake:6 -17.6218889 -34.1982518 -1.0455259 0.0276305
## Peter Lake:9-Peter Lake:6 -2.1386389 -20.2642090 15.9869312 0.9999970
## Peter Lake:7-Paul Lake:7 24.6575817 13.5678793 35.7472841 0.0000000
## Paul Lake:8-Paul Lake:7 -0.3601111 -11.2902412 10.5700189 1.0000000
## Peter Lake:8-Paul Lake:7 16.7471111 5.9617574 27.5324648 0.0000827
## Paul Lake:9-Paul Lake:7 0.9891111 -15.5872518 17.5654741 1.0000000
## Peter Lake:9-Paul Lake:7 16.4723611 -1.6532090 34.5979312 0.1087387
## Paul Lake:8-Peter Lake:7 -25.0176928 -36.1073952 -13.9279904 0.0000000
## Peter Lake:8-Peter Lake:7 -7.9104706 -18.8575073 3.0365661 0.3778093
## Paul Lake:9-Peter Lake:7 -23.6684706 -40.3504838 -6.9864574 0.0004851
## Peter Lake:9-Peter Lake:7 -8.1852206 -26.4074611 10.0370199 0.9089776
## Peter Lake:8-Paul Lake:8 17.1072222 6.3218685 27.8925759 0.0000523
## Paul Lake:9-Paul Lake:8 1.3492222 -15.2271407 17.9255852 0.9999999
## Peter Lake:9-Paul Lake:8 16.8324722 -1.2930979 34.9580424 0.0926020
## Paul Lake:9-Peter Lake:8 -15.7580000 -32.2392597 0.7232597 0.0735733

```

```
## Peter Lake:9-Peter Lake:8 -0.2747500 -18.3133864 17.7638864 1.0000000
## Peter Lake:9-Paul Lake:9 15.4832500 -6.5132124 37.4797124 0.4163366

TP.interaction <- with(PeterPaul.surface, interaction(lakename, month))
TP.anova.2way.int.2 <- aov(data = PeterPaul.surface, tp_ug ~ TP.interaction)
TP.groups <- HSD.test(TP.anova.2way.int.2, "TP.interaction", 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.interaction 10      4.560262 0.05
##
## $means
##              tp_ug      std r      Min      Max      Q25      Q50      Q75
## Paul Lake.5 11.474000 3.928545 6 7.001 17.090 8.1395 11.8885 13.53675
## Paul Lake.6 10.556118 4.416821 17 1.222 16.697 7.4430 10.6050 13.94600
## Paul Lake.7 9.746889 3.525120 18 4.501 21.763 7.8065 9.1555 10.65700
## Paul Lake.8 9.386778 1.478062 18 5.879 11.542 8.4495 9.6090 10.45050
## Paul Lake.9 10.736000 3.615978 5 6.592 16.281 8.9440 10.1920 11.67100
## Peter Lake.5 15.787571 2.719954 7 10.887 18.922 14.8915 15.5730 17.67400
## Peter Lake.6 28.357889 15.588507 18 10.974 53.388 14.7790 24.6840 41.13000
## Peter Lake.7 34.404471 18.285568 17 19.149 66.893 21.6640 24.2070 50.54900
## Peter Lake.8 26.494000 9.829596 19 14.551 49.757 21.2425 23.2250 27.99350
## Peter Lake.9 26.219250 10.814803 4 16.281 41.145 19.6845 23.7255 30.26025
##
## $comparison
## NULL
##
## $groups
##              tp_ug groups
## Peter Lake.7 34.404471      a
## Peter Lake.6 28.357889      ab
## Peter Lake.8 26.494000      abc
## Peter Lake.9 26.219250      abcd
## Peter Lake.5 15.787571      bcd
## Paul Lake.5 11.474000      cd
## Paul Lake.9 10.736000      cd
## Paul Lake.6 10.556118      d
## Paul Lake.7 9.746889      d
## Paul Lake.8 9.386778      d
##
## attr(,"class")
## [1] "group"
```

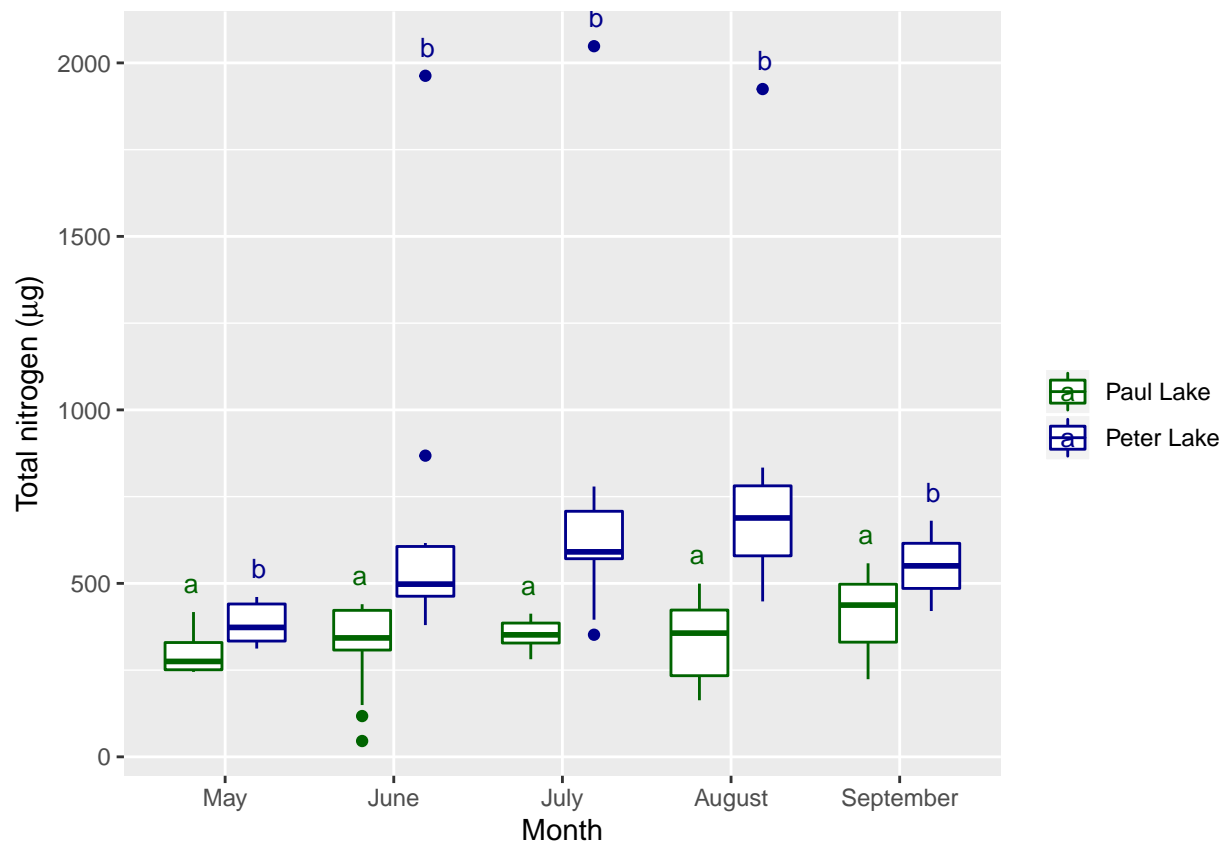
7. Create two plots, with TN (plot 1) or TP (plot 2) as the response variable and month and lake as the predictor variables. Hint: you may use some of the code you used for your visualization assignment. Assign groupings with letters, as determined from your tests. Adjust your axes, aesthetics, and color 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
TN.anova.plot <- ggplot(subset(PeterPaul.surface,month %in% c("5", "6","7","8","9") ),
  aes(y = tn_ug, x = month, color = lakename)) +
  geom_boxplot() +
  scale_color_manual(values = c("dark green", "dark blue")) +
  labs(x = "Month", y = expression(paste("Total nitrogen (", mu,"g)")), color = " ") +
  scale_x_discrete(labels = c("May", "June", "July", "August", "September")) +
  stat_summary(geom = "text", fun.y = max, vjust = -1, size = 3.5,
    label = c("a", "b",
              "a", "b",
              "a", "b",
              "a", "b",
              "a", "b"), position=position_dodge2(.8))

print(TN.anova.plot)
```

Warning: Removed 24 rows containing non-finite values (stat_boxplot).

Warning: Removed 24 rows containing non-finite values (stat_summary).



```
TP.anova.plot <- ggplot(subset(PeterPaul.surface,month %in% c("5", "6","7","8","9")),
  aes(y = tp_ug, x = month, color = lakename)) +
  geom_boxplot() +
  scale_color_manual(values = c("dark green", "dark blue")) +
  labs(x = "Month", y = expression(paste("Total phosphorus (", mu,"g)")), color = " ") +
  scale_x_discrete(labels = c("May", "June", "July", "August", "September")) +
```

```

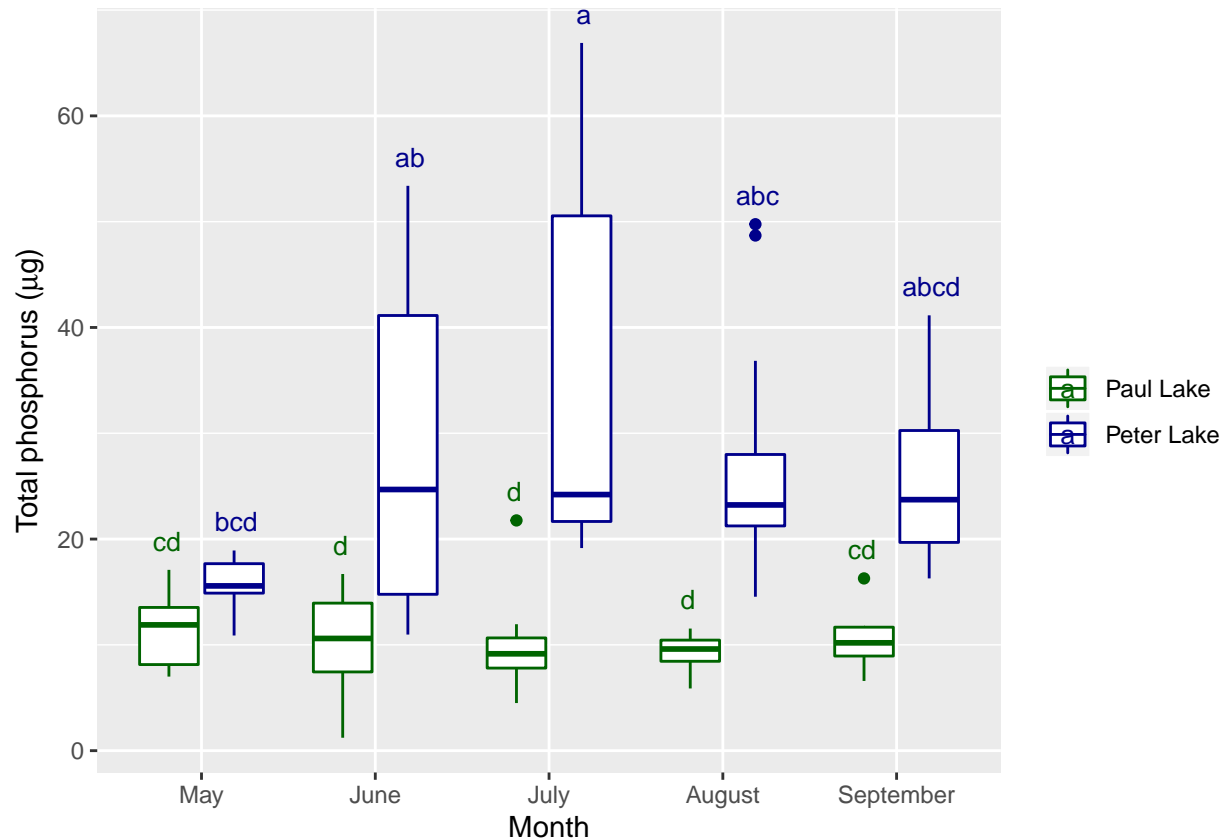
stat_summary(geom = "text", fun.y = max, vjust = -1, size = 3.5,
             label = c("cd", "bcd", "d", "ab", "d",
                       "a", "d", "abc", "cd", "abcd"), position=position_dodge2(.8))

print(TP.anova.plot)

```

Warning: Removed 2 rows containing non-finite values (stat_boxplot).

Warning: Removed 2 rows containing non-finite values (stat_summary).



```

#8
plot_grid(
  TN.anova.plot + theme (legend.position = "top") + labs(x=""),
  TP.anova.plot + theme (legend.position = "none"),
  nrow = 2, align = 'v')

```

Warning: Removed 24 rows containing non-finite values (stat_boxplot).

Warning: Removed 24 rows containing non-finite values (stat_summary).

Warning: Removed 2 rows containing non-finite values (stat_boxplot).

Warning: Removed 2 rows containing non-finite values (stat_summary).

