# Assignment 6: GLMs (Linear Regressios, ANOVA, & t-tests)

## Jess Ozog

## OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics on generalized linear models.

## 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., "Fay_A06_GLMs.Rmd") prior to submission.

The completed exercise is due on Monday, February 28 at 7:00 pm.

## Set up your session

1. Set up your session. Check your working directory. Load the tidyverse, agricolae and other needed packages. Import the *raw* NTL-LTER raw data file for chemistry/physics (`NTL-LTER_Lake_ChemistryPhysics_Raw.csv`). Set date columns to date objects.

2. Build a ggplot theme and set it as your default theme.

```
#1
getwd()
```

```
## [1] "C:/ENV872/Environmental_Data_Analytics_2022/Assignments"
```

```
library(tidyverse)
```

```
## -- Attaching packages --------------------------------------- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v purrr   0.3.4
## v tibble  3.1.4      v dplyr   1.0.7
## v tidyr   1.1.4      v stringr 1.4.0
## v readr   2.0.2      v forcats 0.5.1
```

```
## -- Conflicts ------------------------------------------ tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
library(agricolae)
library(lubridate)

##
## Attaching package: 'lubridate'

## The following objects are masked from 'package:base':
##
##     date, intersect, setdiff, union

chemphys <-  read.csv("../Data/Raw/NTL-LTER_Lake_ChemistryPhysics_Raw.csv",
                      stringsAsFactors = TRUE)
chemphys$sampledate<-as.Date(chemphys$sampledate, format = "%m/%d/%y")

#2
theme1 <- theme_bw(base_size=12) +
  theme(panel.border=element_rect(fill="transparent",color="gray",size=1),
        plot.title = element_text(hjust = 0.5), legend.position="right")
theme_set(theme1)
```

## Simple regression

Our first research question is: Does mean lake temperature recorded during July change with depth across all lakes?

3. State the null and alternative hypotheses for this question: > Answer: H0: Mean lake temperature in July will be equal across all lakes at all depths. Ha: Mean lake temperature in July will not be equal across all lakes at all depths.

4. Wrangle your NTL-LTER dataset with a pipe function so that the records meet the following criteria:

- Only dates in July.
- Only the columns: `lakename`, `year4`, `daynum`, `depth`, `temperature_C`
- Only complete cases (i.e., remove NAs)

5. Visualize the relationship among the two continuous variables with a scatter plot of temperature by depth. Add a smoothed line showing the linear model, and limit temperature values from 0 to 35 °C. Make this plot look pretty and easy to read.

```
#4
chemphys.processed <-
  chemphys %>%
  filter(month(sampledate)==7) %>%
  select(lakename:daynum,depth,temperature_C) %>%
  drop_na(temperature_C)

#5
plot1 <-
  ggplot(chemphys.processed, aes(x=depth, y=temperature_C)) +
  geom_point(size=2.5, color="royalblue2", alpha=0.25) +
  ylim(0,35) +
```
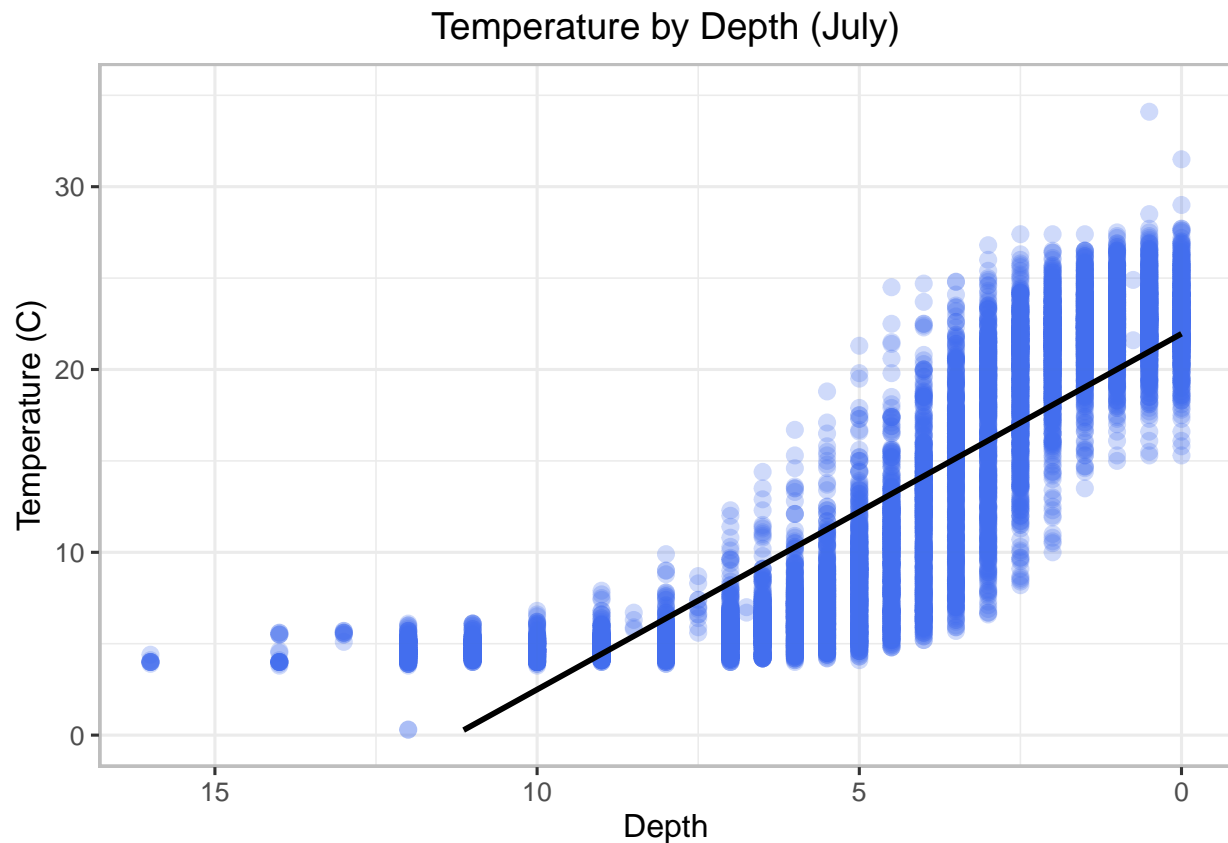
```
  scale_x_reverse() +
  geom_smooth(method=lm, color="black") +
  labs(title="Temperature by Depth (July)", x="Depth", y="Temperature (C)")
plot1
```

## `geom_smooth()` using formula 'y ~ x'

## Warning: Removed 24 rows containing missing values (geom_smooth).



Temperature by Depth (July)

6. Interpret the figure. What does it suggest with regards to the response of temperature to depth? Do the distribution of points suggest about anything about the linearity of this trend?

   Answer: There is an inverse relationship between depth and temperture. As depth increases, temperature decreases. The distribution of points suggest that the points do not display a linear relationship and may need to be log transformed.

7. Perform a linear regression to test the relationship and display the results

```
#7
lm1 <- lm(data = chemphys.processed, temperature_C ~ depth)
summary(lm1)
```

```
##
## Call:
## lm(formula = temperature_C ~ depth, data = chemphys.processed)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -9.5173 -3.0192  0.0633  2.9365 13.5834
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 21.95597    0.06792   323.3   <2e-16 ***
## depth       -1.94621    0.01174  -165.8   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.835 on 9726 degrees of freedom
## Multiple R-squared:  0.7387, Adjusted R-squared:  0.7387
## F-statistic: 2.75e+04 on 1 and 9726 DF,  p-value: < 2.2e-16
```

8. Interpret your model results in words. Include how much of the variability in temperature is explained by changes in depth, the degrees of freedom on which this finding is based, and the statistical significance of the result. Also mention how much temperature is predicted to change for every 1m change in depth.

   Answer: Lake temperature in July is affected by changes in depth. 73.87% of the variability in temperature is explained by changes in depth. For every unit decrease in depth, temperatures decreases by 1.95 degrees (C).(R^2 = 0.7387, df = 9726, p-value < 0.0001)

---

## Multiple regression

Let's tackle a similar question from a different approach. Here, we want to explore what might the best set of predictors for lake temperature in July across the monitoring period at the North Temperate Lakes LTER.

9. Run an AIC to determine what set of explanatory variables (year4, daynum, depth) is best suited to predict temperature.

10. Run a multiple regression on the recommended set of variables.

```
#9
lm2 <- lm(data = chemphys.processed, temperature_C ~ year4 + daynum + depth)
step(lm2)
```

```
## Start:  AIC=26065.53
## temperature_C ~ year4 + daynum + depth
##
##          Df Sum of Sq    RSS   AIC
## <none>                 141687 26066
## - year4   1       101 141788 26070
## - daynum  1      1237 142924 26148
## - depth   1    404475 546161 39189
```

```
## 
## Call:
## lm(formula = temperature_C ~ year4 + daynum + depth, data = chemphys.processed)
## 
## Coefficients:
## (Intercept)         year4        daynum         depth
##    -8.57556       0.01134       0.03978      -1.94644
```

*#10*
```
summary(lm2)
```

```
## 
## Call:
## lm(formula = temperature_C ~ year4 + daynum + depth, data = chemphys.processed)
## 
## Residuals:
##     Min      1Q  Median      3Q     Max
## -9.6536 -3.0000  0.0902  2.9658 13.6123
## 
## Coefficients:
##               Estimate Std. Error  t value Pr(>|t|)
## (Intercept) -8.575564   8.630715   -0.994  0.32044
## year4        0.011345   0.004299    2.639  0.00833 **
## daynum       0.039780   0.004317    9.215  < 2e-16 ***
## depth       -1.946437   0.011683 -166.611  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 3.817 on 9724 degrees of freedom
## Multiple R-squared:  0.7412, Adjusted R-squared:  0.7411
## F-statistic:  9283 on 3 and 9724 DF,  p-value: < 2.2e-16
```

```
AIC(lm1,lm2)
```

```
##     df      AIC
## lm1  3 53762.12
## lm2  5 53674.39
```

11. What is the final set of explanatory variables that the AIC method suggests we use to predict temperature in our multiple regression? How much of the observed variance does this model explain? Is this an improvement over the model using only depth as the explanatory variable?

    Answer: According to the AIC, temperature is explained by all the explanatory variables (year4, daynum, and depth). This full model explains 74.12% of the observed variance, which is slightly more than the model with only depth as an explanatory variable. The full model is an imporvement over the depth model, due to the reduction in the AIC value. With only depth, AIC was 53762.12 and was reduced to 53674.39 with the full model.

## Analysis of Variance

12. Now we want to see whether the different lakes have, on average, different temperatures in the month of July. Run an ANOVA test to complete this analysis. (No need to test assumptions of normality or similar variances.) Create two sets of models: one expressed as an ANOVA models and another expressed as a linear model (as done in our lessons).

```
#12
anova1 <- aov(data = chemphys.processed, temperature_C ~ lakename)
summary(anova1)
```

```
##               Df Sum Sq Mean Sq F value Pr(>F)
## lakename       8  21642  2705.2      50 <2e-16 ***
## Residuals   9719 525813    54.1
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
anova2 <- lm(data = chemphys.processed, temperature_C ~ lakename)
summary(anova2)
```

```
##
## Call:
## lm(formula = temperature_C ~ lakename, data = chemphys.processed)
##
## Residuals:
##     Min     1Q  Median     3Q     Max
## -10.769  -6.614  -2.679   7.684  23.832
##
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)               17.6664     0.6501  27.174  < 2e-16 ***
## lakenameCrampton Lake     -2.3145     0.7699  -3.006 0.002653 **
## lakenameEast Long Lake    -7.3987     0.6918 -10.695  < 2e-16 ***
## lakenameHummingbird Lake  -6.8931     0.9429  -7.311 2.87e-13 ***
## lakenamePaul Lake         -3.8522     0.6656  -5.788 7.36e-09 ***
## lakenamePeter Lake        -4.3501     0.6645  -6.547 6.17e-11 ***
## lakenameTuesday Lake      -6.5972     0.6769  -9.746  < 2e-16 ***
## lakenameWard Lake         -3.2078     0.9429  -3.402 0.000672 ***
## lakenameWest Long Lake    -6.0878     0.6895  -8.829  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.355 on 9719 degrees of freedom
## Multiple R-squared:  0.03953,    Adjusted R-squared:  0.03874
## F-statistic:    50 on 8 and 9719 DF,  p-value: < 2.2e-16
```

13. Is there a significant difference in mean temperature among the lakes? Report your findings.
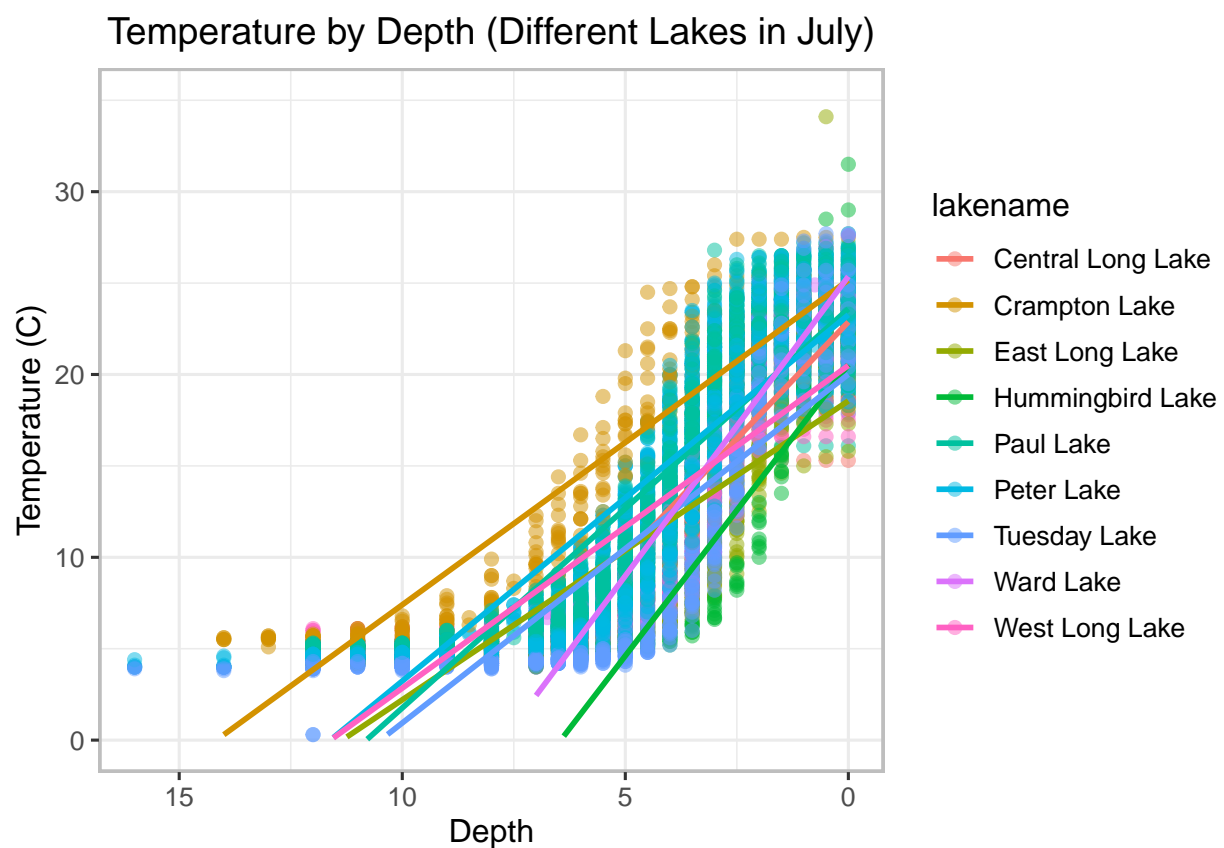
    Answer: There is a signficant difference in mean temperature among lakes ($R^2 = 0.04$, df $= 9719$, p-value $< 0.0001$).

14. Create a graph that depicts temperature by depth, with a separate color for each lake. Add a geom_smooth (method = "lm", se = FALSE) for each lake. Make your points 50 % transparent. Adjust your y axis limits to go from 0 to 35 degrees. Clean up your graph to make it pretty.

```
#14.
plot2 <- ggplot(chemphys.processed,aes(x = depth, y = temperature_C, color=lakename)) +
  geom_point(size = 2, alpha=0.5) +
  ylim(0,35) +
  scale_x_reverse() +
  geom_smooth(method=lm, se = FALSE) +
  labs(title="Temperature by Depth (Different Lakes in July)", x="Depth", y="Temperature (C)")
plot2
```

```
## `geom_smooth()` using formula 'y ~ x'
```

```
## Warning: Removed 73 rows containing missing values (geom_smooth).
```



15. Use the Tukey's HSD test to determine which lakes have different means.

```
#15
lakeGroups <- HSD.test(anova1, "lakename", group=T)
lakeGroups
```

```
## $statistics
##    MSerror   Df     Mean       CV
##    54.1016 9719 12.72087 57.82135
##
```

```
## $parameters
##    test    name.t ntr StudentizedRange alpha
##   Tukey lakename   9         4.387504  0.05
##
## $means
##                  temperature_C      std    r Min  Max    Q25   Q50    Q75
## Central Long Lake     17.66641 4.196292  128 8.9 26.8 14.400 18.40 21.000
## Crampton Lake         15.35189 7.244773  318 5.0 27.5  7.525 16.90 22.300
## East Long Lake        10.26767 6.766804  968 4.2 34.1  4.975  6.50 15.925
## Hummingbird Lake      10.77328 7.017845  116 4.0 31.5  5.200  7.00 15.625
## Paul Lake             13.81426 7.296928 2660 4.7 27.7  6.500 12.40 21.400
## Peter Lake            13.31626 7.669758 2872 4.0 27.0  5.600 11.40 21.500
## Tuesday Lake          11.06923 7.698687 1524 0.3 27.7  4.400  6.80 19.400
## Ward Lake             14.45862 7.409079  116 5.7 27.6  7.200 12.55 23.200
## West Long Lake        11.57865 6.980789 1026 4.0 25.7  5.400  8.00 18.800
##
## $comparison
## NULL
##
## $groups
##                  temperature_C groups
## Central Long Lake     17.66641      a
## Crampton Lake         15.35189     ab
## Ward Lake             14.45862     bc
## Paul Lake             13.81426      c
## Peter Lake            13.31626      c
## West Long Lake        11.57865      d
## Tuesday Lake          11.06923     de
## Hummingbird Lake      10.77328     de
## East Long Lake        10.26767      e
##
## attr(,"class")
## [1] "group"
```

16. From the findings above, which lakes have the same mean temperature, statistically speaking, as Peter Lake? Does any lake have a mean temperature that is statistically distinct from all the other lakes?

> Answer: The lakes that statistically have the same mean temperature as Peter Lake are Paul Lake and Ward Lake. None of the lakes have a mean temperature that is statistically distinct from all the other lakes.

17. If we were just looking at Peter Lake and Paul Lake. What's another test we might explore to see whether they have distinct mean temperatures?

> Answer: If only interested in the differences in mean temperatures at Peter Lake and Paul Lake, a two-sample t-test could be used.