

Assignment 6: Generalized Linear Models

Caroline Watson

OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics (ENV872L) on generalized linear models.

Directions

1. Change “Student Name” on line 3 (above) with your name.
2. Use the lesson as a guide. It contains code that can be modified to complete the assignment.
3. Work through the steps, **creating code and output** that fulfill each instruction.
4. Be sure to **answer the questions** in this assignment document. Space for your answers is provided in this document and is indicated by the “>” character. If you need a second paragraph be sure to start the first line with “>”. You should notice that the answer is highlighted in green by RStudio.
5. When you have completed the assignment, **Knit** the text and code into a single PDF file. You will need to have the correct software installed to do this (see Software Installation Guide) Press the **Knit** button in the RStudio scripting panel. This will save the PDF output in your Assignments folder.
6. After Knitting, please submit the completed exercise (PDF file) to the dropbox in Sakai. Please add your last name into the file name (e.g., “Salk_A06_GLMs.pdf”) prior to submission.

The completed exercise is due on Tuesday, 26 February, 2019 before class begins.

Set up your session

1. Set up your session. Upload the EPA Ecotox dataset for Neonicotinoids and the NTL-LTER raw data file for chemistry/physics.
2. Build a ggplot theme and set it as your default theme.

```
#1
```

```
getwd()
```

```
## [1] "/Users/carolinewatson/Documents/Spring 2019/Environmental Data Analytics/Env_Data_Analytics/Ass
```

```
suppressMessages(library(tidyverse))
```

```
library(viridis)
```

```
## Loading required package: viridisLite
```

```
library(gridExtra)
```

```
##
```

```
## Attaching package: 'gridExtra'
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
##      combine
```

```
library(RColorBrewer)
```

```
library(colormap)
```

```
library(dplyr)
```

```

#uploading EPA Ecotox dataset
ecotox.neonic.data <- read.csv("../Data/Raw/ECOTOX_Neonicotinoids_Mortality_raw.csv")

#fixing format for column headings because they added .. when the data was imported
colnames(ecotox.neonic.data)[8:12] <- c("Duration", "Conc.Type", "Conc.Mean", "Conc.Units", "Pub.Year")

#uploading NTL-LTER dataset
ntllter.chem.phys.data <- read.csv("../Data/Raw/NTL-LTER_Lake_ChemistryPhysics_Raw.csv")

#2
#building ggplot theme and setting it as the default theme
caroline_theme <- theme_classic(base_size = 16) +
  theme(axis.text = element_text(color = "black"),
        legend.position = "right")

theme_set(caroline_theme)

```

Neonicotinoids test

Research question: Were studies on various neonicotinoid chemicals conducted in different years?

3. Generate a line of code to determine how many different chemicals are listed in the Chemical.Name column.
4. Are the publication years associated with each chemical well-approximated by a normal distribution? Run the appropriate test and also generate a frequency polygon to illustrate the distribution of counts for each year, divided by chemical name. Bonus points if you can generate the results of your test from a pipe function. No need to make this graph pretty.
5. Is there equal variance among the publication years for each chemical? Hint: var.test is not the correct function.

```

#3 seeing how many chemicals there are in the chemical name column
summary(ecotox.neonic.data$Chemical.Name)

##  Acetamiprid Clothianidin  Dinotefuran Imidacloprid Imidaclothiz
##           136           74           59           695           9
##  Nitenpyram  Nithiazine  Thiacloprid Thiamethoxam
##           21           22           106           161

#4 #turn data into vector within pipe because otherwise shapiro test doesn't like it
shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Acetamiprid"])

##
##  Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Acetamiprid"]
## W = 0.90191, p-value = 5.706e-08

shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Clothianidin"])

##
##  Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Clothianidin"]
## W = 0.69577, p-value = 4.287e-11

```

```

shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Dinotefuran"])

##
## Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Dinotefuran"]
## W = 0.82848, p-value = 8.83e-07
shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Imidacloprid"])

##
## Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Imidacloprid"]
## W = 0.88178, p-value < 2.2e-16
shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Imidaclothiz"])

##
## Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Imidaclothiz"]
## W = 0.68429, p-value = 0.00093
shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Nitenpyram"])

##
## Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Nitenpyram"]
## W = 0.79592, p-value = 0.0005686
shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Nithiazine"])

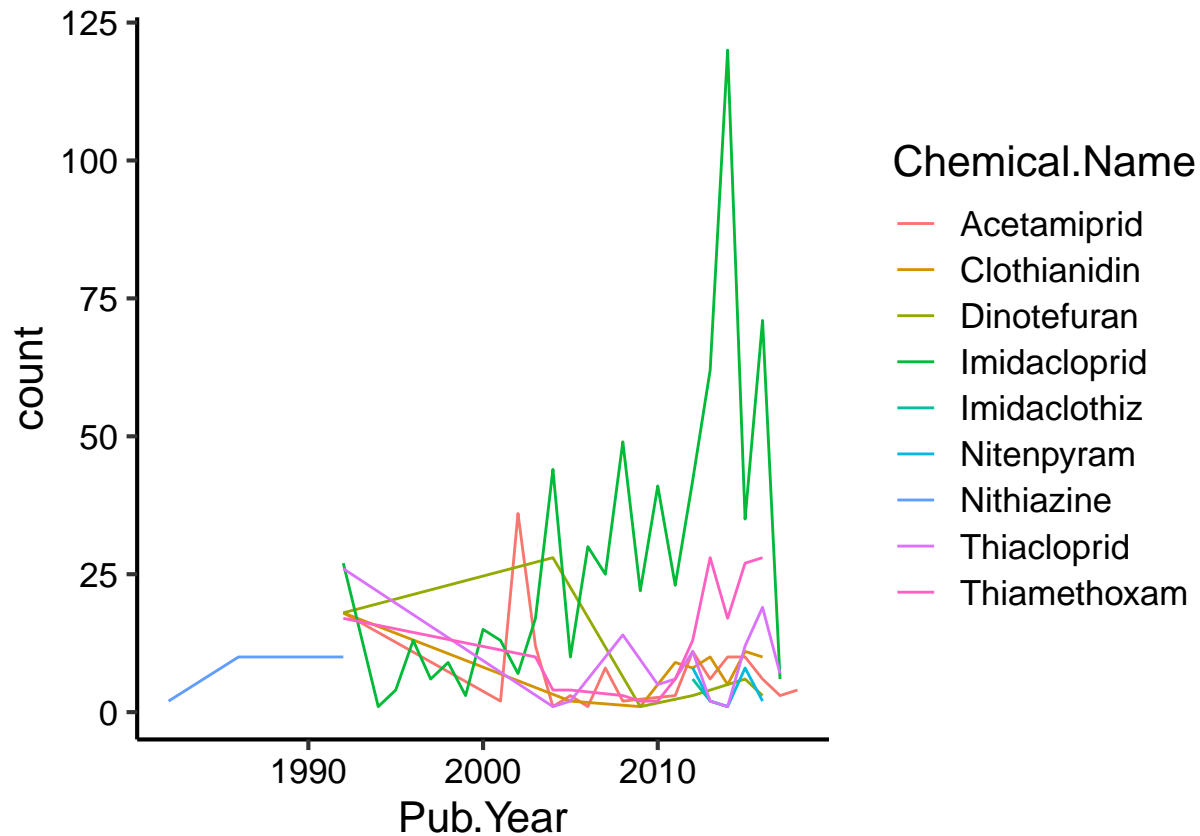
##
## Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Nithiazine"]
## W = 0.75938, p-value = 0.0001235
shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Thiacloprid"])

##
## Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Thiacloprid"]
## W = 0.7669, p-value = 1.118e-11
shapiro.test(ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Thiamethoxam"])

##
## Shapiro-Wilk normality test
##
## data:  ecotox.neonic.data$Pub.Year[ecotox.neonic.data$Chemical.Name == "Thiamethoxam"]
## W = 0.7071, p-value < 2.2e-16
#frequency polygon
freq_pub_chem_plot <- ggplot(ecotox.neonic.data, mapping = aes(x = Pub.Year, color = Chemical.Name)) +

```

```
geom_freqpoly(stat = "count")
print(freq_pub_chem_plot)
```



```
#5
#testing if variances are equal among each chemical
bartlett.test(ecotox.neonic.data$Pub.Year ~ ecotox.neonic.data$Chemical.Name)
```

```
##
## Bartlett test of homogeneity of variances
##
## data: ecotox.neonic.data$Pub.Year by ecotox.neonic.data$Chemical.Name
## Bartlett's K-squared = 139.59, df = 8, p-value < 2.2e-16
```

6. Based on your results, which test would you choose to run to answer your research question?

ANSWER: The variances are not equal, so we would want to run a non-parametric test for a one-way ANOVA and that would be a Kruskal-Wallis test.

7. Run this test below.

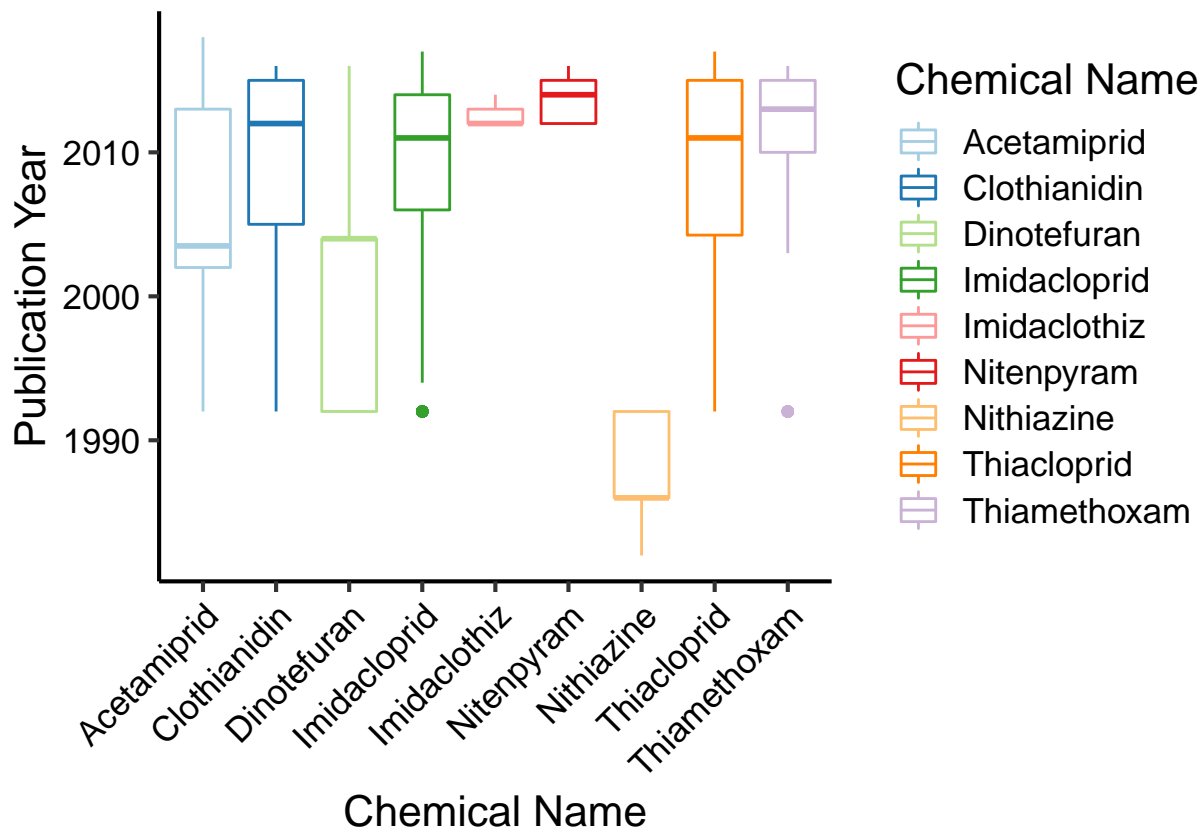
8. Generate a boxplot representing the range of publication years for each chemical. Adjust your graph to make it pretty.

```
#7
#nonparametric Kruskal-Wallis test
pubyear.chem.kw <- kruskal.test(ecotox.neonic.data$Pub.Year ~ ecotox.neonic.data$Chemical.Name)
pubyear.chem.kw
```

```
##
## Kruskal-Wallis rank sum test
```

```
##
## data: ecotox.neonic.data$Pub.Year by ecotox.neonic.data$Chemical.Name
## Kruskal-Wallis chi-squared = 134.15, df = 8, p-value < 2.2e-16

#8
#boxplot of publication years for each chemical
pub_chem_boxplot <- ggplot(ecotox.neonic.data, aes(x = Chemical.Name, y = Pub.Year, color = Chemical.Name)) +
  geom_boxplot() +
  labs(x = "Chemical Name", y = "Publication Year",
       color = "Chemical Name") +
  scale_color_brewer(palette = "Paired") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
print(pub_chem_boxplot)
```



9. How would you summarize the conclusion of your analysis? Include a sentence summarizing your findings and include the results of your test in parentheses at the end of the sentence.

ANSWER: There is a significant difference between publication year and chemical name (Kruskal-Wallis test; $df = 8$, $\chi^2 = 134.2$, $p < 0.0001$).

NTL-LTER test

Research question: What is the best set of predictors for lake temperatures in July across the monitoring period at the North Temperate Lakes LTER?

11. Wrangle your NTL-LTER dataset with a pipe function so that it contains only the following criteria:

- Only dates in July (hint: use the daynum column). No need to consider leap years.
- Only the columns: lakename, year4, daynum, depth, temperature_C

- Only complete cases (i.e., remove NAs)

12. Run an AIC to determine what set of explanatory variables (year4, daynum, depth) is best suited to predict temperature. Run a multiple regression on the recommended set of variables.

```
#11
#wrangling data to only include dates in July, only columns are lakename, year4, daynum, depth, tempera

ntller.chem.phys.data.processed <- ntllter.chem.phys.data %>%
  filter(daynum >= 182 & daynum <= 212) %>%
  select(lakename:daynum, depth, temperature_C) %>%
  na.omit()

#12
#running linear model
yearAIC <- lm(data = ntller.chem.phys.data.processed, temperature_C ~ year4 +
daynum + depth)
summary(yearAIC)

##
## Call:
## lm(formula = temperature_C ~ year4 + daynum + depth, data = ntller.chem.phys.data.processed)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.6517 -2.9937  0.0855  2.9692 13.6171
##
## Coefficients:
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept) -6.455560   8.638808  -0.747   0.4549
## year4         0.010131   0.004303   2.354   0.0186 *
## daynum        0.041336   0.004315   9.580  <2e-16 ***
## depth        -1.947264   0.011676 -166.782  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.811 on 9718 degrees of freedom
## Multiple R-squared:  0.7417, Adjusted R-squared:  0.7417
## F-statistic: 9303 on 3 and 9718 DF, p-value: < 2.2e-16

#running AIC with step model
step(yearAIC)

## Start:  AIC=26016.31
## temperature_C ~ year4 + daynum + depth
##
##           Df Sum of Sq  RSS   AIC
## <none>                 141118 26016
## - year4    1           80 141198 26020
## - daynum   1          1333 142450 26106
## - depth    1         403925 545042 39151
##
## Call:
## lm(formula = temperature_C ~ year4 + daynum + depth, data = ntller.chem.phys.data.processed)
##
## Coefficients:
```

```
## (Intercept)      year4      daynum      depth
##      -6.45556      0.01013      0.04134      -1.94726

#running a multiple regression on the full model again
yearAIC_rerun <- lm(data = ntller.chem.phys.data.processed, temperature_C ~ year4 +
daynum + depth)
summary(yearAIC)

##
## Call:
## lm(formula = temperature_C ~ year4 + daynum + depth, data = ntller.chem.phys.data.processed)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.6517 -2.9937  0.0855  2.9692 13.6171
##
## Coefficients:
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept) -6.455560   8.638808  -0.747   0.4549
## year4        0.010131   0.004303   2.354   0.0186 *
## daynum       0.041336   0.004315   9.580  <2e-16 ***
## depth       -1.947264   0.011676 -166.782 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.811 on 9718 degrees of freedom
## Multiple R-squared:  0.7417, Adjusted R-squared:  0.7417
## F-statistic: 9303 on 3 and 9718 DF,  p-value: < 2.2e-16
```

13. What is the final linear equation to predict temperature from your multiple regression? How much of the observed variance does this model explain?

ANSWER: The full original model has the smallest AIC, thus we would accept this model. The full model explains 74% of the variance.

14. Run an interaction effects ANCOVA to predict temperature based on depth and lakenname from the same wrangled dataset.

```
#14
#interaction effects ANOVA
temp.anova <- lm(data = ntller.chem.phys.data.processed, temperature_C ~ depth*lakenname)
summary(temp.anova)

##
## Call:
## lm(formula = temperature_C ~ depth * lakenname, data = ntller.chem.phys.data.processed)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.6455 -2.9133 -0.2879  2.7567 16.3606
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      22.9455      0.5861  39.147 < 2e-16 ***
## depth           -2.5820      0.2411 -10.711 < 2e-16 ***
## lakennameCrampton Lake      2.2173      0.6804   3.259  0.00112 **
## lakennameEast Long Lake    -4.3884      0.6191  -7.089 1.45e-12 ***
```

```
## lakenamHummingbird Lake      -2.4126      0.8379    -2.879    0.00399 **
## lakenamPaul Lake              0.6105      0.5983     1.020    0.30754
## lakenamPeter Lake             0.2998      0.5970     0.502    0.61552
## lakenamTuesday Lake          -2.8932      0.6060    -4.774    1.83e-06 ***
## lakenamWard Lake              2.4180      0.8434     2.867    0.00415 **
## lakenamWest Long Lake        -2.4663      0.6168    -3.999    6.42e-05 ***
## depth:lakenamCrampton Lake    0.8058      0.2465     3.268    0.00109 **
## depth:lakenamEast Long Lake   0.9465      0.2433     3.891    0.00010 ***
## depth:lakenamHummingbird Lake -0.6026      0.2919    -2.064    0.03903 *
## depth:lakenamPaul Lake        0.4022      0.2421     1.662    0.09664 .
## depth:lakenamPeter Lake       0.5799      0.2418     2.398    0.01649 *
## depth:lakenamTuesday Lake     0.6605      0.2426     2.723    0.00648 **
## depth:lakenamWard Lake        -0.6930      0.2862    -2.421    0.01548 *
## depth:lakenamWest Long Lake   0.8154      0.2431     3.354    0.00080 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.471 on 9704 degrees of freedom
## Multiple R-squared:  0.7861, Adjusted R-squared:  0.7857
## F-statistic: 2097 on 17 and 9704 DF,  p-value: < 2.2e-16
```

15. Is there an interaction between depth and lakenam? How much variance in the temperature observations does this explain?

ANSWER: There is an interaction between depth and lakenam. 79% of the variance is explained in the temperature observations.

16. 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.

#16

```
temp.by.depth <- ggplot(ntlcr.chem.phys.data.processed, aes(x = depth, y = temperature_C, color = lakenam)) +
  geom_point(alpha = 0.5) +
  geom_smooth(method = "lm", se = FALSE) +
  labs(x = "Depth (m)", y = "Temperature (C)", color = "Lake Name") +
  scale_color_brewer(palette = "Paired") +
  ylim(0, 35)
#change colors, change legen label name, put in individual codes/numbers
print(temp.by.depth)
```

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