

Assignment 6: Generalized Linear Models

Ying Wei Jong

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
EPA_Ecotox <- read.csv("./Data/Raw/ECOTOX_Neonicotinoids_Mortality_raw.csv", header=T)
LTER <- read.csv("./Data/Raw/NTL-LTER_Lake_ChemistryPhysics_Raw.csv", header=T)

#2
library(ggplot2)
my.theme <- theme_bw(base_size = 12) +
  theme(axis.text=element_text(color="gray0"), legend.position = "right")
theme_set(my.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
unique(EPA_Ecotox$Chemical.Name)

## [1] Imidacloprid Thiacloprid Thiamethoxam Acetamiprid Clothianidin
## [6] Dinotefuran Nitenpyram Nithiazine Imidaclothiz
## 9 Levels: Acetamiprid Clothianidin Dinotefuran ... Thiamethoxam

#4
YearForEachChem <- by(EPA_Ecotox, EPA_Ecotox$Chemical.Name, function(x) shapiro.test(x$Pub..Year))
sapply(YearForEachChem, print)

##
## Shapiro-Wilk normality test
##
## data: x$Pub..Year
## W = 0.90191, p-value = 5.706e-08
##
##
## Shapiro-Wilk normality test
##
## data: x$Pub..Year
## W = 0.69577, p-value = 4.287e-11
##
##
## Shapiro-Wilk normality test
##
## data: x$Pub..Year
## W = 0.82848, p-value = 8.83e-07
##
##
## Shapiro-Wilk normality test
##
## data: x$Pub..Year
## W = 0.88178, p-value < 2.2e-16
##
##
## Shapiro-Wilk normality test
##
## data: x$Pub..Year
## W = 0.68429, p-value = 0.00093
##
##
## Shapiro-Wilk normality test
##
## data: x$Pub..Year
## W = 0.79592, p-value = 0.0005686
##
##
## Shapiro-Wilk normality test
##
```

```

## data: x$Pub..Year
## W = 0.75938, p-value = 0.0001235
##
##
## Shapiro-Wilk normality test
##
## data: x$Pub..Year
## W = 0.7669, p-value = 1.118e-11
##
##
## Shapiro-Wilk normality test
##
## data: x$Pub..Year
## W = 0.7071, p-value < 2.2e-16

##          Acetamiprid          Clothianidin
## statistic 0.9019051          0.6957727
## p.value   5.705653e-08        4.286524e-11
## method    "Shapiro-Wilk normality test" "Shapiro-Wilk normality test"
## data.name "x$Pub..Year"          "x$Pub..Year"
##          Dinotefuran          Imidacloprid
## statistic 0.8284776          0.881784
## p.value   8.829849e-07        1.381875e-22
## method    "Shapiro-Wilk normality test" "Shapiro-Wilk normality test"
## data.name "x$Pub..Year"          "x$Pub..Year"
##          Imidaclothiz          Nitenpyram
## statistic 0.6842883          0.7959154
## p.value   0.0009299786        0.000568584
## method    "Shapiro-Wilk normality test" "Shapiro-Wilk normality test"
## data.name "x$Pub..Year"          "x$Pub..Year"
##          Nithiazine            Thiacloprid
## statistic 0.7593762          0.7668966
## p.value   0.0001235273        1.117773e-11
## method    "Shapiro-Wilk normality test" "Shapiro-Wilk normality test"
## data.name "x$Pub..Year"          "x$Pub..Year"
##          Thiamethoxam
## statistic 0.7070961
## p.value   1.571879e-16
## method    "Shapiro-Wilk normality test"
## data.name "x$Pub..Year"

#For each chemical, the publication years associated with it is not well-approximated by a normal distr

library(tidyverse)

## -- Attaching packages ----- tidyverse
## v tibble  2.0.1      v purrr   0.2.5
## v tidyr   0.8.2      v dplyr   0.7.8
## v readr   1.3.1      v stringr 1.3.1
## v tibble  2.0.1      v forcats 0.3.0

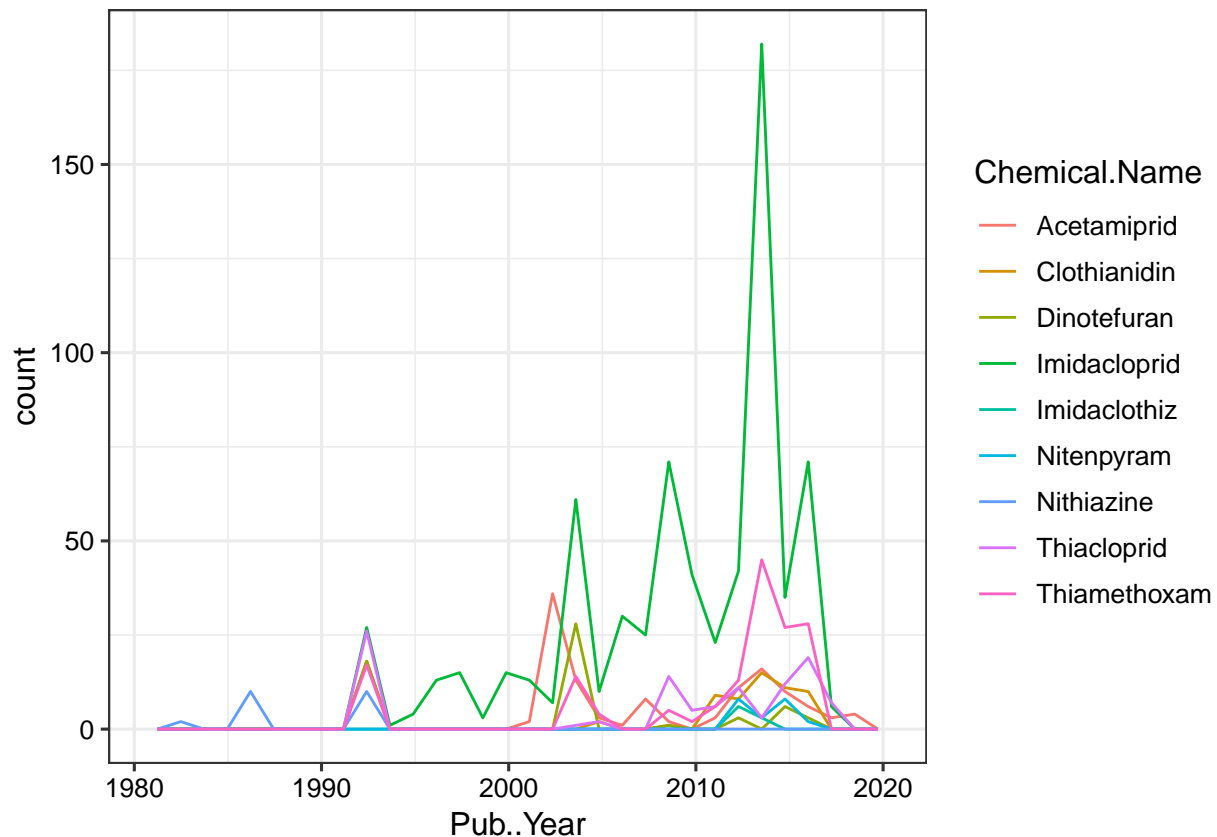
## -- Conflicts ----- tidyverse_confli
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

```

```
#Pipe version didnt work....
#YearForEachChem2 <- EPA_Ecotox %>%
# group_by(Chemical.Name) %>%
# shapiro.test(Pub..Year) #I get an error message: unused argument. Taylor help!

ggplot(EPA_Ecotox, aes(x=Pub..Year, col = Chemical.Name))+
  geom_freqpoly()

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
#5
YearForEachChem3 <- by(EPA_Ecotox, EPA_Ecotox$Chemical.Name, function(x) var(x$Pub..Year))
sapply(YearForEachChem3, print) #They do look like they have non-equal variance

## [1] 59.54809
## [1] 88.28601
## [1] 66.28521
## [1] 39.71249
## [1] 0.5277778
## [1] 2.414286
## [1] 12.81385
## [1] 89.35858
## [1] 53.12112

## Acetamiprid Clothianidin Dinotefuran Imidacloprid Imidaclothiz
## 59.5480937 88.2860052 66.2852133 39.7124914 0.5277778
## Nitenpyram Nithiazine Thiacloprid Thiamethoxam
```

```
##      2.4142857    12.8138528    89.3585804    53.1211180
```

```
#install.packages("car")  
library(car)
```

```
## Loading required package: carData
```

```
##
```

```
## Attaching package: 'car'
```

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

```
##
```

```
##      recode
```

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

```
##
```

```
##      some
```

```
leveneTest(Pub..Year ~ Chemical.Name, data = EPA_Ecotox) #we reject the hypothesis that the groups have
```

```
## Levene's Test for Homogeneity of Variance (center = median)
```

```
##      Df F value    Pr(>F)
```

```
## group    8  7.0203 4.243e-09 ***
```

```
##      1274
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

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

ANSWER: One-way ANOVA. Continuous response with one categorical explanatory variable with more than two categories

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
```

```
lm1 <- lm(Pub..Year ~ Chemical.Name, data = EPA_Ecotox)  
summary(lm1)
```

```
##
```

```
## Call:
```

```
## lm(formula = Pub..Year ~ Chemical.Name, data = EPA_Ecotox)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
```

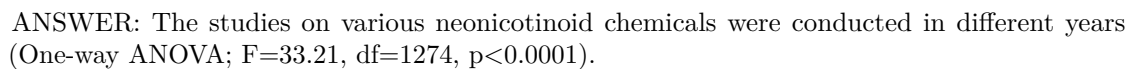
```
## -18.366  -3.993   1.889   4.889  13.441
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)  
## (Intercept)    2005.9926     0.6082 3298.222 < 2e-16 ***  
## Chemical.NameClothianidin     2.0479     1.0246   1.999  0.04584 *  
## Chemical.NameDinotefuran    -3.4333     1.1057  -3.105  0.00194 **  
## Chemical.NameImidacloprid     3.1181     0.6651   4.689 3.05e-06 ***  
## Chemical.NameImidaclothiz     6.4518     2.4412   2.643  0.00832 **  
## Chemical.NameNitenpyram       7.7216     1.6630   4.643 3.78e-06 ***  
## Chemical.NameNithiazine    -17.6290     1.6299 -10.816 < 2e-16 ***  
## Chemical.NameThiacloprid     1.6394     0.9190   1.784  0.07467 .  
## Chemical.NameThiamethoxam     4.3738     0.8261   5.295 1.40e-07 ***
```

```
#8
ggplot(EPA_Ecotox,aes(y=Pub..Year, x=Chemical.Name)) +
  geom_boxplot(aes(fill=Chemical.Name)) +
  xlab("Chemicals") +
  ylab("Published Year")+
  ggtitle("Publised Year Associated With Each Chemical")+
  theme(legend.position = "none")
```



6

- Only dates in July (hint: use the daynum column). No need to consider leap years.
- Only the columns: lakenname, 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
LTER2 <- LTER %>%
  select(lakenname, year4, daynum, depth, temperature_C) %>%
  filter(daynum %in% c(182:212)) %>%
  na.omit(lakenname, year4, daynum, depth, temperature_C)

#12
LTERAIC <- lm(temperature_C ~ year4+daynum+depth,data=LTER2)
step(LTERAIC) #It seems to be best to omit nothing

## 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 = LTER2)
##
## Coefficients:
## (Intercept)      year4      daynum      depth
##   -6.45556     0.01013     0.04134    -1.94726

lm2 <- lm(temperature_C ~ year4+daynum+depth,data=LTER2)
summary(lm2)

##
## Call:
## lm(formula = temperature_C ~ year4 + daynum + depth, data = LTER2)
##
## 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: $\text{temperature} = -6.45 + 0.01\text{year}_4 + 0.041\text{daynum} - 1.95\text{depth}$. This model explains 74.2% of the observed variance

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

```
#14
lm3 <- lm(temperature_C ~ depth * lakenname, data=LTER2)
summary(lm3)

##
## Call:
## lm(formula = temperature_C ~ depth * lakenname, data = LTER2)
##
## 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 ***
## lakennameHummingbird Lake  -2.4126     0.8379   -2.879  0.00399 **
## lakennamePaul Lake         0.6105     0.5983    1.020  0.30754
## lakennamePeter Lake        0.2998     0.5970    0.502  0.61552
## lakennameTuesday Lake     -2.8932     0.6060   -4.774 1.83e-06 ***
## lakennameWard Lake         2.4180     0.8434    2.867  0.00415 **
## lakennameWest Long Lake    -2.4663     0.6168   -3.999 6.42e-05 ***
## depth:lakennameCrampton Lake  0.8058     0.2465    3.268  0.00109 **
## depth:lakennameEast Long Lake  0.9465     0.2433    3.891  0.00010 ***
## depth:lakennameHummingbird Lake -0.6026     0.2919   -2.064  0.03903 *
## depth:lakennamePaul Lake     0.4022     0.2421    1.662  0.09664 .
## depth:lakennamePeter Lake    0.5799     0.2418    2.398  0.01649 *
## depth:lakennameTuesday Lake  0.6605     0.2426    2.723  0.00648 **
## depth:lakennameWard Lake    -0.6930     0.2862   -2.421  0.01548 *
## depth:lakennameWest 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 lakenname? How much variance in the temperature observations does this explain?

ANSWER: Yes, most combinations of depth and lakenname have significant impacts on temperature observations. This model explains about 78.6% of the variance in 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
library(viridis)

## Loading required package: viridisLite

library(RColorBrewer)
ggplot(LTER2, aes(x=depth, y=temperature_C, col=lakename)) +
  geom_point()+
  scale_y_reverse()+
  scale_color_brewer(palette = "YlGnBu")+
  ylab("Temperature (Celsius)")+
  ggtitle("Temperature by Depth Across Different Lakes")
```

