

# Using Citizen Science Species Occurrence Data to Study Biodiversity and Plan for the Future

Dylan Readell

February 2020

## Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Justo donec enim diam vulputate. Arcu vitae elementum curabitur vitae nunc sed velit dignissim sodales. Proin fermentum leo vel orci porta non pulvinar. Nascetur ridiculus mus mauris vitae ultricies leo integer. Facilisis magna etiam tempor orci eu. Vestibulum lorem sed risus ultricies. Euismod elementum nisi quis eleifend quam. Pulvinar sapien et ligula ullamcorper malesuada proin. Aenean pharetra magna ac placerat. Tempor id eu nisl nunc mi ipsum faucibus vitae. Et netus et malesuada fames ac turpis egestas sed. Duis ut diam quam nulla porttitor massa id neque aliquam. Euismod elementum nisi quis eleifend quam adipiscing vitae. Quam viverra orci sagittis eu volutpat odio.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Methods</b>	<b>3</b>
<b>3</b>	<b>Results</b>	<b>12</b>
<b>4</b>	<b>Discussion</b>	<b>13</b>
<b>5</b>	<b>Figures</b>	<b>14</b>

## List of Figures

1	Distribution of common bird species within the LA River Watershed region based on iNaturalist data. . . . .	14
2	The relative abundances of bird orders within the LA River Watershed region based on iNaturalist data. . . . .	15

## 1 Introduction

- To understand the biodiversity of an area, it is important to know the species of that area and how many of each species are present.
- Understanding the biodiversity of an area is useful for preparing for the future. Urban areas like Los Angeles are novel ecosystems. The data collected by citizen scientists like the data available on iNaturalist is extremely useful for identifying what species live in urban areas.
- For some species, seasonal distributions vary greatly. Database information can be utilized to understand how species come and go from an area over the course of a year. This seasonal distribution knowledge can be useful for identifying birds and insects that are still migrating through urban areas, or those that have changed courses to avoid the concrete jungles.
- Literature regarding urban ecosystem planning, especially in the face of climate change, is crucial for understanding the importance of this project. Furthermore, literature regarding the use of historic species occurrence records highlights the significance of why identifying the species in an area is critical for future planning.

## 2 Methods

The following code is a function that is planned to be able to analyze any species occurrence dataset. These datasets can be retrieved from database websites like Global Biodiversity Information Facility (GBIF), Integrated Digitized Biocollections (iDigBio), iNaturalist, and Paleobiology Database (PBDB). The goal of the function is to help users understand the distribution of a desired taxonomic rank throughout the seasons in an area of interest. Furthermore, the function will supply the user with a CSV output file containing the counts of each taxon depending on the taxonomic rank specified by the user. If the user wants to know the count of every family or species in the dataset, then that can be completed with this function.

Regular expressions are utilized in an inner function of the code presented below. The regular expression, along with user input, is used to extract all the occurrences of a desired year and tell the user how many occurrences were from that year.

The function will eventually do much more with the datasets to provide the user with as much information that is necessary for researching the biodiversity of an area.

```

1 import pandas as pd
2 import datetime
3 import matplotlib.pyplot as plt
4 from datetime import datetime
5 from collections import defaultdict
6 import os
7 import csv
8
9 def taxonomic_dataset_analysis():
10     '''
11     This function is designed to read in a species occurrence
12     dataset, analyze the data, and create useful outputs for
13     understanding taxon occurrences throughout seasons. The
14     function also outputs a CSV file that contains the counts of
15     each unique taxon in the dataset in respect to the taxonomic
16     rank specified by the user.
17
18     -----
19
20     inputfile: a species occurrence CSV dataset
21     date_col: the CSV column corresponding to the occurrence date
22     '''
23
24     inputfile = ''
25     inputfile = input('Enter the dataset filename: ')
26
27     # assert error for no filename specified
28     assert len(inputfile) != 0, 'No filename entered.'
29     # assert error that the filename specified is not a string
30     assert type(inputfile) == str, 'Filename is not a string: %r' %
31         inputfile
32
33     print('Beginning seasonal occurrence analysis.')
34
35     # read in the csv data file of species observations
36     data = pd.read_csv(inputfile)
37
38     # input for date column and assert that the date column exists
39     # within dataset
40     date_col = ''
41     date_col = input('Enter the observation date column header name
42     : ')
43     if date_col in data.columns:
44         print('Column found!')
45     assert date_col in data.columns, 'Date column header not found!'
46
47     def seasonal_occurrences(data, date_col, inputfile):
48         '''
49         This inner function uses the occurrence dataset along with
50         the date column specified to create a histogram displaying
51         specified taxonomic rank occurrences by season. The use of this
52         is to visualize the seasonal distribution of species or to
53         understand when people upload species observations to online
54         databases such as iNaturalist.
55
56         -----

```

```

44
45     User input is required to create the plot including the x-
46     label, the y-label, and the title.
47     '''
48
49     # extract the dates column to a list
50     dates = data[date_col].to_list()
51
52     # create an empty list for dates that are separated using
53     datetime
54     dates_separated = list()
55
56     # iterate through dates to separate them using datetime
57     for i in dates:
58         dates_separated.append(datetime.strptime(i, "%Y-%m-%d")
59     )
60
61     # create an empty list for observations by month
62     months = list()
63
64     # iterate through separated dates list to extract just the
65     month as an integer
66     for m in dates_separated:
67         months.append(m.month)
68
69     # create an empty list for observations by season
70     observations_by_season = list()
71
72     # iterate through each observation in the months list
73     # by using an if statement for each season
74     # if the month is in the season's numerical range
75     # the season is appended to the observation_by_season list
76     for x in months:
77         if x in range(1, 4):
78             observations_by_season.append('Winter')
79         if x in range(4, 7):
80             observations_by_season.append('Spring')
81         if x in range(7, 10):
82             observations_by_season.append('Summer')
83         if x in range(10, 13):
84             observations_by_season.append('Autumn')
85
86     # create an empty default dictionary with integers
87     # for the count of each observation by season
88     season_counts = defaultdict(int)
89
90     # iterate through seasons in observations_by_season list
91     # add a count for each season to the season_counts
92     dictionary
93     for season in observations_by_season:
94         season_counts[season] = season_counts[season] + 1
95
96     x_label = ''
97     x_label = input('Enter the label for the x-axis (something
98     about seasons): ')
99     y_label = ''
100    y_label = input('Enter the label for the y-axis (something

```

```

195         about number of occurrences): ')
196         figtitle = ''
197         figtitle = input('Enter the title for your figure: ')
198         print("Sit tight, we're making the figure...")
199
200         # plot a histogram of the season_counts dictionary
201         plt.bar(season_counts.keys(), season_counts.values())
202
203         # parameters for the figure
204         plt.xlabel(x_label)
205         plt.ylabel(y_label)
206         plt.title(figtitle)
207         plt.tight_layout()
208
209         figure_save = ''
210         figure_save = input('What do you want to call your figure
211 save? ')
212         plt.savefig(figure_save)
213         print("Figure successfully saved!")
214
215         print('Figure can be found at', os.path.abspath(figure_save
216 ))
217
218         print('Now creating taxon count CSV file.')
219
220     def taxon_counts(data, date_col, inputfile):
221         '''
222         This inner function is useful for ouputting a CSV file
223         containing the counts of a desired taxonomic rank. For example,
224         if the user specifies the family taxonomic rank, then a CSV
225         file will be created that gives the counts of each family in
226         the occurrence dataset.
227
228         -----
229
230         This function is useful to understand the diversity of an
231         area. An easy way to study the diversity of an area is to know
232         what species are the most common or rare.
233         '''
234
235         # input for the desired taxonomic rank column
236         # counts of each taxon placed in a dictionary
237         taxon_col = ''
238         taxon_col = input('Enter the taxon column header name: ')
239         counts_dic = (data[taxon_col].value_counts()).to_dict()
240
241         # input for the CSV save file
242         outputcsv = ''
243         outputcsv = input('What do you want to call your {} counts
244 save?'.format(taxon_col))
245
246         with open(outputcsv, 'w', newline = '') as csvfile:
247             writer = csv.writer(csvfile)
248             for species, counts in counts_dic.items():
249                 writer.writerow([species, counts])
250
251         print('CSV file can be found at', os.path.abspath(outputcsv

```

```

142     ))
143     # input for a user prompt of whether they want to know the
count of a specific taxon
144     answer = ''
145     answer = input('Do you want to know the count of any
specific taxon? Answer with Y or N.')
146     answer = answer.upper()
147
148     # answering yes brings user to another inner function
149     if answer == 'Y':
150
151         def counts_from_year(inputfile, date_col):
152             '''
153
154             This inner function utilizes a regular expression
and user input to extract occurrences from a user specified
year and output the number of occurrences in total from that
year. The regular expression searches the date column for all
matches to the user's year.
155
156             '''
157
158             dates = []
159
160             # reading in the data and appending the dates to a
list
161             with open(inputfile, 'r') as rf:
162                 data = csv.DictReader(rf, delimiter = ',')
163                 header = data.fieldnames
164                 for dic in data:
165                     dates.append(dic[date_col])
166
167             # input for the desired year
168             # generation of a regular expression based on the
year input
169             year = ''
170             year = input('Enter the year you want the number of
occurrences from: ')
171             pre_regex = '[^-\d]*'
172             regex = year + pre_regex
173             charRe = re.compile(regex)
174
175             year_specified = []
176             year_counts = defaultdict(int)
177
178             # appending each year match to a list
179             for date in dates:
180                 year_specified.append(charRe.match(date))
181
182             # creating a dictionary with the key as the year
and the value as the occurrences
183             for obs in year_specified:
184                 year_counts[obs] += 1
185
186             # this is required because each year match gets
added as a new key to the dictionary

```

```

187         # but the None key in the dictionary has a value
        that is equal to every year that
188         # did not match the user input
189         # the years that did not match is subtracted from
        the total length of the dataset
        count = len(dates) - year_counts.get(None)
190
191
192         print('The year {} had {} species observations.'.
        format(year, count))
193
194         counts_from_year(inputfile, date_col)
195
196         # if user answers no to prompt the analysis is completed
197         if answer == 'N':
198
199             print('Analysis completed.')
200
201         seasonal_occurrences(data, date_col, inputfile)
202         taxon_counts(data, date_col, inputfile)
203
204 taxonomic_dataset_analysis()

```



The following code is intended to produce a map of the distribution of several species at once. Figure 1 shows the distribution of the three most common bird species in the LA River Watershed region presented on iNaturalist. The base layer can be changed for any region, but for this example a shapefile of LA County was used. The function is modular meaning any number of species or other taxon can be used and plotted. A new point variable can be created for each species or taxon. The function is designed to work with any dataset that includes latitude and longitude coordinates. The species are chosen using a grep command and extracting them from the dataset. This function is useful since a website like iNaturalist does not display different species as different point objects simultaneously. The map produced by this function allows the user to visualize how different species are spatially distributed across a region. The function can be modified to extract and plot different aspects of the datasets, as well. For example, instead of extracting and plotting species or other taxon, a user may want to extract iNaturalist user identification numbers and plot the distribution of them.

```

1 library(rgdal)
2 library(ggplot2)
3 library(dplyr)
4
5 dist_map <- function(filename) {
6
7   setwd('~/Documents/UCLA/UCLA Winter 2020/EEB C177/')
8
9   # read in bird occurrence data with coordinates
10  mapdata <- read.csv(filename, stringsAsFactors = FALSE)
11
12  # set a variable for each species that is to be plotted
13  # any number of variables can be assigned for the map
14  mapdata_mallard <- mapdata[grepl('Anas platyrhynchos',
15                                mapdata$taxon_species_name), ]
16  mapdata_dove <- mapdata[grepl('Zenaida macroura',
17                               mapdata$taxon_species_name), ]
18  mapdata_sparrow <- mapdata[grepl('Passer domesticus',
19                                  mapdata$taxon_species_name), ]
20
21  # load in an LA County shapefile and set to a variable
22  county <- readOGR('DRP_COUNTY_BOUNDARY/DRP_COUNTY_BOUNDARY.shp',
23                  layer = 'DRP_COUNTY_BOUNDARY')
24
25  # useful for knowing how the county layer is being displayed
26  # for this shapefile, the projection used is NAD83
27  proj4string(county)
28
29  # find out what class the mapdata is
30  class(mapdata)
31
32  # creating a spatial point for the mapdata
33  # checking the mapdata class again
34  coordinates(mapdata_mallard) <- ~longitude+latitude
35  class(mapdata_mallard)
36  coordinates(mapdata_dove) <- ~longitude+latitude
37  class(mapdata_dove)

```

```

37 coordinates(mapdata_sparrow) <- ~longitude+latitude
38 class(mapdata_sparrow)
39
40 # checking if there is a projection for the mapdata
41 # it returns NA so there is not
42 proj4string(mapdata_mallard)
43 proj4string(mapdata_dove)
44 proj4string(mapdata_sparrow)
45
46 # manually setting the coordinate system for mapdata
47 # to be the same as the county shapefile
48 proj4string(mapdata_mallard) <- CRS("+proj=longlat +datum=NAD83")
49 proj4string(mapdata_dove) <- CRS("+proj=longlat +datum=NAD83")
50 proj4string(mapdata_sparrow) <- CRS("+proj=longlat +datum=NAD83")
51
52 # using spTransform function to project mapdata to the
53 # county shapefile
54 mapdata_mallard <- spTransform(mapdata_mallard, CRS(proj4string(
  county)))
55 mapdata_dove <- spTransform(mapdata_dove, CRS(proj4string(county)))
56 mapdata_sparrow <- spTransform(mapdata_sparrow, CRS(proj4string(
  county)))
57
58 # check to see if the mapdata and county projections match
59 # returns TRUE so they match
60 identical(proj4string(mapdata_mallard), proj4string(county))
61 identical(proj4string(mapdata_dove), proj4string(county))
62 identical(proj4string(mapdata_sparrow), proj4string(county))
63
64 # convert mapdata back to data.frame for ggplot to work with
65 mapdata_mallard <- data.frame(mapdata_mallard)
66 mapdata_dove <- data.frame(mapdata_dove)
67 mapdata_sparrow <- data.frame(mapdata_sparrow)
68
69 # changing the latitude and longitude names to x and y
70 # since we are working with x and y values now
71 names(mapdata_mallard)[names(mapdata_mallard)=="longitude"]<-"x"
72 names(mapdata_mallard)[names(mapdata_mallard)=="latitude"]<-"y"
73
74 names(mapdata_dove)[names(mapdata_dove)=="longitude"]<-"x"
75 names(mapdata_dove)[names(mapdata_dove)=="latitude"]<-"y"
76
77 names(mapdata_sparrow)[names(mapdata_sparrow)=="longitude"]<-"x"
78 names(mapdata_sparrow)[names(mapdata_sparrow)=="latitude"]<-"y"
79
80 # setting parameters for creating the map
81 map <- ggplot() +
82   # plot the county shapefile based on latitude and
  longitude
83   geom_polygon(data = county,
84               aes(x = long, y = lat, group =
  group),
85               fill = 'grey40',
86               color = 'grey90',
87               alpha = 1) +
88   # specify no x and y labels
89   # define the title of the map

```

```

90     labs(x = '',
91          y = '',
92          title = 'Distribution of Common Birds in the LA
River Watershed') +
93     # get rid of x and y ticks and text
94     # make the title bold, adjust its position
95     # adjust the position of the legend title and make it
bold
96     theme(axis.ticks.y = element_blank(),axis.text.y =
element_blank(),
97           axis.ticks.x = element_blank(),axis.text.x =
element_blank(),
98           plot.title = element_text(lineheight=.8, face="bold
", vjust=1, hjust = .5),
99           legend.position = 'right', legend.title = element_
text(face = 'bold')) +
100     # define the title of the legend
101     guides(fill = guide_legend(title = 'Species')) +
102     # plot the first species points
103     # setting the color and color as the name allows it to
appear on the legend
104     # and also allows us to change the appearance using
global setting
105     geom_point(data = mapdata_mallard,
106               aes(x = x, y = y, color = 'Mallard', size = '
Mallard')) +
107     # plot the second species points
108     geom_point(data = mapdata_dove,
109               aes(x = x, y = y, color = 'Mourning Dove',
size = 'Mourning Dove')) +
110     # plot the third species points
111     geom_point(data = mapdata_sparrow,
112               aes(x = x, y = y, color = 'House Sparrow',
size = 'House Sparrow')) +
113     # set the colors for each species
114     scale_color_manual(name = 'Species',
115                       values = c('Mallard'='red', 'Mourning
Dove'='blue', 'House Sparrow'='yellow')) +
116     # set the size of each point for each species
117     scale_size_manual(name = 'Species',
118                      values = c('Mallard'=.4, 'Mourning Dove
'=.4, 'House Sparrow'=.4)) +
119     # set the appearance to be square to avoid distortion
120     coord_equal(ratio = 1)
121
122 map
123 }
124
125 dist_map('birds.csv')

```

The following code is useful for producing a pie chart that displays the abundances of taxon present in a dataset. The iNaturalist bird dataset from the LA River Watershed was used in this code to display the relative abundances of each bird order present in that region. This visualization is shown in Figure 2.

```

1 library(ggplot2)
2
3 pie_chart <- function(filename) {
4
5   # set working directory as directory
6   setwd('~/Documents/UCLA/UCLA Winter 2020/EEB C177/')
7
8   # read in bird occurrence data
9   dataset <- read.csv(filename, stringsAsFactors=FALSE)
10
11   # extract all unique bird orders to a list
12   bird_orders <- list(dataset %>% select(taxon_order_name) %>%
13     distinct %>% arrange(taxon_order_name))
14
15   # creates a pie chart that displays the relative abundances of
16   # orders
17   pie <- ggplot(dataset,
18     aes(factor(1), fill=taxon_order_name)) +
19     # removing the x and y labels, defining the title
20     labs(x='', y='', title='Relative Abundances of Bird Orders in the
21       LA River Watershed') +
22     # making the title bold and adjusting its position
23     # making the legend title bold and placing it on the right
24     theme(plot.title = element_text(lineheight=.8, face="bold", vjust
25       =1, hjust = .5),
26       legend.position = 'right', legend.title = element_text(face
27         = 'bold')) +
28     # defining the title of the legend
29     guides(fill = guide_legend(title='Orders')) +
30     # setting the width of the slices
31     geom_bar(width=1) +
32     # changing chart to polar coordinates to display as pie chart
33     # instead of bars
34     coord_polar(theta='y')
35
36   pie
37 }
38
39 pie_chart('birds.csv')

```

### 3 Results

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Justo donec enim diam vulputate. Arcu vitae elementum curabitur vitae nunc sed velit dignissim sodales. Proin fermentum leo vel orci porta non pulvinar. Nascetur ridiculus mus mauris vitae ultricies leo integer. Facilisis magna etiam tempor orci eu. Vestibulum lorem

sed risus ultricies. Eiusmod elementum nisi quis eleifend quam. Pulvinar sapien et ligula ullamcorper malesuada proin. Aenean pharetra magna ac placerat. Tempor id eu nisl nunc mi ipsum faucibus vitae. Et netus et malesuada fames ac turpis egestas sed. Duis ut diam quam nulla porttitor massa id neque aliquam. Eiusmod elementum nisi quis eleifend quam adipiscing vitae. Quam viverra orci sagittis eu volutpat odio.

## 4 Discussion

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Justo donec enim diam vulputate. Arcu vitae elementum curabitur vitae nunc sed velit dignissim sodales. Proin fermentum leo vel orci porta non pulvinar. Nascetur ridiculus mus mauris vitae ultricies leo integer. Facilisis magna etiam tempor orci eu. Vestibulum lorem sed risus ultricies. Eiusmod elementum nisi quis eleifend quam. Pulvinar sapien et ligula ullamcorper malesuada proin. Aenean pharetra magna ac placerat. Tempor id eu nisl nunc mi ipsum faucibus vitae. Et netus et malesuada fames ac turpis egestas sed. Duis ut diam quam nulla porttitor massa id neque aliquam. Eiusmod elementum nisi quis eleifend quam adipiscing vitae. Quam viverra orci sagittis eu volutpat odio.

## 5 Figures

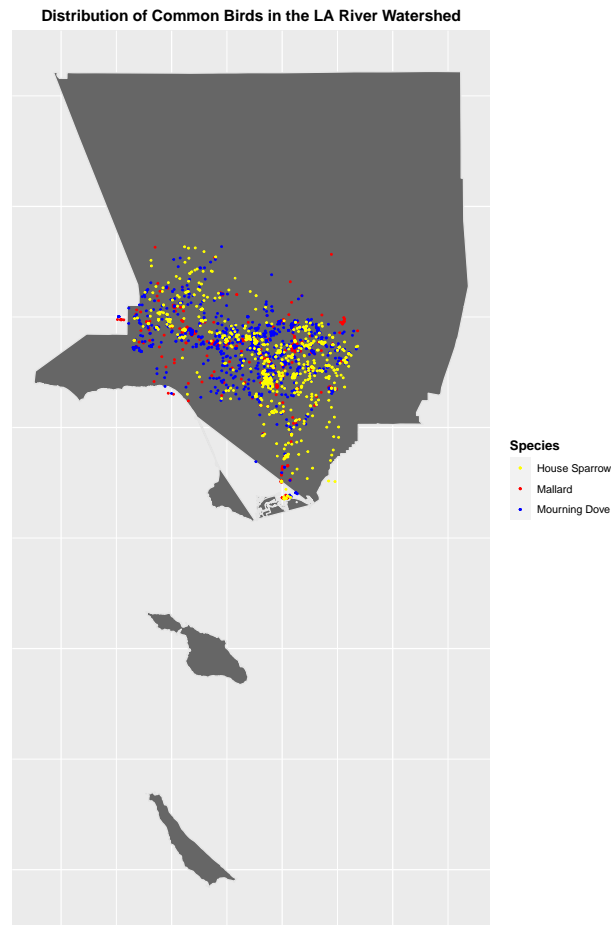


Figure 1: Distribution of common bird species within the LA River Watershed region based on iNaturalist data.

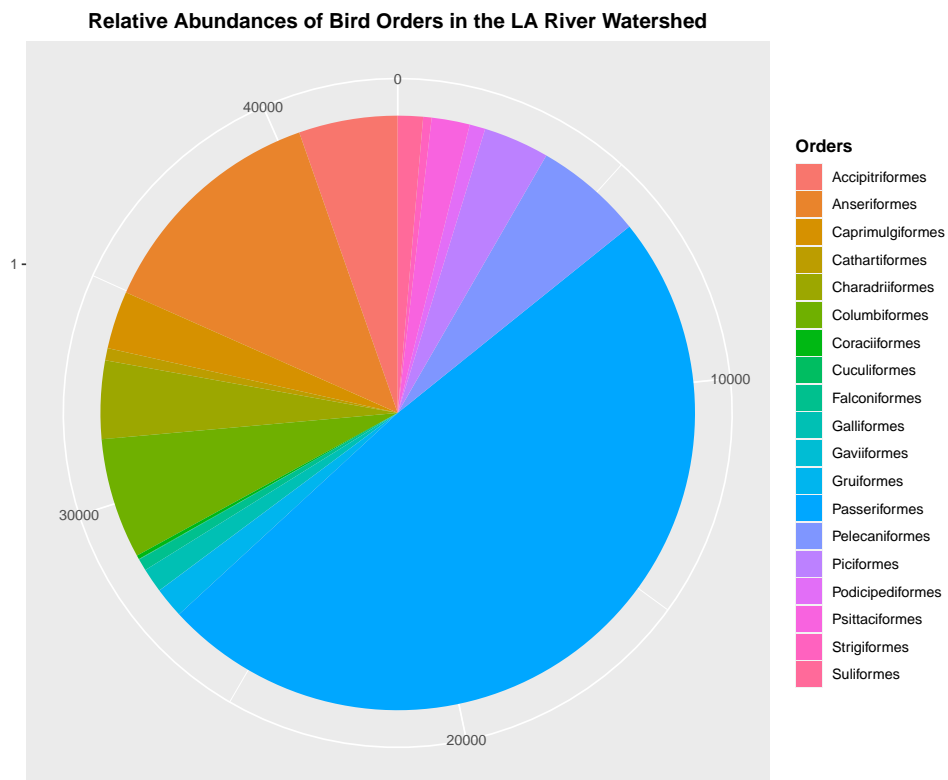


Figure 2: The relative abundances of bird orders within the LA River Watershed region based on iNaturalist data.