

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

1	Introduction	3
2	Methods	3
3	Results	8
4	Discussion	8
5	Figures	9

List of Figures

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
2 import pandas as pd
3 import datetime
4 import matplotlib.pyplot as plt
5 from datetime import datetime
6 from collections import defaultdict
7 import os
8 import csv
9
10 def taxonomic_dataset_analysis():
11     '''
12     This function is designed to read in a species occurrence
13     dataset, analyze the data, and create useful outputs for
14     understanding taxon occurrences throughout seasons. The
15     function also outputs a CSV file that contains the counts of
16     each unique taxon in the dataset in respect to the taxonomic
17     rank specified by the user.
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.

```

```

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

```

```

95     y_label = input('Enter the label for the y-axis (something
about number of occurrences): ')
96     figtitle = ''
97     figtitle = input('Enter the title for your figure: ')
98     print("Sit tight, we're making the figure...")
99
100    # plot a histogram of the season_counts dictionary
101    plt.bar(season_counts.keys(), season_counts.values())
102
103    # parameters for the figure
104    plt.xlabel(x_label)
105    plt.ylabel(y_label)
106    plt.title(figtitle)
107    plt.tight_layout()
108
109    figure_save = ''
110    figure_save = input('What do you want to call your figure
save? ')
111    plt.savefig(figure_save)
112    print("Figure successfully saved!")
113
114    print('Figure can be found at', os.path.abspath(figure_save
))
115
116    print('Now creating taxon count CSV file.')
117
118    def taxon_counts(data, date_col, inputfile):
119        '''
120        This inner function is useful for ouputting a CSV file
containing the counts of a desired taxonomic rank. For example,
if the user specifies the family taxonomic rank, then a CSV
file will be created that gives the counts of each family in
the occurrence dataset.
121
122        -----
123
124        This function is useful to understand the diversity of an
area. An easy way to study the diversity of an area is to know
what species are the most common or rare.
125        '''
126
127        # input for the desired taxonomic rank column
128        # counts of each taxon placed in a dictionary
129        taxon_col = ''
130        taxon_col = input('Enter the taxon column header name: ')
131        counts_dic = (data[taxon_col].value_counts()).to_dict()
132
133        # input for the CSV save file
134        outputcsv = ''
135        outputcsv = input('What do you want to call your {} counts
save?'.format(taxon_col))
136
137        with open(outputcsv, 'w', newline = '') as csvfile:
138            writer = csv.writer(csvfile)
139            for species, counts in counts_dic.items():
140                writer.writerow([species, counts])
141

```

```

142     print('CSV file can be found at', os.path.abspath(outputcsv
143 ))
144     # input for a user prompt of whether they want to know the
145     count of a specific taxon
146     answer = ''
147     answer = input('Do you want to know the count of any
148 specific taxon? Answer with Y or N.')
149     answer = answer.upper()
150
151     # answering yes brings user to another inner function
152     if answer == 'Y':
153         def counts_from_year(inputfile, date_col):
154             '''
155             This inner function utilizes a regular expression
156             and user input to extract occurrences from a user specified
157             year and output the number of occurrences in total from that
158             year. The regular expression searches the date column for all
159             matches to the user's year.
160
161             '''
162             dates = []
163
164             # reading in the data and appending the dates to a
165             list
166             with open(inputfile, 'r') as rf:
167                 data = csv.DictReader(rf, delimiter = ',')
168                 header = data.fieldnames
169                 for dic in data:
170                     dates.append(dic[date_col])
171
172             # input for the desired year
173             # generation of a regular expression based on the
174             year input
175             year = ''
176             year = input('Enter the year you want the number of
177 occurrences from: ')
178             pre_regex = '[^-\d]*'
179             regex = year + pre_regex
180             charRe = re.compile(regex)
181
182             year_specified = []
183             year_counts = defaultdict(int)
184
185             # appending each year match to a list
186             for date in dates:
187                 year_specified.append(charRe.match(date))
188
189             # creating a dictionary with the key as the year
190             and the value as the occurrences
191             for obs in year_specified:
192                 year_counts[obs] += 1
193
194             # this is required because each year match gets

```

```

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

```

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. 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.

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. 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.

5 Figures