

Bye Bye Birdie: Creating a Multiple Linear
Regression to Predict the Flight Initiation
Distance of the Dark-Eyed Junco (*Junco
hymenalis*)

Omer Lavian

February 29, 2020

Contents

1	Introduction	3
2	Methods	3
3	Results	5

1 Introduction

Anthropogenic change has undoubtedly had an effect on animal behavior [4]. One element of behavior that may be affected by human presence is flight initiation distance (FID). FID is the distance an animal initially moves away from an apparent threat. In birds, previous research has suggested that flight initiation distance is significantly associated with the animal's initial distance from an approaching threat [3]. Data on initial distance, as well as a variety of other factors, have been taken by the Yeh Lab at UCLA along with FID measurements in a sparrow known as the dark-eyed junco (*Junco hyemalis*). While these data have been analyzed, to my knowledge, they have not been combined in such a way as to allow one to predict the flight initiation distance of a junco given a certain set of factors. My goal, therefore, is to utilize the Yeh Lab's FID data to perform a multiple linear regression which will make possible predictions of FID for this species [2]. The overall goal of this project is to gain a better understanding of the ways in which different aspects of urbanization interact to affect flight initiation distance in dark-eyed juncos.

2 Methods

This section is incomplete, but overall the project involves quantifying the relationships between different variables and FID and then combining that into a multiple linear regression. To quantify the relationships, I started by making a function that prints the Pearson's correlation coefficient for interval variables that are input into the function and FID. I have also made a function that performs dummy coding on categorical variables with two categories [1].

Listing 1: Function comparing FID to other variables.

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import csv

def multiple_corrs_excel(filename, colx, coly=list):
    """
    Computes correlation coefficients of a list of
    columns in an Excel file to one particular column.
    Also creates a plot for each correlation.

    Example:
    >>> statdoer_version2("Spam_Data.xlsx", "Spam
    Consumption Frequency", ["Spam Quality", "Love for
    Eggs", "Frequency of Farts in Your General
    Direction"])
```

*The correlation coefficient of Spam Consumption
Frequency and Spam Quality is 0.75.*

PLOT APPEARS HERE

*The correlation coefficient of Spam Consumption
Frequency and Love for Eggs is 0.84.*

PLOT APPEARS HERE

*The correlation coefficient of Spam Consumption
Frequency and Frequency of Farts in Your General
Direction is 1.00.*

PLOT APPEARS HERE

'''

```
DataFrame = pd.read_excel(filename)
```

```
for item in colx:
    y = float(DataFrame[colx].corr(DataFrame[item]))
    print("The correlation coefficient of {} and {}
          is {:.{digits}f}.\n".format(colx, item,
          number = y, digits=2))
    corrplot = plt.plot(DataFrame[colx], DataFrame[
        item], "bo")
    plt.xlabel(colx)
    plt.ylabel(item)
    plt.title(colx + "_vs_" + item)
    plt.show()
```

Listing 2: Function that performs dummy coding on categorical variable with two categories

```
import pandas as pd
import csv
import re
```

```
def dummy_coder_excel():
    #(inputfile, variable, cat1, cat2, outputfile)
    '''
    Converts categorical variables with two categories
    into dichotomous variables and stores the output
    in a separate text file. This allows categorical
```

```

        variables to be used in a multiple linear
        regression.
    "category1" is stored as a 0, while "category2" is
    stored as a 1.

    This is a draft function. It will be modified to be
    more generalizable.

    '''

    #use regex to make function robust against typos; if
    wrote hop or fly in past of present tense or in
    upper or lowecase, then fxn should still work
    #no typos appear in this case, but future variability
    in data entry will have less of an effect on the
    function due to this use of regex
    Dataframe = pd.read_excel("FID.Data.xlsx")
    dummied_list = []
    #iterate through items in category and assign values
    of 0 or 1 to each entry
    for item in Dataframe["Flew.hop"]:
        matchtest1 = re.match("[Ff]*", str(item))
        matchtest2 = re.match("[Hh]*", str(item))
        if bool(matchtest1) == True:
            dummied_list.append(0)
        if bool(matchtest2) == True:
            dummied_list.append(1)
        else:
            continue
    #move codes to separate file for future use
    with open("FlewHop_Dummied.txt", "w+") as output:
        for item in dummied_list:
            output.write(str(item) + "\n")

```

3 Results

The end result of this project will (hopefully) be a multiple linear regression that combines all of the relationships with FID. Ideally, I would like to create a function that takes user input for a variety of factors and prints out a predicted FID.

References

- [1] Dummy variable (statistics) - Wikiversity.
- [2] Multiple Regression with Categorical Variables.
- [3] Daniel T. Blumstein. Flight-Initiation Distance in Birds Is Dependent on Intruder Starting Distance. *The Journal of Wildlife Management*, 67(4):852, October 2003.
- [4] B. B. M. Wong and U. Candolin. Behavioral responses to changing environments. *Behavioral Ecology*, 26(3):665–673, May 2015.