

Data Analysis using Python 3

Data Analyst is now a very lucrative job handle to have in the tech world, with the skills and patience of a data analyst one can analyze large chunks of Data in very short span of time which would have been otherwise impossible or would take up a long time.

Early January this year I did PH526x: Using Python for Research – (Harvard University) via edX, which taught me how to analyze .csv data files using Python 3. I scored well enough and passed the course. So, here I am presenting you, what I learned.



It was fun and exciting to work throughout the course and learn new stuff. My repository – [Github repo](#) holds the codes for study cases I had in the course. Professor Jukka-Pekka “JP” Onnela is an excellent professor who would make you learn stuff with his excellent teaching style. It was an awesome experience.

So, Let’s analyze one of the case called: Bird Migration Analysis [repo](#)

Aim: Track the movement of three gulls namely – Eric, Nico & Sanne

Dataset: <https://inbo.carto.com/u/lifewatch/datasets> ; used dataset – [csv](#)

Summary: One fascinating area of research uses GPS to track movements of animals. It is now possible to manufacture a small GPS device that is solar charged, so you don’t need to change batteries and use it to track flight patterns of birds. The

Dependencies :

- We will divide our case study into five parts:

- ## PART (1/5): Latitude and Longitude

The screenshot shows an Excel spreadsheet with the following columns: A1, B1, C1, D1, E1, F1, G1, H1, I1, J1, K1, L1, M1, N1, O1, P1, Q1, R1, S1, T1, U1. The data is organized into rows, with the first row being the header. The 'latitude' and 'longitude' columns are circled in red. The 'bird_name' column contains names like '0.15 Eric', '2.48386 Eric', etc. The spreadsheet is titled 'bird_tracking.csv'.

A1	B1	C1	D1	E1	F1	G1	H1	I1	J1	K1	L1	M1	N1	O1	P1	Q1	R1	S1	T1	U1
1	altitude	date_time	device_info	direction	latitude	longitude	speed_2d	bird_name												
2	71	2013-08-15	851	-150.47	49.41986	2.120733	0.15	Eric												
3	68	2013-08-15	851	-136.151	49.41988	2.120746	2.48386	Eric												
4	68	2013-08-15	851	166.7975	49.42031	2.120885	0.596657	Eric												
5	73	2013-08-15	851	32.76936	49.42036	2.120859	0.310161	Eric												
6	69	2013-08-15	851	45.19123	49.42033	2.120887	0.193132	Eric												
7	54	2013-08-15	851	-46.3445	49.42037	2.120804	2.904772	Eric												
8	57	2013-08-15	851	-56.3699	49.42035	2.120901	3.080584	Eric												
9	60	2013-08-15	851	-79.1702	49.42034	2.120809	2.198650	Eric												
10	59	2013-08-15	851	-57.8824	49.42032	2.120806	2.640076	Eric												
11	107	2013-08-15	851	119.6047	49.42029	2.121301	4.592213	Eric												
12	61	2013-08-15	851	-119.371	49.42034	2.12085	0.230217	Eric												
13	56	2013-08-15	851	-130.033	49.42086	2.121236	0.622415	Eric												
14	57	2013-08-15	851	-110.181	49.42096	2.121362	1.89783	Eric												
15	68	2013-08-15	851	-105.55	49.42573	2.152945	1.275343	Eric												
16	74	2013-08-15	851	-50.4790	49.42578	2.152959	1.634656	Eric												
17	78	2013-08-15	851	-75.1377	49.42575	2.153052	1.52315	Eric												
18	77	2013-08-15	851	132.0973	49.42574	2.153097	1.188486	Eric												
19	72	2013-08-15	851	-69.2591	49.42579	2.152993	1.006479	Eric												
20	74	2013-08-15	851	-7.47807	49.42582	2.153009	1.438506	Eric												
21	93	2013-08-15	851	71.0265	49.42942	2.162817	8.908179	Eric												
22	84	2013-08-15	851	156.082	49.43667	2.193876	0.850209	Eric												

The code:

```
bird_migration_trajectories_lat.long.py
1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import numpy as np
4
5  birddata = pd.read_csv("bird_tracking.csv") #make sure,you are in the right directory , check (>>>pwd)
6
7  # >>>birddata.info() #look at basic info abot the data frame
8  # >>>birddata.head() #look for first 5 rows in the data set
9  # >>>birddata.tail() #look for last 5 rows in the data set
10
11  bird_names = pd.unique(birddata.bird_name) #look at the unique names of the birds in the csv_file
12  # >>>print(bird_name)
13
14  ix = birddata.bird_name == "Eric" #storing the indices of the bird Eric
15  x,y = birddata.longitude[ix], birddata.latitude[ix]
16  plt.figure(figsize = (7,7))
17  plt.plot(x,y,"b.")
18  # >>>plt.show() #if you want to check trajectory of "Eric" only
19
20  ''' To look at all the birds trajectories, we plot each bird in the same plot '''
21  plt.figure(figsize = (7,7))
22  for bird_name in bird_names:
23      ix = birddata.bird_name == bird_name #storing the indices of the bird Eric
24      x,y = birddata.longitude[ix], birddata.latitude[ix]
25      plt.plot(x,y,".", label=bird_name)
26  plt.xlabel("Longitude")
27  plt.ylabel("Latitude")
28  plt.legend(loc="lower right")
29  plt.show()
30
```

In the code, Firstly we import the modules – pandas,matplotlib, and numpy. Then we import the csv file from the default directory (check default directory, `>>>pwd` or else change the path to the directory holding the csv file using `>>>cd directory_address`) into the variable birddata.

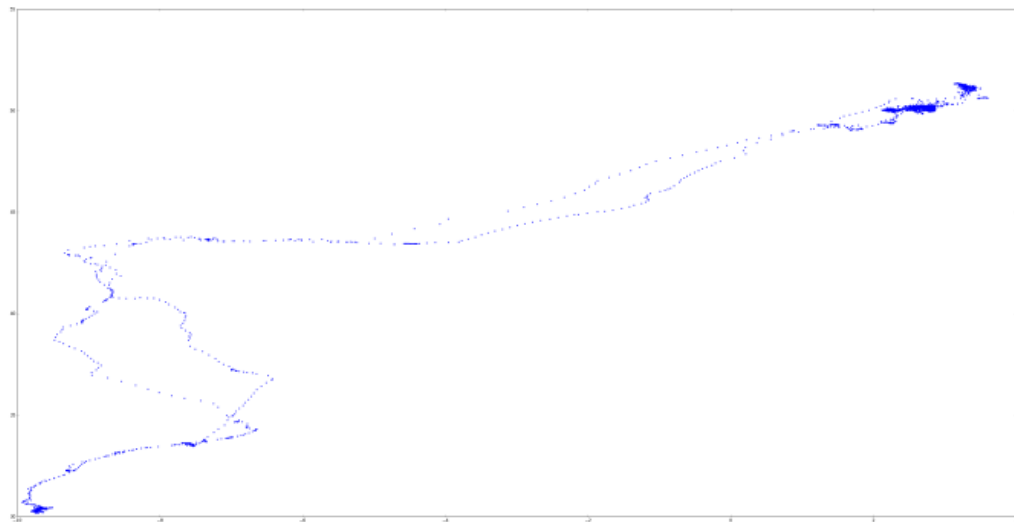
◆ `bird_names = pd.unique(birddata.bird_name)` is used to find all the unique bird names from the csv file and save it to the variable bird_names using pandas dataframe and [unique\(\)](#) function.

Next, we are going to select the latitude and longitude data only for the gull named “Eric”.We code,

◆ `ix = birddata.bird_name == “Eric”`
`x,y = birddata.longitude[ix], birddata.latitude[ix]`
`plt.figure(figsize = (7,7))`
`plt.plot(x,y,”b.”)`

Here, we are specifying the variable ix to contain the data of the column named "bird_name" from the csv file having the name of the bird as "Eric". Next, we are specifying x to hold longitude data and y to hold latitude data of "Eric". We use the matplotlib function, figure() to initialize its size as 7 x 7 and plot it using the plot() function ([learn matplotlib](#)). The parameters inside the function plot() i.e x, y and "b." are specifying to use longitude data along x axis, latitude along y and b=blue, . = circles in the visualization.

Output: [enlarged view](#)



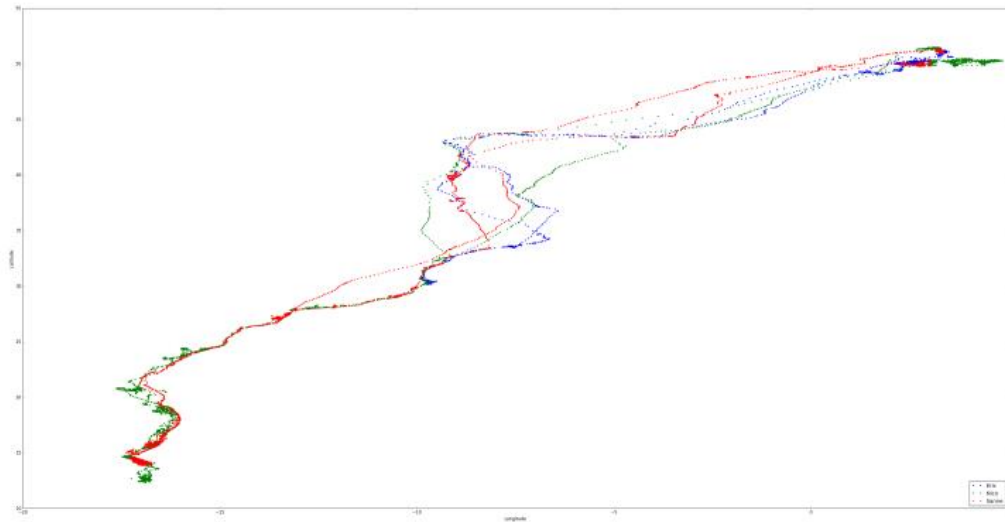
But Now, to look at all the bird's trajectories, we plot each bird in the same figure. We code,

```
◆ plt.figure(figsize = (7,7))
  for bird_name in bird_names:
      ix = birddata.bird_name == bird_name
      x,y = birddata.longitude[ix], birddata.latitude[ix]
      plt.plot(x,y,".", label=bird_name)
      plt.xlabel("Longitude")
      plt.ylabel("Latitude")
      plt.legend(loc="lower right")
      plt.show()
```

Here, we are plotting the location for all the three gulls namely Eric, Nico and Sanne. We create a 7 x 7 figure using plt.plot(figsize = (7,7)). We store unique bird's name in the variable ix, longitude and latitude data in the variables x and y respectively and we iterate over every data using a for loop. Next we, plot the data x and y using "." = circular marks and we add a label named "bird_data". We also use labels Longitude and Latitude along x and y axis respectively using xlabel() and ylabel() functions. legend() is used to locate the info bar in the plot, which is initialized to

lower right. Finally, we use the `show()` function to get the visualized data for all the three gulls.

Output : [enlarged view](#)



PART (2/5):_2D Speed Vs Frequency

In this second part of the case study, we are going to visualize 2D speed Vs Frequency for the gull named “Eric”.

The Code:

```
bird_migration_speed.py
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 ix = birddata.bird_name == "Eric" #storing the indices of the bird Eric
6 speed = birddata.speed_2d[ix]
7 # >>>plt.hist(speed) #[F.A.I.L.S] ,since we have non numeric numbers in the speed array
8 # >>>plt.hist(speed[:10]) #plot a histogram using the first 10 observations of speed [works]
9 # >>>np.isnan(speed) #>>>np.isnan(speed).any() #I am looking for any non number objects in the speed column using numpy's isnan function.
10 # >>>np.sum(np.isnan(speed)) # I find out the count of non numeric entries, False=0 & True =1 from isnan()
11
12 # >>>ind = np.isnan(speed) #return true if entry is not numerical or (not a NaN)
13 # >>>plt.hist(speed[~ind]) #we will include only those entries for which ind != True
14
15 plt.figure(figsize = (8,4))
16 # >>>speed = birddata.speed_2d[birddata.bird_name == "Eric"] #step 5 & 6 combined
17 ind = np.isnan(speed)
18 plt.hist(speed[~ind], bins=np.linspace(0,30,20), normed=True)
19 plt.xlabel(" 2D speed (m/s) ")
20 plt.ylabel(" Frequency ")
21 plt.show()
22
23
24 ...
25 We can also plot a similar histogram using the pandas module instead of pyplot.
26 The benefit of using pandas is that we do not have to deal with NaNs explicitly.
27 Instead, all of that happens under the hood.
28
29 NaNs - Not-a-Number
30
31 >>>birddata.speed_2d.plot(kind='hist', range=[0,30])
32 >>>plt.xlabel("2D speed")
33 >>>plt.savefig("hist_birdmig_speed.pdf")
34 ...
35
```

◆ `ix = birddata.bird_name == "Eric"`

`speed = birddata.speed_2d[ix]`

Here, we load bird data for the gull “Eric” into the variable `ix` and speed data of the same gull “Eric” into the variable `speed`.

◆ `plt.figure(figsize = (8,4))`

`ind = np.isnan(speed)`

`plt.hist(speed[~ind], bins=np.linspace(0,30,20), normed=True)`

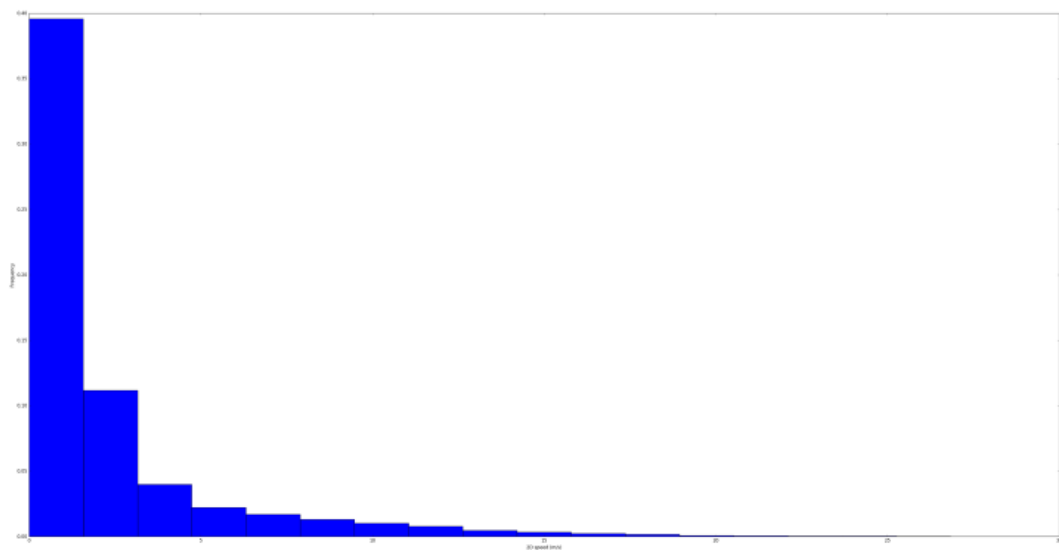
`plt.xlabel(" 2D speed (m/s) ")`

`plt.ylabel(" Frequency ")`

`plt.show()`

We plot a 8 x 4 figure and allot isnan speed data into `ind`. We find out the count of non numeric entries, `False=0` & `True =1` using the `isnan()` function. Next, we plot a histogram using the `hist()` function. The parameters `speed[~ind]` indicates that we will include only those entries for which `ind != True`, `bins=np.linspace(0,30,20)` indicates the bins along x axis will vary from 0 to 30 with 20 bins within them linearly spaced. Lastly, we plot 2D speed in m/s along x-axis and Frequency along y-axis using the `xlabel()` and `ylabel()` functions respectively and plot the data using `plt.show()`.

Output: [enlarged view](#)



PART (3/5): Time and Date

The third part is associated with date and time. We are going to visualize the time (in days) required by Eric to reach constant distances. If he requires same time to cover almost same distances, then the curve will be linear.

```
bird_migration_date.time.py
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import datetime
4
5 # >>>birddata.columns #check the columns of dataset
6 # >>>birddata.date_time[0:3] #check first few entries of date_time
7 # >>>datetime.datetime.today() #returns the current Date (yy-mm-dd) & time (h:m:s)
8
9 # >>>time_1 = datetime.datetime.today()
10 # >>>time_2 = datetime.datetime.today() #enter this code waiting few seconds
11 # >>>time_2-time_1 #you will get the measure of elapsed seconds as output, the resulting object is called date time time delta object
12
13 # >>>date_str = birddata.date_time[0]
14 # >>>date_str
15 # >>>date_str[:-3] # slices/removes the UTC +00 coordinated time stamps
16 # >>>datetime.datetime.strptime(date_str[:-3], "%Y-%m-%d %H:%M:%S") #the time stamp strings from date_str are converted to datetime object to be worked upon.
17
18 timestamps = []
19 for k in range(len(birddata)):
20     timestamps.append(datetime.datetime.strptime(birddata.date_time.iloc[k][:-3], "%Y-%m-%d %H:%M:%S"))
21
22     """
23     The next step for me is to construct a panda series object
24     and insert the timestamp from my Python list into that object.
25     I can then append the panda series as a new column in my bird data data frame.
26     """
27     birddata["timestamp"] = pd.Series(timestamps, index = birddata.index)
28     # >>>birddata.timestamp[4] - birddata.timestamp[3] #measure time difference between row 4 & 3
29
30     """
31     What I'd like to do next is to create a list that captures the amount of time
32     that has elapsed since the beginning of data collection.
33     """
34     times = birddata.timestamp[birddata.bird_name == "Eric"]
35     elapsed_time = [time-times[0] for time in times]
36
37     #But how can we measure time in certain units, like hours or days?
38     # >>>elapsed_time[1000]/datetime.timedelta(days=1) #output is the no of days have passed between these two points
39     # >>>len(elapsed_time) # check the length of entries
40     # >>>elapsed_time[19000]/datetime.timedelta(hours=1) #output is the no of hours have passed between these two points
41
42     plt.plot(np.array(elapsed_time)/datetime.timedelta(days=1))
43     plt.xlabel(" Observation ")
44     plt.ylabel(" Elapsed time (days) ")
45     plt.show()
```

We import the libraries matplotlib, pandas, and datetime.

```
▽ timestamps = []
   for k in range(len(birddata)):
```

```
       timestamps.append(datetime.datetime.strptime(birddata.date_time.iloc[k][:-3], "%Y-%m-%d %H:%M:%S"))
```

we create an empty list called **timestamps** and append date-time data of the birds to it.


```
>>>datetime.datetime.today() #returns the current Date (yy-mm-dd) & time (h:m:s).
```

```
>>>date_str[:3] #slices/removes the UTC +00 coordinated time stamps.
```

```
>>>datetime.datetime.strptime(date_str[:3], "%Y-%m-%d %H:%M:%S") ,the time-  
stamp strings from date_str are converted to datetime object to be worked upon.  
"%Y-%m-%d %H:%M:%S" is the Year-Month-Date and Hour-Minute-Second format"
```

The next step for us is to construct a panda series object and insert the timestamp from our Python list into that object. We can then append the panda series as a new column in my bird data, data frame.

```
▽ birddata["timestamp"] = pd.Series(timestamps, index = birddata.index)
```

What we'd like to do next is to create a list that captures the amount of time that has elapsed since the beginning of data collection.

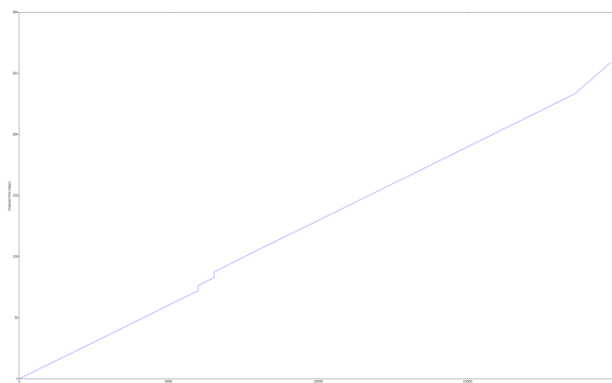
```
▽ times = birddata.timestamp[birddata.bird_name == "Eric"]  
  elapsed_time = [time-times[0] for time in times]
```

we calculated the elapsed time for the gull Eric

```
▽ plt.plot(np.array(elapsed_time)/datetime.timedelta(days=1))  
  plt.xlabel(" Observation ")  
  plt.ylabel(" Elapsed time (days) ")  
  plt.show()
```

We plot the observation(reference points at constant distances) along x axis vs elapsed time(in days) along y axis.We label our plot using xlabel() and ylabel() as Observation and Elapsed time (days) respectively along x and y axis.We Observe the curve.

Output : [enlarged view](#)



PART (4/5): Daily Mean Speed

We are going to visualize Daily mean speed of the gull named “Eric” for the total number of days of recorded flight.

```
bird_migration_daily_mean_speed.py
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import datetime
4 import numpy as np
5
6 #code from bird_migration_date.time *****
7 timestamps = []
8 for k in range(len(birddata)):
9     timestamps.append(datetime.datetime.strptime(birddata.date_time.iloc[k][-3], "%Y-%m-%d %H:%M:%S"))
10 birddata["timestamp"] = pd.Series(timestamps, index = birddata.index)
11
12 #code from bird_migration_date.time *****
13 data = birddata[birddata.bird_name == "Eric"]
14 times = data.timestamp
15 elapsed_time = [time-times[0] for time in times]
16 elapsed_days = np.array(elapsed_time)/datetime.timedelta(days=1)
17
18 next_day = 1
19 inds = []
20 daily_mean_speed = []
21 for (i,t) in enumerate(elapsed_days):
22     if t < next_day:
23         inds.append(i)
24     else:
25         daily_mean_speed.append(np.mean(data.speed_2d[inds]))
26         next_day += 1
27         inds = []
28
29 plt.figure(figsize = (8,6))
30 plt.plot(daily_mean_speed, "rs-")
31 plt.xlabel(" Day ")
32 plt.ylabel(" Mean Speed (m/s) ");
33 plt.show()
34
```

Up to line 16, we borrowed the code from part (3/5).

Next, we [enumerate](#) the [elapsed_days](#) and hold its returned tuple of an index and elapsed days in *i* and *t* respectively. Until and unless the elapsed day has not reached the next day we append the index to the empty list *inds*.

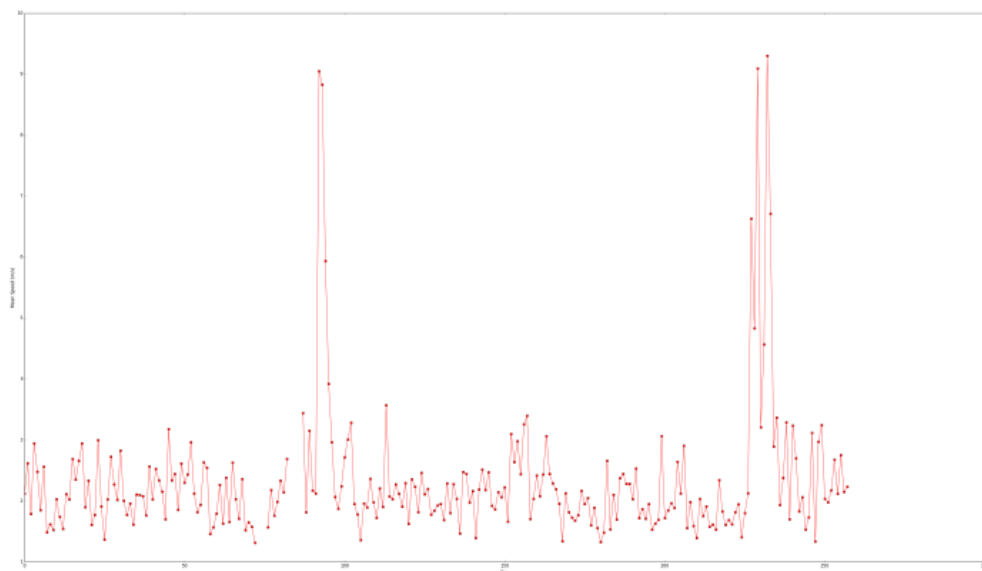
Or else If we reach the next day, we append speed data to daily mean speed and increment [next_day](#) by 1.

Lastly, we plot the figure of size 8 x 6

```
plt.plot(daily_mean_speed, "rs-")  
plt.xlabel(" Day ")  
plt.ylabel(" Mean Speed (m/s) ");  
plt.show()
```

In the plot code, r represent red, s represents square (for the turning points) and – represents the visualization style of the curve. We next, label the x and y axis using the xlabel() and ylabel() as "Day" and " Mean Speed (m/s) " respectively. Lastly, we show() the final plot.

Output: [enlarged view](#)



PART (5/5): Cartographic View

In this last part, i.e part 5, we are going to track the Birds over Political Map.

```
bird_migration_cartographic.py
1 import cartopy.crs as ccrs
2 import cartopy.feature as cfeature
3 import matplotlib.pyplot as plt
4
5 proj = ccrs.Mercator() #To move forward, we need to specify a specific projection that we're interested in using.
6
7 plt.figure(figsize=(10,10))
8 ax = plt.axes(projection=proj)
9 ax.set_extent((-25.0, 20.0, 52.0, 10.0))
10 ax.add_feature(cfeature.LAND)
11 ax.add_feature(cfeature.OCEAN)
12 ax.add_feature(cfeature.COASTLINE)
13 ax.add_feature(cfeature.BORDERS, linestyle=':')
14 for name in bird_names:
15     ix = birddata['bird_name'] == name
16     x,y = birddata.longitude[ix], birddata.latitude[ix]
17     ax.plot(x,y,'.', transform=ccrs.Geodetic(), label=name)
18 plt.legend(loc="upper left")
19 plt.show()
20
```

We import the cartopy and matplotlib module along with its salient libraries.

▽ `proj = ccrs.Mercator()`

To move forward, we need to specify a specific projection that we're interested in using. So we are using the cartopy [Mercator\(\)](#) function and initializing it to proj.

▽ `plt.figure(figsize=(10,10))`
`ax = plt.axes(projection=proj)`
`ax.set_extent((-25.0, 20.0, 52.0, 10.0))`
`ax.add_feature(cfeature.LAND)`
`ax.add_feature(cfeature.OCEAN)`
`ax.add_feature(cfeature.COASTLINE)`
`ax.add_feature(cfeature.BORDERS, linestyle=':')`

We plot a 10 x 10 figure and draw an axis with projection along the variable proj. Next, we add the political features like Land, Ocean, Coastline and borders into our plot. We automatically get the Political shapes and features mentioned above according to the gps locations (i.e, latitude and longitude) present in our data.

```

▽ for name in bird_names:
    ix = birddata['bird_name'] == name
    x,y = birddata.longitude[ix], birddata.latitude[ix]
    ax.plot(x,y,'.', transform=ccrs.Geodetic(), label=name)

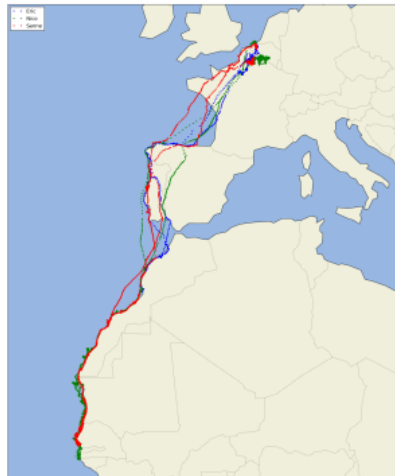
plt.legend(loc="upper left")
plt.show()

```

For every gull we plot its latitude and longitude data and use the Geodic() function to conform with the Geographical features.

Lastly, we visualize the mapped data of the gulls.

Output : [enlarged view](#)



Resources :

- [edX – Course](#)
- [Python Functions](#)
- [GitHub repo](#)
- [Tutorials Point – Python](#)
- [Python Doc](#)
- [Google Developers – Python](#)
- [Learnpythonthehardway.org](#)
- [Codecademy – Python](#)