



Data manipulation: Time Series & Geographical visualization

AAA-Python Edition



Plan

- 1- Date and time types
- 2- Time series basics
- 3- Date ranges, frequencies and shifting
- 4- Timezone
- 5- Periods and resampling
- 6- Rolling and expanding
- 7- Geographical visualization with basemap



1- Date and time types

Datetime module

```
1 # importing datetime module,  
2 import datetime as dt
```

```
1 # datetime type  
2 # to have the current time  
3 now = dt.datetime.now()
```

datetime.datetime(2018, 10, 27, 11, 37, 0, 875665)

year month day minute second

hour

microsecond

```
1 # to access the year  
2 now.year
```

2018

```
1 # datetime.timedelta type  
2 diff = dt.datetime.now() - dt.datetime(2018,10,26)
```

datetime.timedelta(1, 41841, 327812)

```
1 # difference in days  
2 diff.days
```

1

```
1 # timedelta type  
2 # operations with timedelta  
3 dt.datetime.now()+dt.timedelta(2)*2
```

datetime.datetime(2018, 10, 31, 11, 39, 23, 519331)

```
1 # string conversions  
2 dateS = dt.datetime(2018,11,5)  
3 # from datetime to string using str  
4 str(dateS)
```

'2018-11-05 00:00:00'

```
1 # from datetime to string using strftime  
2 dt.datetime.now().strftime("%Y/%m/%d, %H:%M")
```

'2018/10/27, 11:54'

datetime.datetime(2018, 10, 27, 13, 5)

```
1 # from string to datetime  
2 from dateutil.parser import parse  
3 parse("2018/10/27, 13h05")  
4
```

3



2- Time series basics

Indexing

```
# indexing series with datetime indexes
```

```
ser1 = S(range(3), index=[dt.datetime(2018,6,1),dt.datetime(2018,6,3),dt.datetime(2018,6,2)] )
```

ser1

2018-06-01	0
2018-06-03	1
2018-06-02	2

dtype: int64

```
1 # DatetimeIndex type
2 ser1.index
```

```
DatetimeIndex(['2018-06-01', '2018-06-03', '2018-06-02'], dtype='datetime64[ns]', freq=None)
```

ser2

2018-06-03	0
2018-06-04	1
2018-06-02	2

dtype: int64

```
1 # adding series values is done by aligning dates in the index
2 ser1 + ser2
```

The values to be added must have the same index.

ser2 + ser2

2018-06-01	NaN
2018-06-02	4.0
2018-06-03	1.0
2018-06-04	NaN

dtype: float64

If one of the indexes values is missing from one of the two series, the corresponding addition result will be a “**Nan**” value.



2- Time series basics

Selecting

```
1 # datetimeindex values type
2 times= ser1.index[1]
```

Timestamp('2018-06-03 00:00:00')

ser1

2018-06-01	0
2018-06-03	1
2018-06-02	2
dtype: int64	

```
1 # selecting a value by a timestamp
2
3 ser1[times]
```

1

```
1 # selecting a value by a datetime value
2 ser1[dt.datetime(2018,6,3)]
```

Selecting the same value

```
1 # selecting a value (or a slice) by a string: a date
2 ser1["6/3/2018"]
```

2018-06-03	1
dtype: int64	

```
1 # selecting a value ( or a slice)by a string: a date
2 ser1["20180603"]
```

The order of month, year, day is the same: year-month-day (from left to right)

The order is different from previous one: month-day-year (from left to right)

```
1 # selecting a value (or a slice )by a string: a year
2 ser1["2018"]
```

```
1 # selecting a slice by 2 strings
2 ser1["20180602 ":"20180603 "]
```

2018-06-03	1
2018-06-02	2
dtype: int64	

included

2018-06-01	0
2018-06-03	1
2018-06-02	2
dtype: int64	



3- Date ranges, frequencies and shifting

Date ranges

```
1 # creating dateitmeindex using a date range
2 dr = pd.date_range("2018-10-05", "2018-10-10")
```

included

included

```
DatetimeIndex(['2018-10-05', '2018-10-06', '2018-10-07', '2018-10-08',
               '2018-10-09', '2018-10-10'],
              dtype='datetime64[ns]', freq='D')
```

By default, the frequency is "one day"

```
1 # creating the same datetimeindex differently
2 # with start and period key argument
3 dr = pd.date_range(start="2018-10-05", periods=6)
```

included

```
1 # creating the same datetimeindex differently
2 # with start and period key argument
3 dr = pd.date_range(end="2018-10-10", periods=6)
```

Frequencies (using date_range)

```
1 # creating a datetimeindex with an hour frequency
2 dr2 = pd.date_range(start="2018-10-05", periods=6, freq="H")
```

By default, time starts at 00:00:00

```
DatetimeIndex(['2018-10-05 00:00:00', '2018-10-05 01:00:00',
               '2018-10-05 02:00:00', '2018-10-05 03:00:00',
               '2018-10-05 04:00:00', '2018-10-05 05:00:00'],
              dtype='datetime64[ns]', freq='H')
```

+ 1 hour

Frequencies (using resample)

```
1 # resample the series by a 12 hours frequency
2 ser2.resample("12h").mean()
```

New values

2018-06-02 00:00:00	2.0
2018-06-02 12:00:00	NaN
2018-06-03 00:00:00	0.0
2018-06-03 12:00:00	NaN
2018-06-04 00:00:00	1.0
Freq: 12H, dtype: float64	



3- Date ranges, frequencies and shifting

Frequencies (using date_range)

```
# creating a datetimeindex with 3 hours frequency
```

```
dr3 = pd.date_range(start="2018-10-05-15-00", periods =6, freq="3h")
```

```
DatetimeIndex(['2018-10-05 15:00:00+00:00', '2018-10-05 18:00:00+00:00',  
              '2018-10-05 21:00:00+00:00', '2018-10-06 00:00:00+00:00',  
              '2018-10-06 03:00:00+00:00', '2018-10-06 06:00:00+00:00'],  
              dtype='datetime64[ns, tzlocal()]', freq='3H')
```

+ 3 hours

```
# creating a datetimeindex with 30 minutes frequency
```

```
dr4 = pd.date_range(start="2018-10-05-15-00", periods =6, freq="30min")
```

```
DatetimeIndex(['2018-10-05 15:00:00+00:00', '2018-10-05 15:30:00+00:00',  
              '2018-10-05 16:00:00+00:00', '2018-10-05 16:30:00+00:00',  
              '2018-10-05 17:00:00+00:00', '2018-10-05 17:30:00+00:00'],  
              dtype='datetime64[ns, tzlocal()]', freq='30T')
```

+ 30 minutes

```
# creating a datetimeindex with month begin frequency
```

```
dr5 = pd.date_range(start="2018-10-05", periods =6, freq="MS")
```

First day of the month

Not included

```
DatetimeIndex(['2018-11-01', '2018-12-01', '2019-01-01', '2019-02-01',  
              '2019-03-01', '2019-04-01'],  
              dtype='datetime64[ns]', freq='MS')
```

[By Amina Delali]



3- Date ranges, frequencies and shifting

shifting

```
1 # shifting values without shifting indexes
2 ser1.shift(1)
```

The index didn't change

```
2018-06-01    NaN
2018-06-03     0.0
2018-06-02     1.0
dtype: float64
```

ser1

```
2018-06-01    0
2018-06-03    1
2018-06-02    2
dtype: int64
```

```
1 # shifting values with indexes
2 ser1.shift(1, freq="D")
```

The index will change

```
2018-06-02    0
2018-06-04    1
2018-06-03    2
dtype: int64
```

```
1 from pandas.tseries.offsets import Day, Week
2 from datetime import date
3
4 #shift using time offsets
5 now = dt.datetime.now()
6
7 print(now.strftime("Today = %Y/%m/%d"))
8 # timestamp created by shifting by 1 day from now
9 tomorrow = now + Day()
10 print(tomorrow.strftime("Tomorrow = %Y/%m/%d"))
```

offsets

```
Today = 2018/10/27
Tomorrow = 2018/10/28
```

```
1 # shift to the next week using rollforward
2 Week().rollforward(now).strftime("Next week = %Y/%m/%d")
```

'Next week = 2018/11/03'

```
1 # shift to the previous week using rollback
2 Week().rollback(now).strftime("Previous week = %Y/%m/%d")
```

'Previous week = 2018/10/20'



4- Timezone

Timezone conversions

```
1 # from a naive (without time zone) serieses to a localized one
2 print("ser1 timezone:",ser1.index.tz)
3 ser3=ser1.tz_localize("Africa/Algiers")
4 ser3.index.tz
```

ser1

```
2018-06-01    0
2018-06-03    1
2018-06-02    2
dtype: int64
```

ser3

```
2018-06-01 00:00:00+01:00    0
2018-06-03 00:00:00+01:00    1
2018-06-02 00:00:00+01:00    2
dtype: int64
```

```
2018-05-31 23:00:00+00:00    0
2018-06-02 23:00:00+00:00    2
2018-06-01 23:00:00+00:00    4
dtype: int64
```

ser1 timezone: None

<DstTzInfo 'Africa/Algiers' LMT+0:12:00 STD>

ser3 +ser4

<UTC>

```
# from a localized serieses to another time zone series
print("ser3 timezone:",ser3.index.tz)
ser4= ser3.tz_convert("America/New_York")
ser4.index.tz
```

```
# additionning 2 series with different time zones
# the result will be in UTC time zone
(ser3 + ser4).index.tz
```

ser3

```
2018-06-01 00:00:00+01:00    0
2018-06-03 00:00:00+01:00    1
2018-06-02 00:00:00+01:00    2
dtype: int64
```

ser3 timezone: Africa/Algiers

<DstTzInfo 'America/New_York' LMT-1 day, 19:04:00 STD>

ser4

```
2018-05-31 19:00:00-04:00    0
2018-06-02 19:00:00-04:00    1
2018-06-01 19:00:00-04:00    2
dtype: int64
```

```
# generating a localized date_range
```

```
ser5 = S(np.random.rand(1),index= pd.date_range(start= "2018-10-27",periods=1,tz="Africa/Algiers"))
```

```
2018-10-27 00:00:00+01:00    0.041752
Freq: D, dtype: float64
```

[By Amina Delali]



5- Periods and resampling

Periods

```
1 # Periods represent durations in time.  
2 # aPer will represent a duration equal to a month  
3 aPer = pd.Period("January 2018", freq="M")
```

Period('2018-01', 'M')

```
1 # adding values to a period will shift the period to that value * frequency  
2 aPer + 5
```

Period('2018-06', 'M')

```
PeriodIndex(['2018-01', '2018-02', '2018-03', '2018-04', '2018-05', '2018-06',  
            '2018-07', '2018-08', '2018-09', '2018-10', '2018-11', '2018-12'],  
            dtype='period[M]', freq='M')
```

```
1 # period_range will create a range of periods and also a PeriodIndex  
2 aPerR = pd.period_range("January 2018", "December 2018", freq="M")
```

```
1 # period indexes can be used as Series indexes  
2 # (they can also be created using pd.PeriodIndex)  
3 ser6 = S(range(12), index=aPerR)
```

2018-01	0
2018-02	1
2018-03	2
2018-04	3
2018-05	4
2018-06	5
2018-07	6
2018-08	7
2018-09	8
2018-10	9
2018-11	10
2018-12	11
Freq: M, dtype: int64	



5- Periods and resampling

Periods conversions

```
PeriodIndex(['2018-01', '2018-02', '2018-03', '2018-04', '2018-05', '2018-06',  
            '2018-07', '2018-08', '2018-09', '2018-10', '2018-11', '2018-12'],  
            dtype='period[M]', freq='M')
```

They didn't change

Before

Used to convert frequencies

After

```
1 # we can convert periods from frequency to another  
2 aPerR.asfreq("D", how="start")
```

```
PeriodIndex(['2018-01-01', '2018-02-01', '2018-03-01', '2018-04-01',  
            '2018-05-01', '2018-06-01', '2018-07-01', '2018-08-01',  
            '2018-09-01', '2018-10-01', '2018-11-01', '2018-12-01'],  
            dtype='period[D]', freq='D')
```

A day value added (but the temporal distance between the two periods remains the same)

```
1 # time series indexed by timestamps (datetime indexes) can be  
2 # converted to series indexed by periodindexes  
3 ser1.to_period("M")
```

Duplicated values

```
2018-06    0  
2018-06    1  
2018-06    2  
Freq: M, dtype: int64
```



5- Periods and resampling

Resampling

```
1 # resampling can be done in 2 directions
2 # from high to low frequency = downsampling
3 # ( from shorter to longer time duration)
4 ser6.resample("A-DEC").sum()
```

From 1 January to 31 December

In this case we can do aggregations: the short periods fall into longer ones

```
1 # from low to high frequency = upsampling
2 # ( from longer to shorter time duration)
3 ser6.resample("D").asfreq()
```

```
ser6.resample("D").asfreq().dropna()
```

No need for aggregation

Only the first day of each month has defined values (corresponding to the month value of the previous series. All the other values are Nan

```
2018    66
Freq: A-DEC, dtype: int64
```

```
2018-01-01    0.0
2018-01-02    NaN
2018-01-03    NaN
2018-01-04    NaN
2018-01-05    NaN
...
2018-12-31    NaN
Freq: D, Length: 365, dtype: float64
```

```
2018-01-01    0.0
2018-02-01    1.0
2018-03-01    2.0
2018-04-01    3.0
2018-05-01    4.0
2018-06-01    5.0
2018-07-01    6.0
2018-08-01    7.0
2018-09-01    8.0
2018-10-01    9.0
2018-11-01   10.0
2018-12-01   11.0
Freq: D, dtype: float64
```



6- Rolling and expanding

rolling

ser7

2018-01-01	0
2018-01-02	1
2018-01-03	2
2018-01-04	3
2018-01-05	4
2018-01-06	5
2018-01-07	6
2018-01-08	7
2018-01-09	8
2018-01-10	9

Freq: D, dtype: int64

```
# creating ser7 series using a period index
ser7 = S(range(10),index=pd.period_range("20180101",periods=10))
```

```
1 # rolling.mean() (in this case)will calculate the mean
2 # of each 2 consecutive values
3 ser7.rolling(2).mean()
```

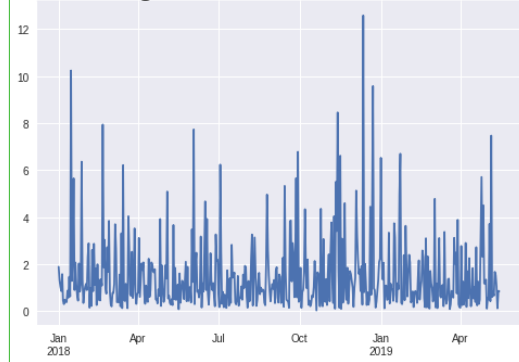
$$== (2+3)/2$$

$$== (6+7)/11$$

2018-01-01	NaN
2018-01-02	0.5
2018-01-03	1.5
2018-01-04	2.5
2018-01-05	3.5
2018-01-06	4.5
2018-01-07	5.5
2018-01-08	6.5
2018-01-09	7.5
2018-01-10	8.5

Freq: D, dtype: float64

Plotting of ser8 Series



Plotting of ser8.rolling(50).mean()



```
1 # rolling.mean in this case will calculate
2 # the mean of each 50 consecutive values
3 ser8.rolling(50).mean().plot()
```



6- Rolling and expanding

expanding

```
1 # expanding.mean (in this case )wil caculate
2 # the mean of 1, then 2, then 3, ..., until 10 consecutive values
3 ser7.expanding().mean()
```

ser7

2018-01-01	0
2018-01-02	1
2018-01-03	2
2018-01-04	3
2018-01-05	4
2018-01-06	5
2018-01-07	6
2018-01-08	7
2018-01-09	8
2018-01-10	9

Freq: D, dtype: int64

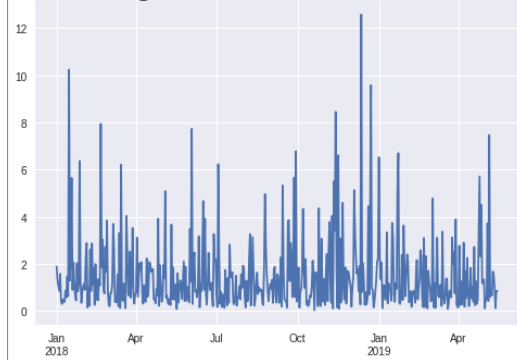
$$==(0+1+2)/3$$

$$==(0+1+2+3+4+5+6+7)/8$$

2018-01-01	0.0
2018-01-02	0.5
2018-01-03	1.0
2018-01-04	1.5
2018-01-05	2.0
2018-01-06	2.5
2018-01-07	3.0
2018-01-08	3.5
2018-01-09	4.0
2018-01-10	4.5

Freq: D, dtype: float64

Plotting of ser8 Series



Plotting of ser8.expanding().mean()



```
1 # apply expanding.mean to the ser8 series
2 ser8.expanding(50).mean().plot()
```



7- Geographical Visualization with basemap

Installation

- To install **basemap** in google colab, you have to:
 - Install the following libraries:
 - libproj-dev, proj-data, proj-bin

```
2 !apt-get install libproj-dev proj-data proj-bin
```

- Libgeos-dev

```
1 !apt-get install libgeos-dev
```

- Finally, install **basemap** from a GitHub repository:

```
1 !pip install https://github.com/matplotlib/basemap/archive/v1.1.0.tar.gz
```

- The import it from **matplotlib toolkit**:

```
1 # necessary imports
2 %matplotlib inline
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from mpl_toolkits.basemap import Basemap
```




7- Geographical Visualization with basemap

Usage

- Before projecting using **basemap**, you have to select the **projection** to use.
- In this example, we will use an “**orthographic projection**” (it shows half the globe at a time).

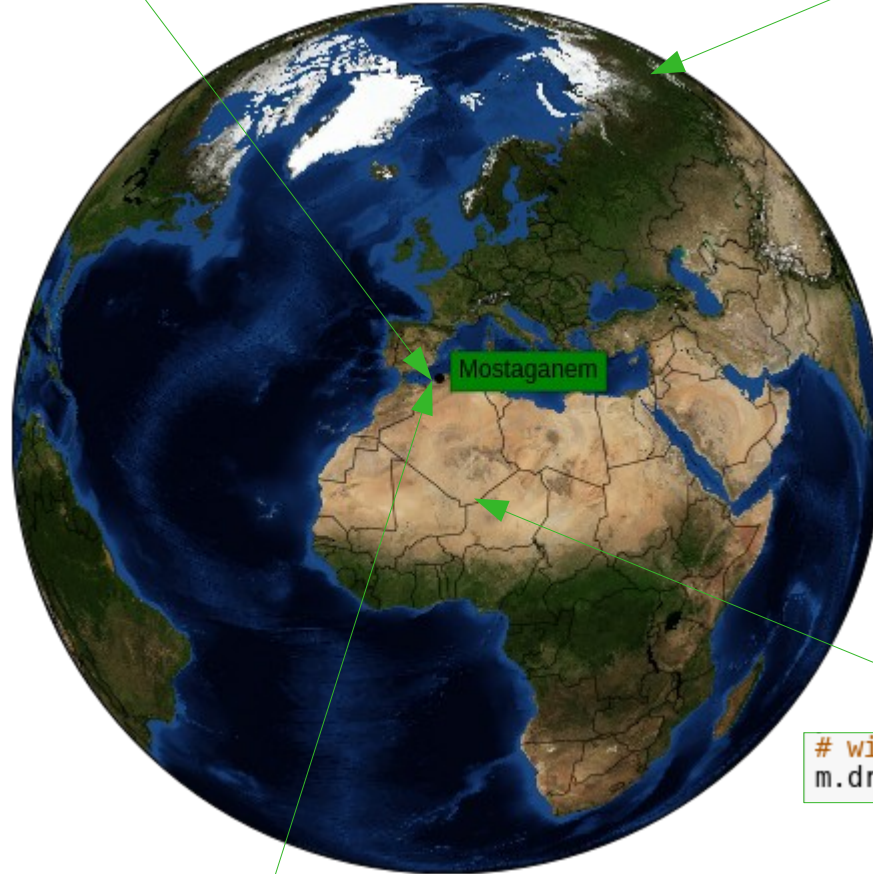
```
1 # you can specify your own country center latitude and longitude coordinates (in decimal degrees)
2 country_lat=28.0339
3 country_long=1.6596
4 # you can specify longitude and latitude for your own city (in decimal degrees)
5 city_lon =0.1401
6 city_lat = 36.0131
7
8 # a figure with a defined size
9 plt.figure(figsize=(8, 8))
10 # selecting orthographic projection
11 m = Basemap(projection='ortho', lat_0=country_lat, lon_0=country_long)
12 # display the blumarble image as map background
13 m.blumarble();
14 # convert the geographic coordinates to projection coordinates
15 x, y = m(city_lon, city_lat)
16 #plotting the mark corresponding to the city coordinates
17 plt.plot(x, y, 'ok', markersize=5)
18 # plotting a text showing the name of the city
19 plt.text(x+300000, y, "Mostaganem", bbox=dict(facecolor="green"), fontsize=12, color="black");
20 # will draw political country boundaries
21 m.drawcountries()
```




7- Geographical Visualization with basemap

Usage

```
4 # you can specify longitude and latitude for your own city (in decimal degrees)
5 city_lon = 0.1401
6 city_lat = 36.0131
12 # display the blumarble image as map background
13 m.blumarble();
```



```
# will draw political country boundaries
m.drawcountries()
```

[By Amina Delali]

```
#plotting the mark corresponding to the city coordinates
plt.plot(x, y, 'ok', markersize=5)
```



7- Geographical Visualization with basemap

Plotting

```
1 cities = pd.read_csv('AAA-Ped-Week3/A3P-w3-dz.csv')
2 # display only the first row
3 cities.head(1)
```

The data used for plotting

	city	lat	lng	country	iso2	admin	capital	population	population_proper
0	Algiers	36.763056	3.050556	Algeria	DZ	Alger	primary	3354000.0	1977663.0

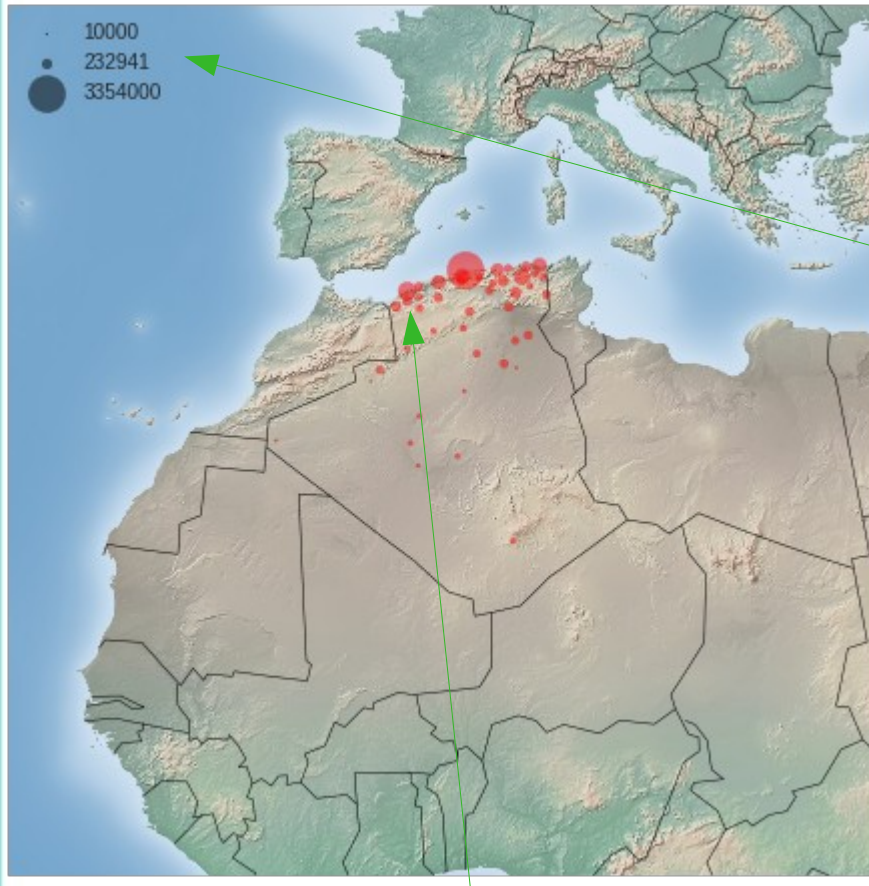
```
1 # extract latitudes and longitudes values from corresponding columns
2 lat = cities['lat'].values
3 lon = cities['lng'].values
4 # extract the population values
5 population = cities['population'].values
6
7 fig = plt.figure(figsize=(8, 8))
8 # use of , you have to specify the width and the height( in projections units : meters)
9 # or specify the four corner coordinates
10 m = Basemap(projection='lcc', lat_0=country_lat, lon_0=country_long, width=5E6, height=5E6)
11 # Project a shaded relief image onto the map
12 m.shadedrelief()
13 m.drawcountries(color='black')
14 # scatter city data, with size reflecting population
15 m.scatter(lon, lat, latlon=True, s= population/10000, c="red", alpha=0.5 )
16
17 # selecting 3 values: 10000, mean, maximum of population values as legend
18 legend_values= [10000, int(cities.population.mean()), int(cities.population.max())]
```

An other projection was used:
The Lambert conformal conic
projection



7- Geographical Visualization with basemap

Plotting



A shaded relief image projected on the map

```
12 m.shadedrelief()
```

```
# plotting the legend
for a in legend_values:
    plt.scatter([], [], c='k', alpha=0.5, s=a/10000, label=str(a) )
plt.legend(loc='upper left');
```

The lon, lat parameter represent the actual longitudes, latitudes coordinates

The sizes will be calculated from population size column

```
14 # scatter city data, with size reflecting population
15 m.scatter(lon, lat, latlon=True, s= population/10000, c="red", alpha=0.5 )
```

[By Amina Delali]



References

- Google Colab. Basemap-install. On-line at https://colab.research.google.com/drive/1_Xw_MEIrl0leP-v8vlmhUj6BLJLKmoV. Accessed on 31-10-2018.
- Wes McKinney. Python for data analysis: Data wrangling with Pandas, NumPy, and IPython. O'Reilly Media, Inc, 2018.
- simplemaps. Algeria cities database. On-line at <https://simplemaps.com/data/dz-cities>. Accessed on 31-10-2018.
- Jake VanderPlas. Python data science handbook: essential tools for working with data. O'Reilly Media, Inc, 2017.



Thank you!

FOR ALL YOUR TIME