

Project 1: Explore Weather Trends

Name: Tuan Bui

Date: April 28, 2020

1 . Outlines:

- What tools did you use for each step? (Python, SQL, Excel, etc) I used SQL to extract data and then used Python for manipulate and plot the data.
- How did you calculate the moving average? Suppose that we have a series x indexed by i . Let w be the window of the moving average, then the value of w moving average series at index n is computed as

$$MA_n = \frac{\sum_{i=n-w+1}^n x_i}{w}$$

- What were your key considerations when deciding how to visualize the trends? The key considerations that I took when visualize the trends are:
 - remember to take out the NaN values appear after calculating moving average.
 - use dash lines for trends so that they don't make the plot too dense.
 - first record in HoChiMinh city and global are not the same

2. Manipulating data

Load data

To load data from the database, use the following SQL query and then download the dataset as CSV file

```
select * from global_data
```

```
select * from city_list
```

```
select * from city_data
```

Now use Python to load the dataset form hard disk

```
In [1]: 1 import pandas as pd
        2 city_list = pd.read_csv('city_list.csv', header = 0)
        3 city_data = pd.read_csv('city_data.csv', header = 0)
        4 global_data = pd.read_csv('global_data.csv', header = 0)
```

```
In [2]: 1 # Extract temperature data for Ho Chi Minh city, Vietnam
        2 HCM = city_data[city_data['city']=='Ho Chi Minh City']
          [['year','avg_temp']]
        3 HCM.head()
```

Out[2]:

	year	avg_temp
25763	1825	27.11
25764	1826	NaN
25765	1827	NaN
25766	1828	NaN
25767	1829	NaN

Check if there are some missing values

```
In [3]: 1 print(sum(HCM['avg_temp'].isna()))
        2 print(sum(global_data['avg_temp'].isna()))

15
0
```

There are some missing values in HCM dataset so we impute the NaN using filling forward method

```
In [4]: 1 HCM.fillna(method = 'ffill',inplace = True)
```

Calculate moving average

```
In [5]: 1 # set window for moving average
        2 w = 10
        3 HCM['MA'] = HCM['avg_temp'].rolling(window=w).mean()
        4 global_data['MA'] =
          global_data['avg_temp'].rolling(window=w).mean()
        5 HCM = HCM.iloc[w:]
        6 global_data = global_data.iloc[w:]
```

3. Visualization: Line chart

```

In [7]: 1 import matplotlib.pyplot as plt
2 import numpy as np
3 from numpy.polynomial import Polynomial as P
4
5 f,(ax,ax2) = plt.subplots(2,1,sharex=True, facecolor='w',figsize=
(13,8))
6 ax.plot(HCM['year'],HCM['MA'],global_data['year'],global_data['MA']
)
7 ax.axvline(1834, 0, 30, label='HCM first record',ls = '--',c='r')
8
9 # Calculate and plot the trend lines
10 Z_g = np.polyfit(global_data['year'],global_data['MA'],1)
11 p_g = np.polyld(Z_g)
12 ax.plot(global_data['year'],p_g(global_data['year']),"b--")
13
14 Z_l = np.polyfit(HCM['year'],HCM['MA'],1)
15 p_l = np.polyld(Z_l)
16 ax.plot(HCM['year'],p_l(HCM['year']),"y--")
17
18 ax.legend(['HCM 10-year moving average','10-year moving
average','HCM first record'],loc='upper left')
19
20 ax.set_ylim(26,29)
21 ax.locator_params(axis='y', nbins=6)
22 ax2.plot(HCM['year'],HCM['MA'],global_data['year'],global_data['MA']
])
23 ax2.axvline(1834, 0, 30, label='HCM first record',ls = '--',c='r')
24
25 # Calculate and plot the trend lines
26 global_34 = global_data[global_data['year']>=1834]
27 Z_g = np.polyfit(global_34['year'],global_34['MA'],1)
28 p_g = np.polyld(Z_g)
29 ax2.plot(global_34['year'],p_g(global_34['year']),"b--")
30
31 Z_l = np.polyfit(HCM['year'],HCM['MA'],1)
32 p_l = np.polyld(Z_l)
33 ax2.plot(HCM['year'],p_l(HCM['year']),"y--")
34
35 # Min and Max
36 min_HCM = min(HCM['MA'])
37 max_HCM = max(HCM['MA'])
38 min_g = min(global_34['MA'])
39 max_g = max(global_34['MA'])
40
41 ax.axhline(min_HCM,xmin = 0.95,xmax=0.96,ls = '-',c='b')
42 ax.axhline(max_HCM,xmin = 0.95,xmax=0.96,ls = '-',c='b')
43 ax.plot((2018,2018), (min_HCM,max_HCM),ls = '--',c='b')
44 ax.text(2020,27.3,round(max_HCM-min_HCM,2),rotation=90)
45
46 ax2.axhline(min_g,xmin = 0.95,xmax=0.96,ls = '-',c='b')
47 ax2.axhline(max_g,xmin = 0.95,xmax=0.96,ls = '-',c='b')
48 ax2.plot((2018,2018), (min_g,max_g),ls = '--',c='b')
49 ax2.text(2020,8.7,round(max_g-min_g,2),rotation=90)
50
51 ax2.set_ylim(7,10)
52 ax2.locator_params(axis='y', nbins=6)

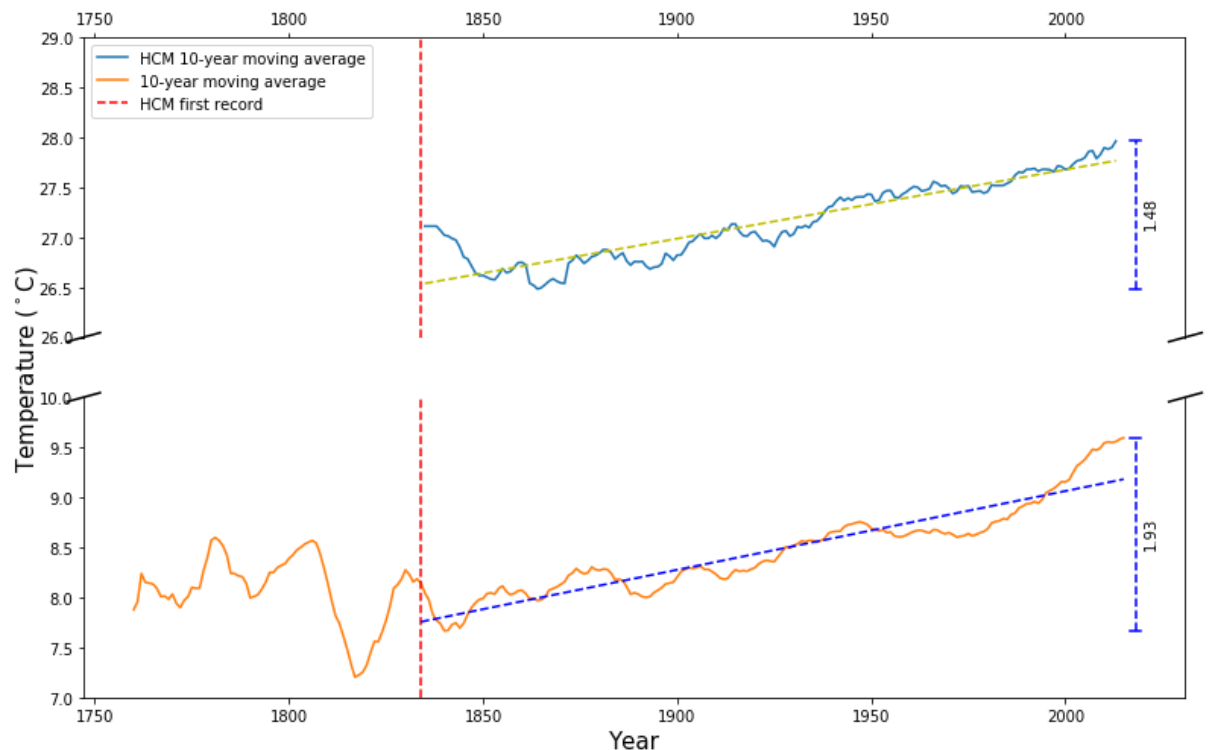
```

```

53 ax.spines['bottom'].set_visible(False)
54 ax2.spines['top'].set_visible(False)
55 ax.yaxis.tick_left()
56 ax.tick_params(labeltop=True)
57 ax.xaxis.tick_top()
58
59 d = .015 # how big to make the diagonal lines in axes coordinates
60 # arguments to pass to plot, just so we don't keep repeating them
61 kwargs = dict(transform=ax.transAxes, color='k', clip_on=False)
62 ax.plot((-d, +d), (-d, +d), **kwargs) # top-left diagonal
63 ax.plot((1 - d, 1 + d), (-d, +d), **kwargs) # top-right diagonal
64
65 kwargs.update(transform=ax2.transAxes) # switch to the bottom axes
66 ax2.plot((-d, +d), (1 - d, 1 + d), **kwargs) # bottom-left
    diagonal
67 ax2.plot((1 - d, 1 + d), (1 - d, 1 + d), **kwargs) # bottom-right
    diagonal
68
69 f.suptitle('Global and Ho Chi Minh temperature change over
    time',fontsize=20)
70
71 f.add_subplot(111, frameon=False)
72 # hide tick and tick label of the big axis
73 plt.tick_params(labelcolor='none', top=False, bottom=False,
    left=False, right=False)
74 plt.xlabel("Year",fontsize=15)
75 plt.ylabel("Temperature ($^\circ$C)",fontsize=15)
76
77 plt.show()

```

Global and Ho Chi Minh temperature change over time



4. Observations:

Similarities:

- The temperature in Ho Chi Minh and global are both increasing over time
- The fluctuation in both graphs seems to be very similar. At a particular period, temperatures were having the same behavior in both Ho Chi Minh city and global

Differences:

- Observing the slope of two trend lines we see that global's temperature is increasing faster than Ho Chi Minh city's
- Since the first record in Ho Chi Minh city to date, temperature change in Ho Chi Minh (1.48) is less than global (1.93)