

In [2]:

```
1 import dtale
2 import pandas as pd
3
4 data1 = pd.read_csv('temperature.csv',header=0)
5 data2 = pd.read_csv('population.csv',header=0)
6
7 temp = pd.DataFrame(data1)
8 pop = pd.DataFrame(data2)
```

In [ ]:

```
1 # Just for convenience - Population dataset
2 dtale.show(pop)
```

In [ ]:

```
1 # Just for convenience - Temperature dataset
2 dtale.show(temp)
```

In [3]:

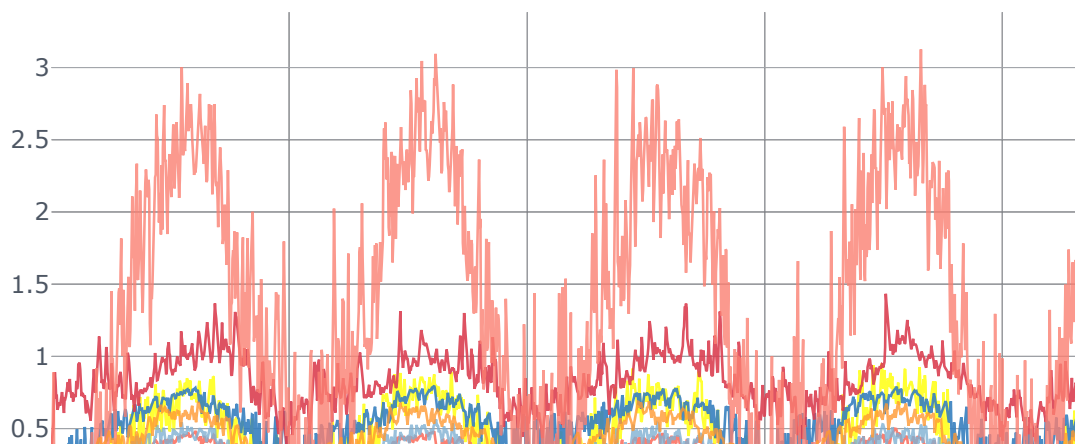
```
1 #####
2 ##### PART 1 #####
3 #####
4
5 import cufflinks as cf
6 from IPython.display import display,HTML
7 # making all charts public and setting a global theme
8 cf.set_config_file(sharing='public',theme='white',offline=True)
9
10 # Manually cleaning up few city names in the temperature dataset
11 # In order to correctly match cities when merging
12 temp['name'] = temp['name'].str.split('/').str[0]
13 temp.loc[temp['name']=='NYC','name'] = 'New York'
14 temp.loc[temp['name']=='Wash DC','name'] = 'Washington'
15 temp.loc[temp['name']=='Chicago O\'Hare','name'] = 'Chicago'
16 temp.loc[temp['name']=='St Louis','name'] = 'St. Louis'
17
18 # Population weight: (city_population / total_population)
19 pop['ratio'] = pop['population']/(pop['population'].sum())
20
21 # Merging population and temperature dataframes
22 merged = pop.merge(temp, left_on='City', right_on='name', how='outer').dropna()
23
24 # Multiplying the corresponding population weight to each city's temperature
25 merged['temp_mean_c'] = merged['temp_mean_c']*merged['ratio']
26 merged['temp_min_c'] = merged['temp_min_c']*merged['ratio']
27 merged['temp_max_c'] = merged['temp_max_c']*merged['ratio']
28
29 # Further cleaning the merged dataframe
30 merged = merged.drop(columns=['name', 'ratio', 'continent', 'country_code', 'cou
31 merged = merged.drop(columns=['State', 'population', 'Lon', 'Lat', 'station_code
32
33 # Change location_date column to datetime type
34 merged['location_date'] = pd.to_datetime(merged['location_date'])
35
36 # Creating a datetime dataframe in order to merge and identify which dates were
37 df = pd.DataFrame({'dates':pd.date_range('2015-01-01','2021-04-20')})
38
39 merged_min = merged[['City','location_date','temp_min_c']]
40 merged_mean = merged[['City','location_date','temp_mean_c']]
41 merged_max = merged[['City','location_date','temp_max_c']]
42
43 ##### Mininum temperature #####
44 merged_min = merged_min.set_index('location_date')
45 merged_min = merged_min.pivot_table(values='temp_min_c', index=merged_min.index,
46 merged_min = merged_min.merge(df, left_on='location_date', right_on='dates', how
47 merged_min = merged_min.set_index('dates').sort_values(by='dates', ascending=True
48 missing_min = merged_min # to be used for later
49 # for a missing data, render a reasonably predictive data by getting the average
50 merged_min = merged_min.where(merged_min.notnull(), other=(merged_min.fillna(met
51
52 ##### Maximum temperature #####
53 merged_max = merged_max.set_index('location_date')
54 merged_max = merged_max.pivot_table(values='temp_max_c', index=merged_max.index,
55 merged_max = merged_max.merge(df, left_on='location_date', right_on='dates', how
56 merged_max = merged_max.set_index('dates').sort_values(by='dates', ascending=True
57 missing_max = merged_max # to be used for later
58 # for a missing data, render a reasonably predictive data by getting the average
59 merged_max = merged_max.where(merged_max.notnull(), other=(merged_max.fillna(met
```

```

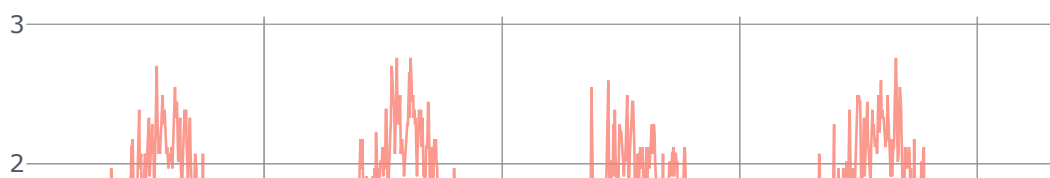
60
61 ##### Mean temperature #####
62 merged_mean = merged_mean.set_index('location_date')
63 merged_mean = merged_mean.pivot_table(values='temp_mean_c', index=merged_mean.in
64 merged_mean = merged_mean.merge(df, left_on='location_date', right_on='dates', h
65 merged_mean = merged_mean.set_index('dates').sort_values(by='dates', ascending=T
66 missing_mean = merged_mean # to be used for later
67 # for a missing data, render a reasonably predictive data by getting the average
68 merged_mean = merged_mean.where(merged_mean.notnull(), other=(merged_mean.fillna
69
70 # Plotting
71 ##### You can click the label on the legend to unsee the chosen line on the
72 merged_mean.ihplot(kind='line', title='Population-weighted Mean Temperature (°C)
73 merged_min.ihplot(kind='line', title='Population-weighted Minumum Temperature (°C
74 merged_max.ihplot(kind='line', title='Population-weighted Maximum Temperature (°C
75
76

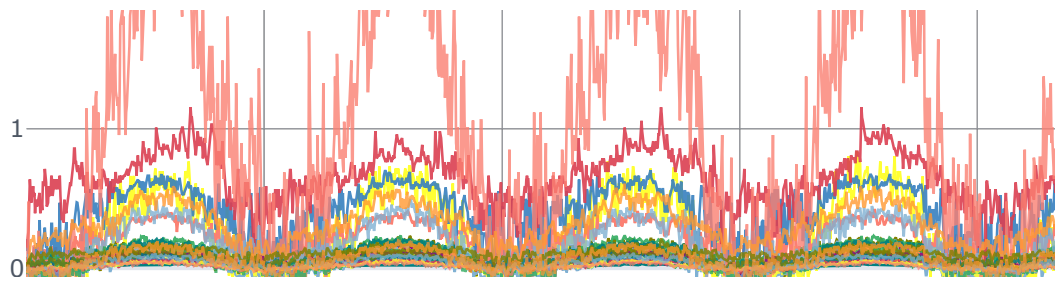
```

Population-weighted Mean Temperature (°C) timeseries by US c

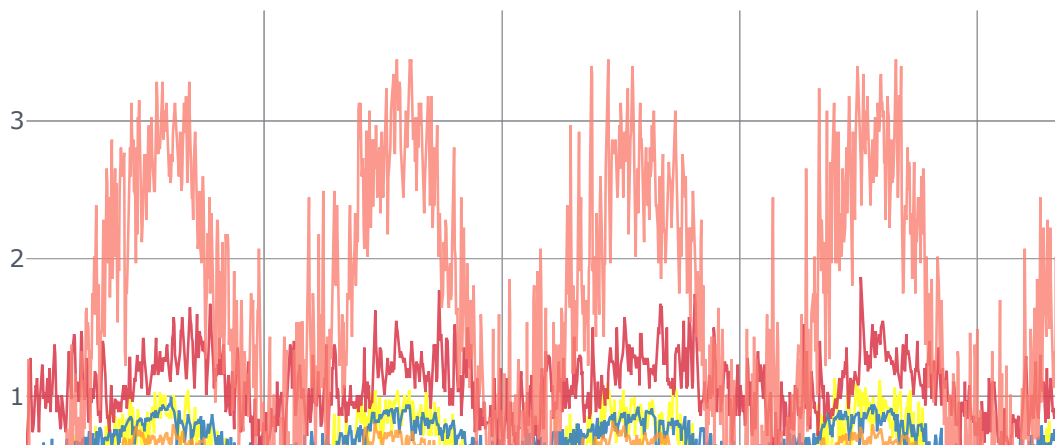


Population-weighted Minumum Temperature (°C) timeseries by





Population-weighted Maximum Temperature ( $^{\circ}\text{C}$ ) timeseries by



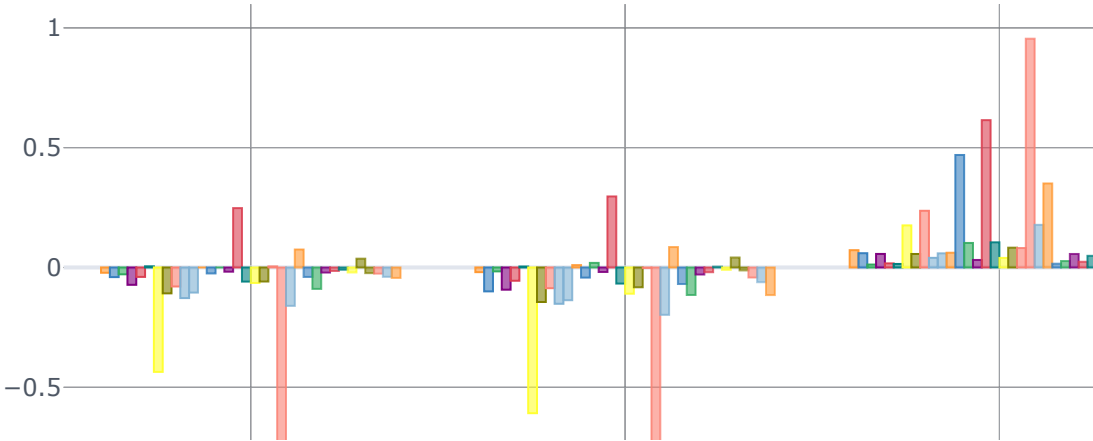
In [4]:

```
1 #####
2 ##### PART 2 #####
3 #####
4 ##### (1) Season #####
5
6 s_mean = merged_mean.reset_index(level='dates')
7 s_min = merged_min.reset_index(level='dates')
8 s_max = merged_max.reset_index(level='dates')
9
10
11 def get_season(row):
12     if row['dates'].month >= 3 and row['dates'].month <= 5:
13         return 'Spring'
14     elif row['dates'].month >= 6 and row['dates'].month <= 8:
15         return 'Summer'
16     elif row['dates'].month >= 9 and row['dates'].month <= 11:
17         return 'Fall'
18     else:
19         return 'Winter'
20
21 # Comparing mean temperature for each season
22 s_mean['Season'] = s_mean.apply(get_season, axis=1)
23 seasmean = s_mean.groupby(s_mean['Season']).mean()
24
25 # Comparing min temperature for each season
26 s_min['Season'] = s_min.apply(get_season, axis=1)
27 seasmin = s_min.groupby(s_min['Season']).min()
28 seasmin = seasmin.drop(columns='dates')
29
30 # Comparing max temperature for each season
31 s_max['Season'] = s_max.apply(get_season, axis=1)
32 seasmax = s_max.groupby(s_max['Season']).max()
33 seasmax = seasmax.drop(columns='dates')
34
35 # Plot
36 seasmean.iplot(kind='bar', title='Population-weighted Seasonal Mean Temperature
37 seasmin.iplot(kind='bar', title='Population-weighted Seasonal Min Temperature (°
38 seasmax.iplot(kind='bar', title='Population-weighted Seasonal Max Temperature (°
```

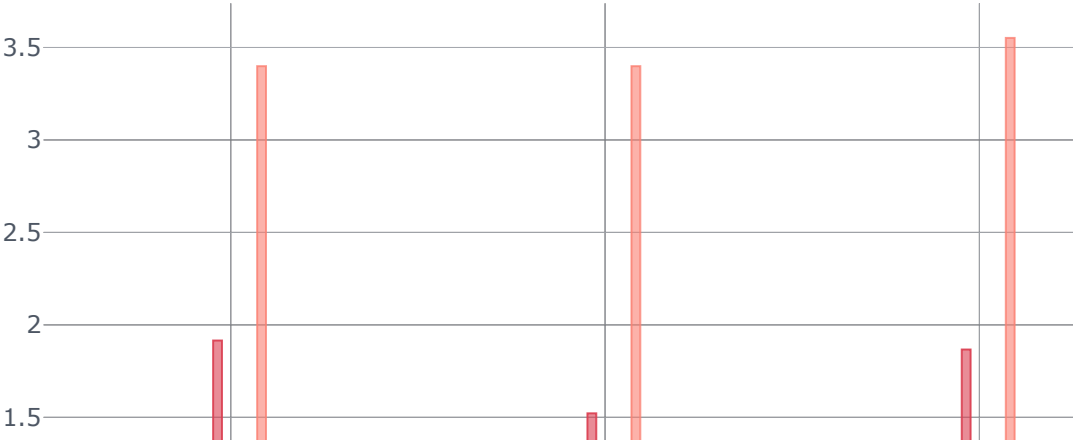
Population-weighted Seasonal Mean Temperature (°C) timeserie



Population-weighted Seasonal Min Temperature (°C) timeseries



Population-weighted Seasonal Max Temperature (°C) timeseries



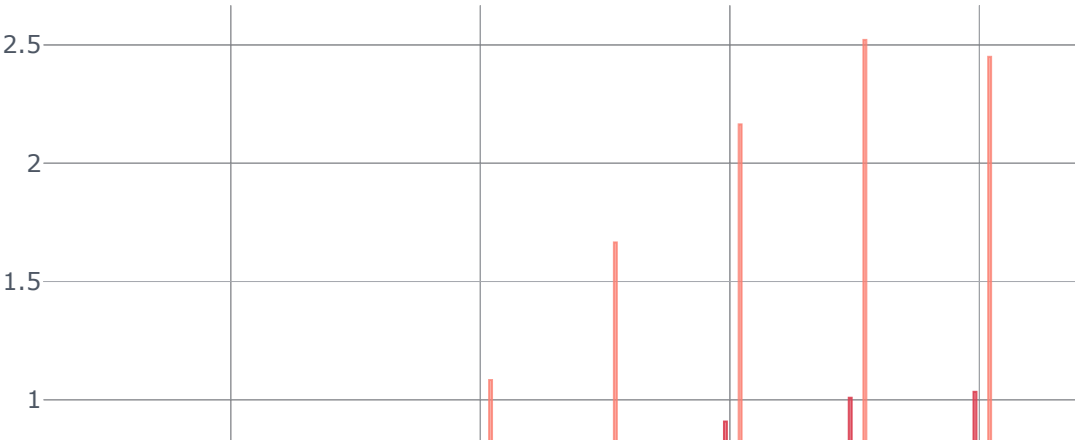


In [5]:

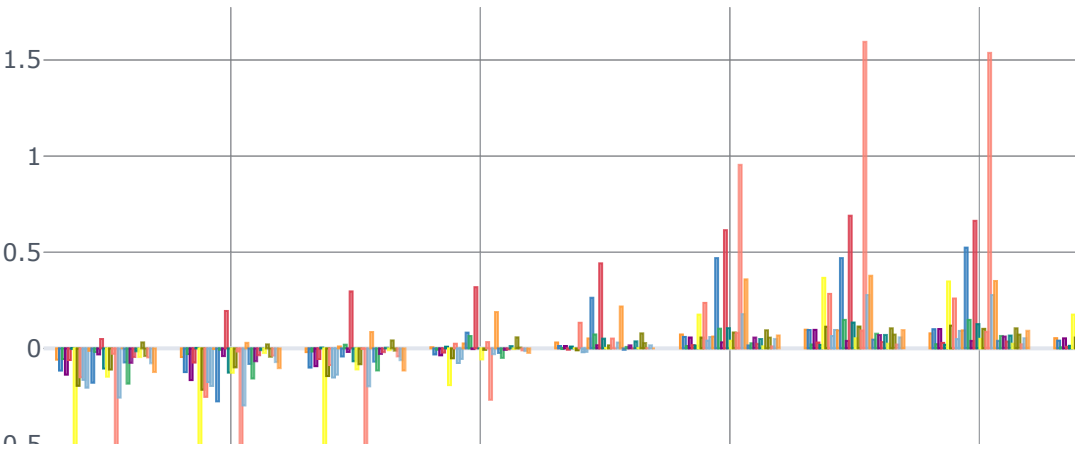
```
1 #####
2 ##### PART 2 #####
3 #####
4 ##### (2) Month #####
5
6 m_mean = merged_mean.reset_index(level='dates')
7 m_min = merged_min.reset_index(level='dates')
8 m_max = merged_max.reset_index(level='dates')
9
10
11 def get_month(row):
12     if row['dates'].month == 1:
13         return 1
14     elif row['dates'].month == 2:
15         return 2
16     elif row['dates'].month == 3:
17         return 3
18     elif row['dates'].month == 4:
19         return 4
20     elif row['dates'].month == 5:
21         return 5
22     elif row['dates'].month == 6:
23         return 6
24     elif row['dates'].month == 7:
25         return 7
26     elif row['dates'].month == 8:
27         return 8
28     elif row['dates'].month == 9:
29         return 9
30     elif row['dates'].month == 10:
31         return 10
32     elif row['dates'].month == 11:
33         return 11
34     else:
35         return 12
36
37 # Comparing mean temperature for each month
38 m_mean['month'] = m_mean.apply(get_month, axis=1)
39 monmean = m_mean.groupby(m_mean['month']).mean()
40 monmean = monmean.sort_index(ascending=True)
41
42 # Comparing min temperature for each month
43 m_min['month'] = m_min.apply(get_month, axis=1)
44 monmin = m_min.groupby(m_min['month']).min()
45 monmin = monmin.sort_index(ascending=True)
46 monmin = monmin.drop(columns='dates')
47
48 # Comparing max temperature for each month
49 m_max['month'] = m_max.apply(get_month, axis=1)
50 monmax = m_max.groupby(m_max['month']).max()
51 monmax = monmax.sort_index(ascending=True)
52 monmax = monmax.drop(columns='dates')
53
54 # Plot
55 monmean.iplot(kind='bar', title='Population-weighted Monthly Mean Temperature (°C)')
56 monmin.iplot(kind='bar', title='Population-weighted Monthly Min Temperature (°C)')
57 monmax.iplot(kind='bar', title='Population-weighted Monthly Max Temperature (°C)')
58
```



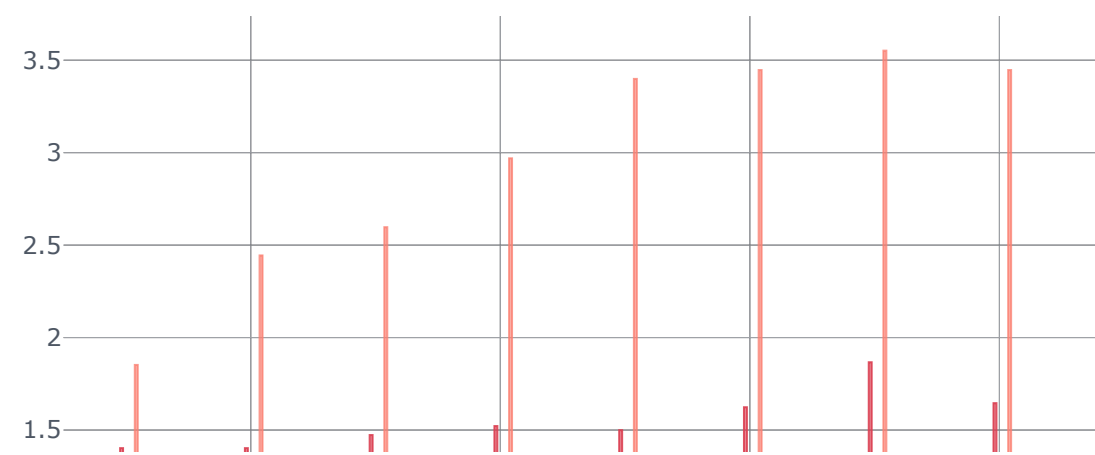
Population-weighted Monthly Mean Temperature (°C) timeseries



Population-weighted Monthly Min Temperature (°C) timeseries



Population-weighted Monthly Max Temperature (°C) timeseries



In [6]:

```
1 #####
2 ##### PART 2 #####
3 #####
4 ##### (3) Missing Data #####
5
6 # Dichotomizing missing/non-missing values
7 mizz_mean = missing_mean.isna()
8 # Plotting - Each scatter marker displays the date and the name of the city of
9 mizz_mean = mizz_mean.plot(kind='scatter', title='Missing data (marked true)',
```

## Missing data (marked true)

