

Experiment No:12

Aim: Located dataset (e.g., sample_weather.txt) for working on weather data which reads the text input files and finds average for temperature, dew point and wind speed.

Theory: Analysis of Weather data using Pandas, Python, and Seaborn:



The [most recent post](#) on this site was an analysis of how often people cycling to work actually get trained on in different cities around the world. You can check it out [here](#).

The analysis was completed using data from the [Wunderground](#) weather website, Python, specifically the Pandas and [Seaborn](#) libraries. In this post, I will provide the Python code to replicate the work and analyze information for your own city. During the analysis, I used [Python Jupyter notebooks](#) to interactively explore and cleanse data; there's a simple setup if you elect to use something like [the Anaconda Python distribution](#) to install everything you need.

If you want to skip data downloading and scraping, all of the data I used is available to [download here](#).

Scraping Weather Data

Wunderground.com has a "Personal Weather Station (PWS)" network for which fantastic historical weather data is available – covering temperature, pressure, wind speed and direction, and of course rainfall in mm – all available on a per-minute level. Individual stations can be examined at specific URLs, for example [here](#) for station "IDUBLIND35".

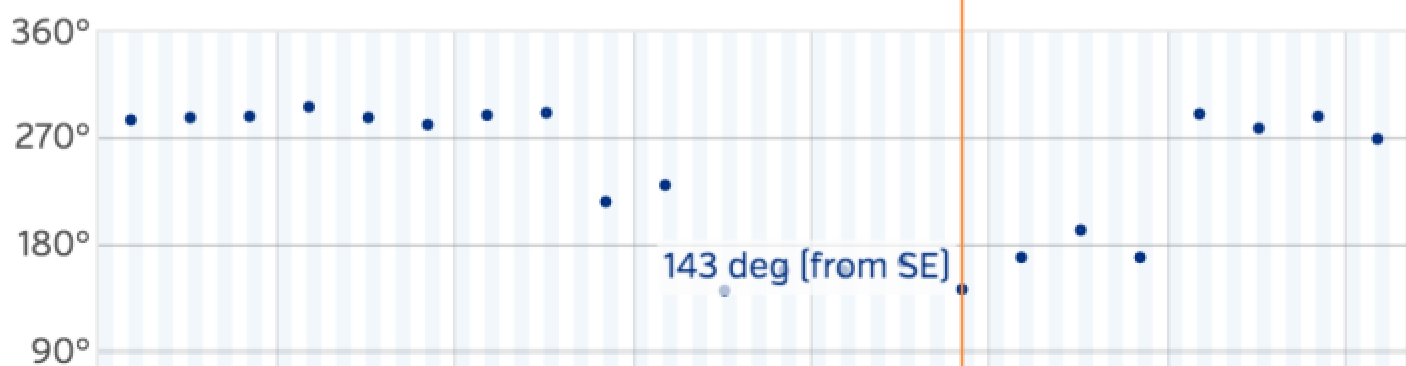
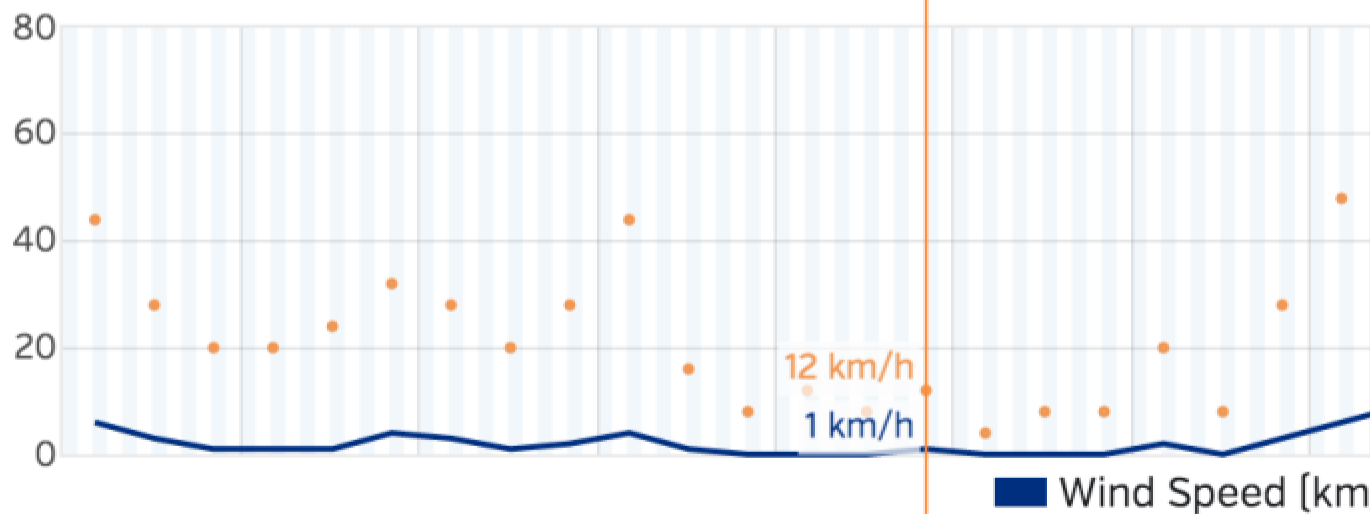
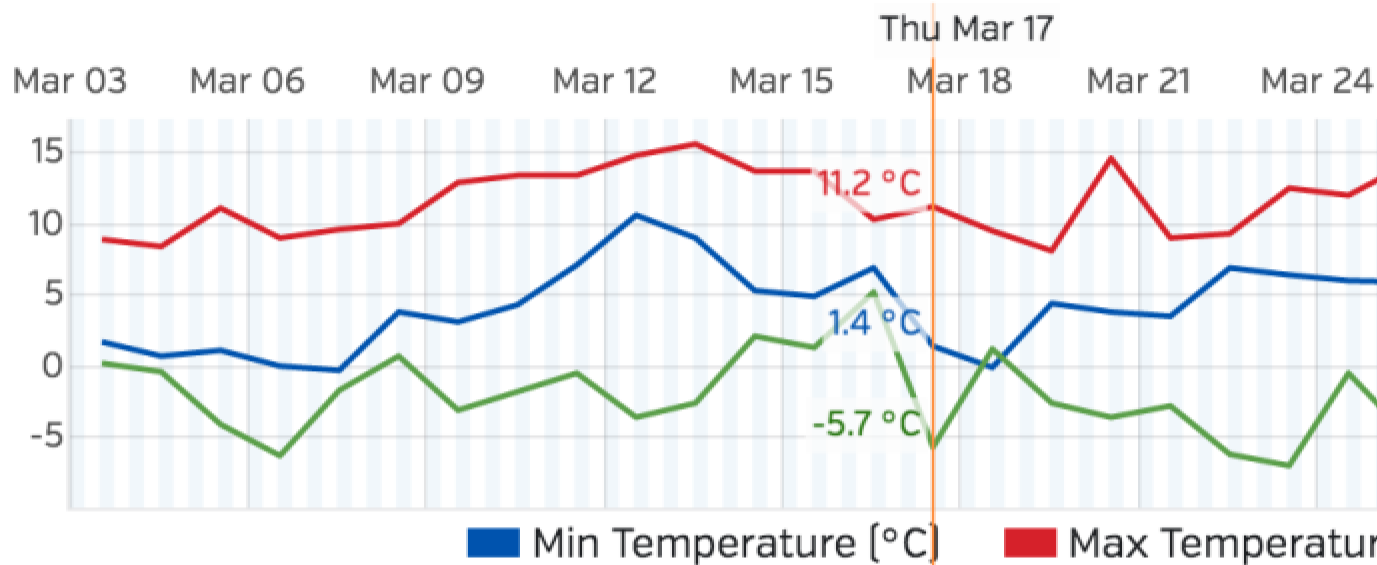
There's no official API for the PWS stations that I could see, but there is a very good [API for forecast data](#). However, CSV format data with hourly rainfall, temperature, and pressure information can be downloaded from the website with some simple Python scripts.

Graphs

Table

Weather History Graph

March 3, 2016 - April 2, 2016



Wunderground have an excellent site with interactive graphs to look at weather data on a daily, monthly, and yearly level. Data is also available to download in CSV format, which is great for data science purposes.

```
import requests

import pandas as pd

from dateutil import parser, rrule

from datetime import datetime, time, date

import time

def getRainfallData(station, day, month, year):
    """
    Function to return a dataframe of minute-level weather data for a single Wunderground PWS station.
    Args:
    station(string): Station code from the Wunderground website
    day(int): Day of month for which data is requested
    month(int): Month for which data is requested
    year(int): Year for which data is requested
    Returns:
    Pandas Dataframe with weather data for specified station and date.
    """
    url = "http://www.wunderground.com/weatherstation/WXDailyHistory.asp?ID={station}&day={day}&month={month}&year={year}&graphspan=day&format=1"
    full_url = url.format(station=station, day=day, month=month, year=year)

    # Request data from wunderground
    response = requests.get(full_url, headers={'User-agent': 'Mozilla/5.0 (Windows NT 6.1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/41.0.2228.0 Safari/537.36'})
```

```

print("Workingon{}".format(station))d
ata[station] = []

for date in dates:

# Print period status update messages

if date.day % 10 == 0:

print("Workingondate:{}forstation{}".format(date,station))don

e = False

while done == False:

try:

weather_data=getRainfallData(station,date.day,date.month,date.year)do

ne= True

except ConnectionError as e:

#MaygetratelimitedbyWunderground.com,backoffifso.print

("Got connection error on {}".format(date))

print("Willretryin{}seconds".format(backoff_time))ti

me.sleep(10)

# Add each processed date to the overall

datadata[station].append(weather_data)

# Finally combine all of the individual days and output to CSV for

analysis.pd.concat(data[station]).to_csv("data/{ }_weather.csv".form

```

```

monthly.replace({"Rainy": {True: "Wet", False: "Dry"}},
inplace=True)
monthly['month_name'] = monthly['month'].apply(lambda x: calendar.mo
nth_abbr[x]) # Get aggregate stats for each day in the dataset on rain in general - for
heatmaps.rainy_days = data.groupby(['day']).agg({
"rain": {"rain": lambda x: (x >
0.0).any(), "rain_amount": "sum"},
"total_rain": {"total_rain":
"max"}, "get_wet_cycling": {"get_wet_cycling"
: "any"}
})
# clean up the aggregated data to a more easily analysed
set: rainy_days.reset_index(drop=False, inplace=True) # remove the 'day' as the index
rainy_days.rename(columns={"": "date"}, inplace=True) # The old index column didn't have a name
- add "date" as name
rainy_days.columns = rainy_days.columns.droplevel(level=0) # The aggregation left us with
a multi-index
# Remove the top level of this index.
rainy_days['rain'] = rainy_days['rain'].astype(bool) # Change the "rain" column to True/False values
# Add the number of rainy hours per day to the rainy_days dataset. temp =
data.groupby(["day", "hour_of_day"])['raining'].any()
temp =
temp.groupby(level=[0]).sum().reset_index()
temp.rename(columns={'raining': 'hours_raining'},
inplace=True)
temp['day'] = temp['day'].apply(lambda x: x.to_datetime().date())
(rainy_days, station)

```

At this point, we have two basic data frames which we can use to visualise patterns for the citybeing analysed.

BarchartofMonthlyRainyCycles

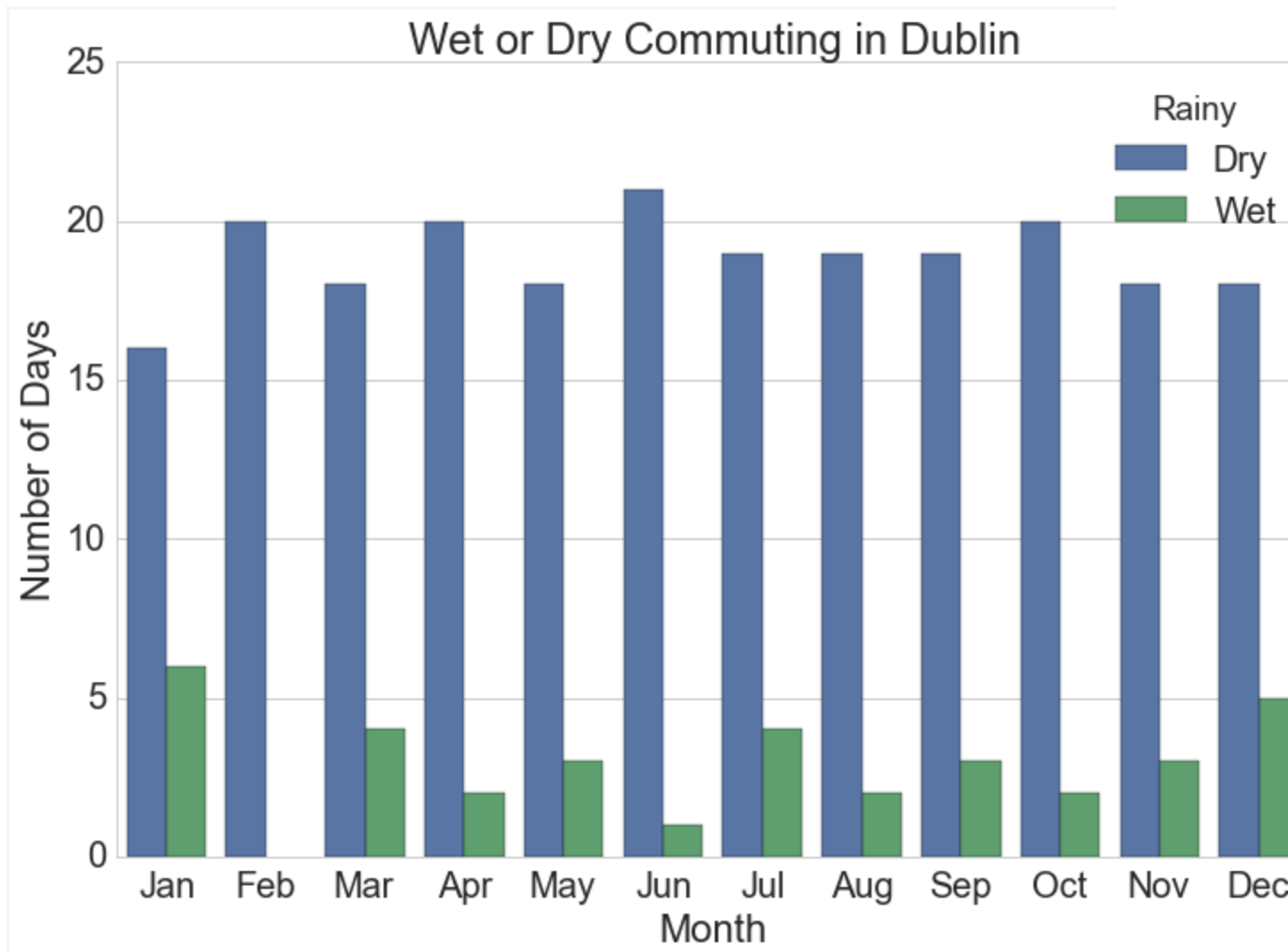
```
chart.# Monthly plot of rainy days
```

```
sns.barplot(x="month_name",y="Days",hue="Rainy",data=monthly.sort_values(['month','Rainy']))

plt.xlabel("Month")plt.ylabel

("NumberofDays")

plt.title("WetorDryCommutingin{ }".format(station))
```



Number of days monthly when cyclists get wet commuting at typical work times in Dublin, Ireland.

Heatmaps of Rainfall and Rainy Hours per day

The heatmaps shown on the blog post are generated using the “calmap” python library, installable using pip. Simply import the library, and form a Pandas series with a `DatetimeIndex` and the library takes care of the rest. I had some difficulty here with font sizes, so had to increase the size of the plot overall to counter.

```
import calmap

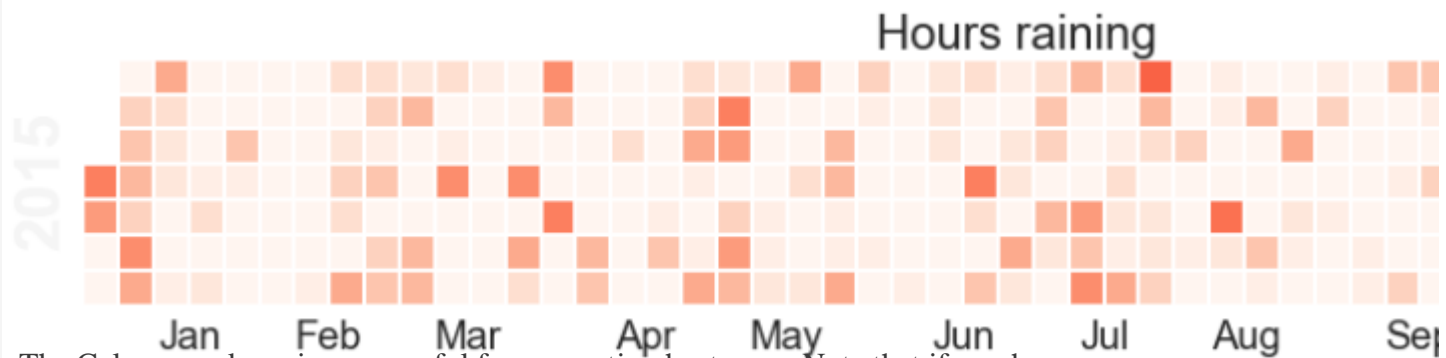
temp =
rainy_days.copy().set_index(pd.DatetimeIndex(analysis['rainy_days']['date']))#temp.s
et_index('date',inplace=True)

fig, ax = calmap.calendarplot(temp['hours_raining'], fig_kws={"figsize":(15,4)})
```

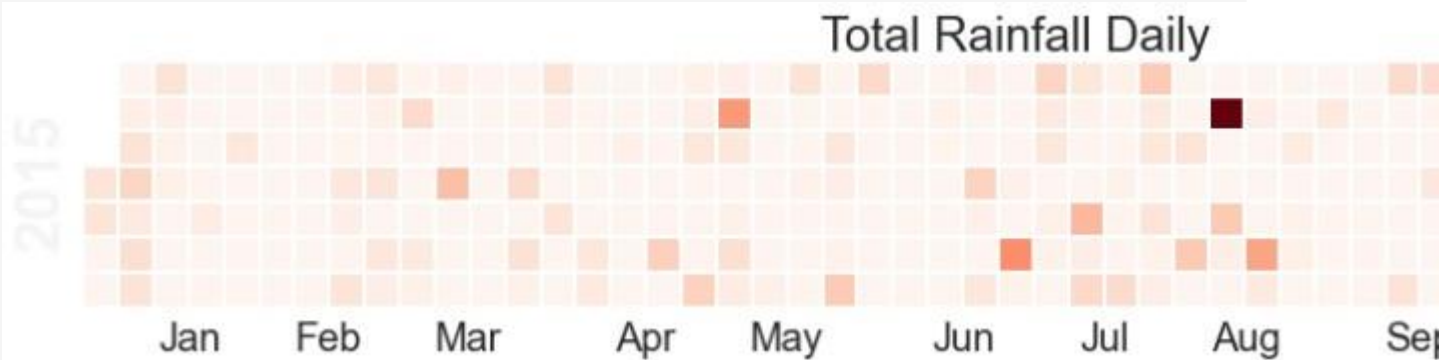
```
plt.title("Hours raining")
```

```
fig, ax = calmap.calendarplot(temp['total_rain'],
```

```
fig_kws={"figsize":(15,4)})plt.title("TotalRainfall Daily")
```



The Calmap package is very useful for generating heatmaps. Note that if you have highly outlying points of data, these will skew your color mapping considerably – I'd advise removing or reducing them for visualisation purposes.



Comparison of Every City in Dataset

To compare every city in the dataset, summary stats for each city were calculated in advance and then the plot was generated using the seaborn library. To achieve this as quickly as possible, I wrapped the entire data preparation and cleansing phase described above into a single function called “analyse data”, used this function on each city’s dataset, and extracted out the pieces of information needed for the plot

Here’s the wrapped analyse_data function:

```
def analyse_station(data_raw, station):  
    """  
    Function to analyse weather data for a period from one weather station. Args:  
    data_raw (pd.DataFrame): Pandas Dataframe made from CSV downloaded from wunderground.com  
    station (String): Name of station being analysed (for comments) Returns:  
    dict: Dictionary with analysis in keys:  
    data: Processed and cleansed data  
    monthly: Monthly aggregated statistics on rain fall etc.  
    wet_cycling: Data on working days and whether you get wet or not  
    commuting_rainy_days: Daily total rainfall for each day in dataset.  
    """  
    # Give the variables some friendlier names and convert types as necessary. data_r  
    aw['temp'] = data_raw['TemperatureC'].astype(float) data_raw['rain'] =  
    data_raw['HourlyPrecipMM'].astype(float) data_raw['total_rain'] =  
    data_raw['dailyrainMM'].astype(float) data_raw['date'] = data_raw['DateUTC'].apply(parser.parse)
```

```

data_raw['humidity'] =
data_raw['Humidity'].astype(float) data_raw['wind_direction']=d
ata_raw['WindDirectionDegrees'] data_raw['wind']=
data_raw['WindSpeedKMH']

# Extract out only the data we need.

data=data_raw.loc[:,['date','station','temp','rain','total_rain','humidity','wind']]

data = data[(data['date'] >= datetime(2015,1,1)) & (data['date']
<=datetime(2015,12,31))]

# There's an issue with some stations that record rainfall ~-2500 where data is missing.

if (data['rain'] < -500).sum() > 10:
print("There's more than 10 messed up days for
{ }".format(station))# remove the bad samples

data = data[data['rain'] > -500]

# Assign the "day" to every date entry

data['day'] = data['date'].apply(lambda x: x.date())

# Get the time, day, and hour of each timestamp in the
datasetdata['time_of_day'] = data['date'].apply(lambda x:
x.time())data['day_of_week'] = data['date'].apply(lambda x:
x.weekday())data['hour_of_day'] =
data['time_of_day'].apply(lambda x: x.hour)# Mark the month for
each entry so we can look at monthly patternsdata['month'] =
data['date'].apply(lambda x: x.month)

# Is each time stamp on a working day (Mon-Fri)

data['working_day']=(data['day_of_week']>=0)&(data['day_of_week']<=4)#Class
ifyintomorningoreveningtimes(assumingtravelbetween8.15-9amand5.15-6pm)

data['morning']=(data['time_of_day']>=time(8,15))&(data['time_of_day']<=tim
e(9,0))

```

```

data['evening'] = (data['time_of_day'] >= time(17,15)) & (data['time_of_day']
<=time(18,0))

# If there's any rain at all, mark
data['raining'] = data['rain']
> 0.0

# You get wet cycling if its a working day, and its raining at the travel
times!data['get_wet_cycling'] = (data['working_day']) & ((data['morning'] &
data['rain']) |(data['evening'] & data['rain']))

# Looking at the working days only:
wet_cycling=data[data['working_day']==True].groupby('day')['get_wet_cycling'].any()
wet_cycling =
pd.DataFrame(wet_cycling).reset_index()# Group by
month for display
wet_cycling['month'] = wet_cycling['day'].apply(lambda x: x.month)
monthly =
wet_cycling.groupby('month')['get_wet_cycling'].value_counts().reset_index()monthly.re
name(columns={"get_wet_cycling": "Rainy", 0: "Days"},
inplace=True)monthly.replace({"Rainy": {True: "Wet", False: "Dry"}},
inplace=True)monthly['month_name']=monthly['month'].apply(lambda x:calendar.mont
h_abbr[x])

# Get aggregate stats for each day in the
dataset.rainy_days = data.groupby(['day']).agg({
"rain": {"rain": lambda x: (x >
0.0).any(), "rain_amount": "sum"}

```

```
plt.suptitle("What percentage of your cycles to work do you need a raincoat?",
y=1.05,fontsize=32)

plt.title("Based on Wunderground.com weather data for 2015",
fontsize=18)plt.xticks(rotation=60)plt.savefig("images/city_comparison_wet_c
ommutes.png",bbox_inches='tight')
```

