

Name: Chaoran Li
Student ID: 2021489307 (UTD ID) cxl190012 (NET ID)
Course Number: CS 6364.002

Homework 1 Writeup Part

Task 1: Test Python Environment

```
Task 1: Test Python Environment  
hello world
```

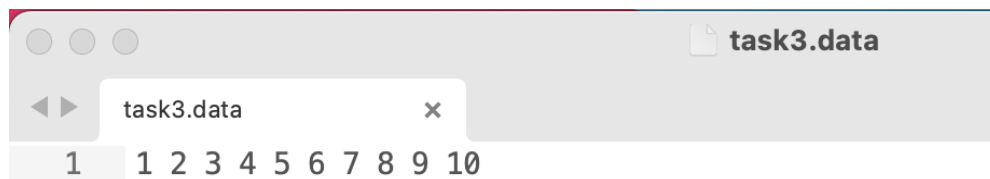
The code is very simple and the above result showed that our python environment works well.

Task 2: Define Object

```
Task 2: Define Object  
items: [1, 2, 3, 4, 5]
```

For this task, I directly assigned a list to object items and returned it to main. Hence, I can use it in below tasks.

Task 3: File Reading



```
Task 3: File Reading  
items1: [1, 2, 3, 4, 5]  
items2: [6, 7, 8, 9, 10]
```

After creating task3.data, I firstly read whole file as a string and splitted it with space. Then, I directly splitted the list into two half because I could not see other instructions for how to split the elements.

Task 4: Data Structure

First, wrote a program to use all three functions: items(), keys() and values(). These three functions are iterable and can be used to search in dictionary:

Code:

```
def task4_show_dict_functions(dic):  
    print("Show three dictionary functions:")  
    print("items(): " + str(type(dic.items())))  
    for (k, v) in dic.items():  
        print("{0} -> {1}".format(k, v))  
    print("keys(): " + str(type(dic.keys())))  
    for k in dic.keys():  
        print(k)  
    print("values(): " + str(type(dic.values())))  
    for v in dic.values():  
        print(v)
```

Then, I chose to tranverse items() to print required result.

Print the required results:

school: UAlbany

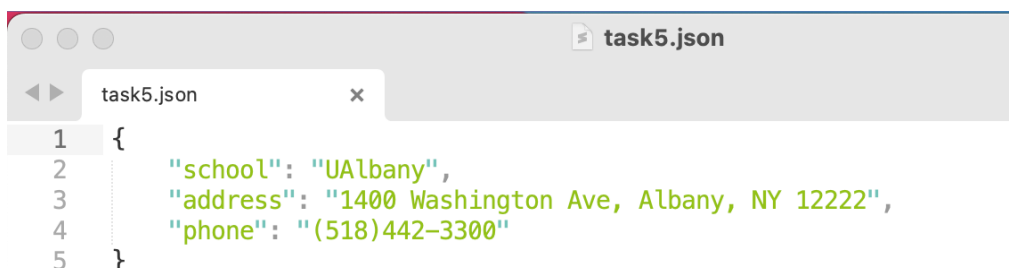
address: 1400 Washington Ave, Albany, NY 12222

phone: (518)442-3300

Task 5: Data Serialization

Task 5: Data Serialization

```
{  
    "school": "UAlbany",  
    "address": "1400 Washington Ave, Albany, NY 12222",  
    "phone": "(518)442-3300"  
}  
Save dict to task5.json.  
<class 'dict'>, {'school': 'UAlbany', 'address': '1400 Washington Ave, Albany, NY 12222', 'phone': '(518)442-3300'}  
Read dict from task5.json.
```



```
task5.json  
1 {  
2     "school": "UAlbany",  
3     "address": "1400 Washington Ave, Albany, NY 12222",  
4     "phone": "(518)442-3300"  
5 }
```

I first saved the dictionary to task5.json. Then, I read it back. As the above result shown, the return result is a dictionary and is just the same as we assigned before.

Task 6: Data Serialization

For this, I still use json module to solve this question. I also learned from my friend that he used pickle module to solve this question. I can give another solution by pickle if necessary. I firstly use json.dumps() to save objects to file. To meet the requirement that I should save objects to file one by one, I firstly create an empty file. Then, I append the string to file each time. I use “\n” to split objects from each other. Then, when I need to load them back, I read the file by line and find target object from it. Below is the output of my code.

```
Task 6: Data Serialization
[1, 2, 3, 4, 5]
{'school': 'UAlbany', 'address': '1400 Washington Ave, Albany, NY 12222', 'phone': '(518)442-3300'}
Save data structures to task6.data.
items: <class 'list'>, [1, 2, 3, 4, 5]
data: <class 'dict'>, {'school': 'UAlbany', 'address': '1400 Washington Ave, Albany, NY 12222', 'phone': '(518)442-3300'}
Read data structures from task6.data.
```



Task 7: Data Preprocessing

I firstly read some tweets and use tweet.keys() to analyze the keywords of all tweets.

```
dict_keys(['lang', 'favorited', 'truncated', 'text', 'created_at', 'retweeted', 'source', 'place', 'user', 'retweet_count', 'id', 'favorite_count'])
dict_keys(['lang', 'favorited', 'truncated', 'text', 'created_at', 'retweeted', 'source', 'user', 'urls', 'retweet_count', 'id', 'favorite_count'])
dict_keys(['lang', 'favorited', 'truncated', 'text', 'created_at', 'retweeted', 'user_mentions', 'source', 'user', 'urls', 'retweet_count', 'id', 'favorite_count'])
```

Since I found ‘id’ is the keyword, below work is just reading all tweets and collecting all ids.

```
[429129916446031872, 429117247307923456, 429315798893461505, 429079910091476992, 429030956888899584, 429152645702352896, 429016075921940481, 428898232690044928, 4289566
```

```
Present with first 20 tweet ids:
429129916446031872, 429117247307923456, 429315798893461505, 429079910091476992, 429030956888899584,
429152645702352896, 429016075921940481, 428898232690044928, 428956892602191873, 429111068368326656,
429167217976958976, 429379822314614785, 429235096919359488, 428986795586383872, 429238790310203392,
429509723516194817, 429509617756811264, 429509135437017089, 429509033045659648, 429508763708817408,
...
```

I chose first 20 tweet ids for representation since the result is quite long.

Task 8: Data Preprocessing: tweets filtering



In Task 7, we could find the keyword for this task is 'created_at'. I chose first 10 recently tweets and found them all created at Feb 01, 2014.

Task 9: File operation

For this task, assume we do not need to sort tweets in each group. We found 18 different labels which also means 18 files in task9-output folder. Some of them only contains one tweet.

Task 9: File operations

```
01-31-2014-05, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-31-2014-18, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-31-2014-02, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-30-2014-23, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-31-2014-07, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-30-2014-22, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-30-2014-14, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-30-2014-18, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-31-2014-04, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-31-2014-08, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-31-2014-22, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-31-2014-12, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-30-2014-20, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
01-31-2014-13, [{'lang': 'en', 'favorited': False, 'truncated': False, 't
02-01-2014-07, [{'lang': 'en', 'favorited': False, 'retweeted_status': {'
02-01-2014-06, [{'lang': 'en', 'favorited': False, 'retweeted_status': {'
02-01-2014-05, [{'lang': 'en', 'favorited': False, 'retweeted_status': {'
```

Task 10: NLP and Sentiment Analysis

Actually, I spent a lot of time on this. Because I use python 3.8 in my environment, but pattern only support python 3.6 for pip install.

You can use this link to help you install pattern for your environment:

<https://github.com/clips/pattern>

Besides, my IDE required me to install nltk even I never use it in my code. My guess is that module pattern relies on nltk module. And in python 3, you should use:

from pattern.text.en import positive, sentiment

instead of

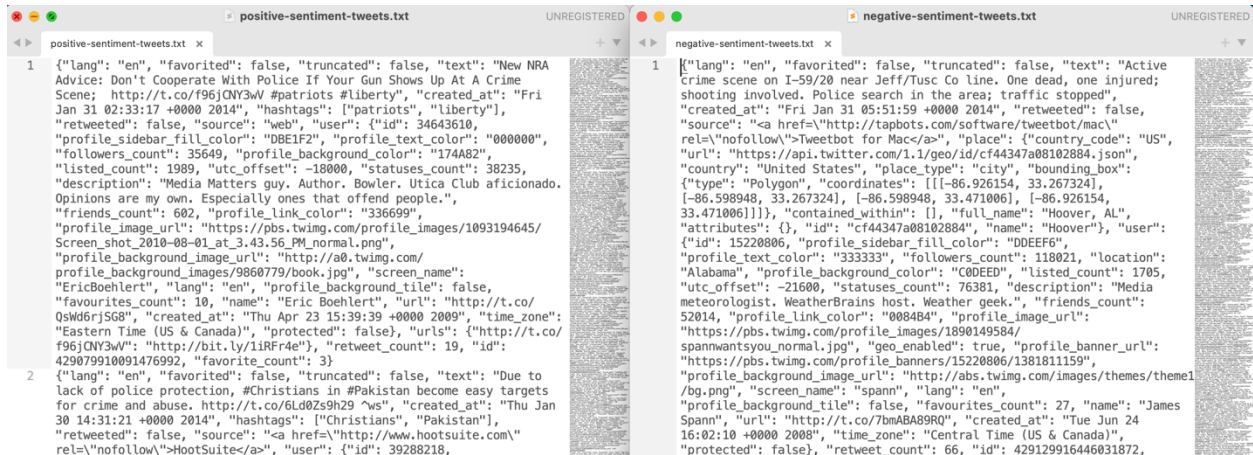
from pattern.en import positive, sentiment

Show the sentiment score of each tweet:

Task 10: NLP and Sentiment Analysis

```
(-0.07777777777777778, 0.46666666666666666)
(-0.4, 0.7)
(0.0, 0.06666666666666667)
(0.13636363636363635, 0.45454545454545453)
(-0.15000000000000002, 0.15000000000000002)
(0.0, 0.0)
(0.0, 0.0)
(0.15416666666666667, 0.6041666666666667)
(0.13636363636363635, 0.45454545454545453)
(-0.18333333333333338, 0.8666666666666667)
(0.0, 0.0)
(0.0, 0.0)
(0.35714285714285715, 0.5714285714285714)
(0.0, 0.0)
(-0.1953125, 0.07083333333333333)
```

Then, I use positive function to separate tweets into positive-sentiment-tweets.txt and negative-sentiment-tweets.txt



The image shows two side-by-side text editors. The left editor, titled 'positive-sentiment-tweets.txt', contains two JSON objects. The first object (line 1) is for a tweet about the NRA, and the second (line 2) is about Christians in Pakistan. The right editor, titled 'negative-sentiment-tweets.txt', contains one JSON object (line 1) for a tweet about a crime scene. Both JSON objects include fields for language, sentiment, text, creation time, user information, and various counts.

```
1 {"lang": "en", "favorited": false, "truncated": false, "text": "New NRA Advice: Don't Cooperate With Police If Your Gun Shows Up At A Crime Scene; http://t.co/196jQNY3w #patriots #Liberty", "created_at": "Fri Jan 31 02:33:17 +0000 2014", "hashtags": ["patriots", "liberty"], "retweeted": false, "source": "web", "user": {"id": 34643610, "profile_sidebar_fill_color": "08E1F2", "profile_text_color": "000000", "followers_count": 35649, "profile_background_color": "174A82", "listed_count": 1989, "utc_offset": -18000, "statuses_count": 38235, "description": "Media Matters guy. Author. Bowler. Utica Club aficionado. Opinions are my own. Especially ones that offend people.", "friends_count": 602, "profile_link_color": "336699", "profile_image_url": "https://pbs.twimg.com/profile_images/1093194645/Screen_shot_2010-08-01_at_3.43.56_PM_normal.png", "profile_background_image_url": "http://a0.twimg.com/profile_background_images/9868779/book.jpg", "screen_name": "EricBoehlert", "lang": "en", "profile_background_tile": false, "favourites_count": 10, "name": "Eric Boehlert", "url": "http://t.co/QsWd6rjSG8", "created_at": "Thu Apr 23 15:39:39 +0000 2009", "time_zone": "Eastern Time (US & Canada)", "protected": false, "urls": {"http://t.co/196jQNY3w": "http://bit.ly/1iRFr4e"}, "retweet_count": 19, "id": 429079910091476992, "favorite_count": 3}

2 {"lang": "en", "favorited": false, "truncated": false, "text": "Due to lack of police protection, #Christians in #Pakistan become easy targets for crime and abuse. http://t.co/6Ld0Zs9h29 ^ws", "created_at": "Thu Jan 30 14:31:21 +0000 2014", "hashtags": ["Christians", "Pakistan"], "retweeted": false, "source": "<a href='\"http://www.hootsuite.com\"' rel='\"nofollow\"'>HootSuite</a>", "user": {"id": 39288218,
```

```
1 [{"lang": "en", "favorited": false, "truncated": false, "text": "Active Crime scene on I-59/20 near Jeff/Tusc Co line. One dead, one injured; shooting involved. Police search in the area; traffic stopped", "created_at": "Fri Jan 31 05:51:59 +0000 2014", "retweeted": false, "source": "<a href='\"http://tapbots.com/software/tweetbot/mac\"' rel='\"nofollow\"'>Tweetbot for Mac</a>", "place": {"country_code": "US", "url": "https://api.twitter.com/1.1/geo/id/cf44347a08102884.json", "country": "United States", "place_type": "city", "bounding_box": {"type": "Polygon", "coordinates": [[[[-86.926154, 33.267324], [-86.598948, 33.471006], [-86.598948, 33.267324], [33.471006]]]}], "contained_within": [], "full_name": "Hoover, AL", "attributes": {}, "id": "cf44347a08102884", "name": "Hoover", "user": {"id": 15220806, "profile_sidebar_fill_color": "DDEEF6", "profile_text_color": "333333", "followers_count": 118021, "location": "Alabama", "profile_background_color": "C0DEED", "listed_count": 1705, "utc_offset": -21600, "statuses_count": 76381, "description": "Media meteorologist. WeatherBrains host. Weather geek.", "friends_count": 52014, "profile_link_color": "0084B4", "profile_image_url": "https://pbs.twimg.com/profile_images/1890149584/spannwantsyou_normal.jpg", "geo_enabled": true, "profile_banner_url": "https://pbs.twimg.com/profile_banners/15220806/1381811159", "profile_background_image_url": "http://abs.twimg.com/images/themes/theme1/bg.png", "screen_name": "spann", "lang": "en", "profile_background_tile": false, "favourites_count": 27, "name": "James Spann", "url": "http://t.co/7bmABAB9RQ", "created_at": "Tue Jun 24 16:02:10 +0000 2008", "time_zone": "Central Time (US & Canada)", "protected": false, "retweet_count": 66, "id": 429129916446031872,
```