Big Data and Automated Content Analysis

Week 2 – Thursday Data harvesting and storage

Damian Trilling

d.c.trilling@uva.nl @damian0604 www.damiantrilling.net

Afdeling Communicatiewetenschap Universiteit van Amsterdam

9 April 2020



Today

This week's exercise

Discussing the code

Reading a JSON file into a dict, looping over the dict

Task 1: Print all titles of all videos

```
import json

with open("/home/damian/pornexercise/xhamster.json") as fi:

data=json.load(fi)

for k,v in data.items()):
 print (v["title"])
```

Reading a JSON file into a dict, looping over the dict

Task 1: Print all titles of all videos

```
import json

with open("/home/damian/pornexercise/xhamster.json") as fi:
   data=json.load(fi)

for k,v in data.items()):
   print (v["title"])
```

NB: You have to know (e.g., by reading the documentation of the dataset) that the key is called title

Reading a JSON file into a dict, looping over the dict

Task 1: Print all titles of all videos

```
import json

with open("/home/damian/pornexercise/xhamster.json") as fi:
   data=json.load(fi)

for k,v in data.items()):
   print (v["title"])
```

NB: You have to know (e.g., by reading the documentation of the dataset) that the key is called title

NB: data is in fact a dict of dicts, such that each value v is another dict.

For each of these dicts, we retrieve the value that corresponds to the key title

What to do if you do not know the structure of the dataset?

Inspecting your data: use the functions type() and len() and/or the dictionary method .keys()

```
1 len(data)
```

- 2 type(data)
- 3 data.keys()

What to do if you do not know the structure of the dataset?

Inspecting your data: use the functions type() and len() and/or
the dictionary method .keys()

```
l len(data)
type(data)
data.keys()
```

len() returns the number of items of an object; type() returns the type
object; .keys() returns a list of all available keys in the dictionary

What to do if you do not know the structure of the dataset?

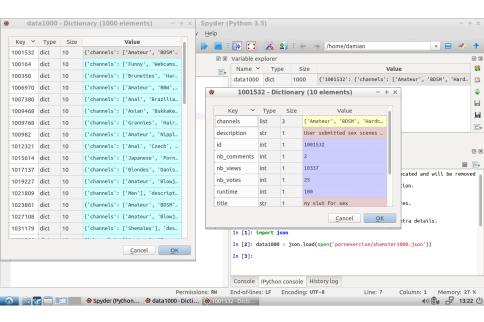
Inspecting your data: use the module pprint

```
1 from pprint import pprint
```

pprint(data)

(but better do this on a smaller subset of the data!)

```
{'1002968': {'channels': ['Tits'],
             'description': 'NA',
             'id': 1002968.
             'nb comments': 6,
             'nb views': 26683,
             'nb votes': 68.
             'runtime': 1083,
             'title': 'big tit chick gets fucked',
             'upload date': '2012-01-22',
             'uploader': '894621d2f89ab1b859c17b9b81f182e82a375c55'},
 '1003007': {'channels': ['Hardcore', 'Spanking', 'Squirting'],
             'description': 'Antonia plays with her 16 inch toy and her '
                            'crop, till she sprays cum into her own mouth!',
             'id': 1003007.
             'nb comments': 5.
             'nb views': 5008,
             'nb votes': 35,
             'runtime': 213.
             'title': 'antonia just playin around',
             'upload date': '2012-01-22',
```



For the sake of completeness...

.items() returns a key-value *pair*, that's why we need to assign *two* variables in the for statement.

These alternatives would also work:

```
for v in data.values()):
    print(v["title"])

for k in data: #or: for k in data.keys():
    print(data[k]["title"])
```

Do you see (dis-)advantages?

Working with a subset of the data

What to do if you want to work with a smaller subset of the data? Taking a random sample of 10 items in a dict:

```
import random
mydict_short = dict(random.sample(mydict.items(),10))
```

Taking the first 10 elements in a list:

```
mylist_short = mylist[:10]
```

Initializing variables, merging two lists, using a counter

Task 2: Average tags per video and most frequently used tags

```
from collections import Counter
    alltags=[]
    i=0
    for k.v in data.items():
6
       i += 1
       alltags.extend(v["channels"])
7
8
    print(len(alltags), "tags are describing", i, "different videos")
    print("Thus, we have an average of",len(alltags)/i,"tags per video")
10
11
    c=Counter(alltags)
12
    print (c.most_common(100))
13
```

(there are other, more efficient ways of doing this)

Nesting blocks, using a defaultdict to count, error handling

Task 3: What porn category is most frequently commented on?

```
from collections import defaultdict
2
    commentspercat=defaultdict(int)
3
    for k.v in data.items():
           for tag in v["channels"]:
5
               try:
6
                   commentspercat[tag]+=int(v["nb_comments"])
7
8
               except:
9
                  pass
10
    print(commentspercat)
    # if you want to print in a fancy way, you can do it like this:
11
    for tag in sorted(commentspercat, key=commentspercat.get, reverse=True):
12
       print( tag,"\t", commentspercat[tag])
13
```

A defaultdict is a normal dict, with the difference that the type of each value is pre-defined and it doesn't give an error if you look up a non-existing key

Nesting blocks, using a defaultdict to count, error handling

Task 3: What porn category is most frequently commented on?

```
from collections import defaultdict
2
    commentspercat=defaultdict(int)
3
    for k,v in data.items():
           for tag in v["channels"]:
5
6
               try:
                   commentspercat[tag]+=int(v["nb_comments"])
7
8
               except:
9
                  pass
10
    print(commentspercat)
    # if you want to print in a fancy way, you can do it like this:
11
    for tag in sorted(commentspercat, key=commentspercat.get, reverse=True):
12
       print( tag,"\t", commentspercat[tag])
13
```

A defaultdict is a normal dict, with the difference that the type of each value is pre-defined and it doesn't give an error if you look up a non-existing key

NB: In line 7, we assume the value to be an int, but the datasets sometimes contains the string "NA" instead of a string representing an int. That's why we need the try/except construction

Adding elements to a list, sum() and len()

Task 4: Average length of descriptions

```
length=[]
for k,v in data.items():
    length.append(len(v["description"]))

print ("Average length",sum(length)/len(length))
```

Extending (merging) vs appending

```
Merging:
```

```
11 = [1,2,3]
2 	 12 = [4,5,6]
3 # either:
   11 = 11 + 12
5 # or:
   11.extend(12)
   print(11)
   gives [1,2,3,4,5,6]
   Appending:
1 \quad 11 = [1,2,3]
2 	 12 = [4,5,6]
   11.append(12)
   print(11)
```

gives [1,2,3,[4,5,6]] 12 is seen as *one* element to append to 11

Tokenizing with .split()

Task 5: Most frequently used words

```
allwords=[]
for k,v in data.items():
    allwords+=v["description"].split()

c2=Counter(allwords)
print(c2.most_common(100))
    .split() changes a string to a list of words.
"This is cool''.split()
results in
["This", "is", "cool"]
```

Concluding remarks

Make sure you fully understand the code!

Re-read the corresponding chapters

PLAY AROUND!!!