



17.7. Extracting from Nested Data

A common problem, especially when dealing with data returned from a web site, is to extract certain elements from deep inside a nested data structure. In principle, there's nothing more difficult about pulling something out from deep inside a nested data structure: with lists, you use [] to index or a for loop to get them them all; with dictionaries, you get the value associated with a particular key using [] or iterate through all the keys, accessing the value for each. But it's easy to get lost in the process and think you've extracted something different than you really have. Because of this, we have created a usable technique to help you during the debugging process.

Follow the system described below and you will have success with extracting nested data. The process involves the following steps:

1. Understand the nested data object.
2. Extract one object at the next level down.
3. Repeat the process with the extracted object

Understand. Extract. Repeat.

To illustrate this, we will walk through extracting information from data formatted in a way that it's return by the Twitter API. This nested dictionary results from querying Twitter, asking for three tweets matching "University of Michigan". As you'll see, it's quite a daunting data structure, even when printed with nice indentation as it's shown below.

```
Save & Run Original - 1 of 1 Show in CodeLens
```

```
1 res = {  
2     "search_metadata": {  
3         "count": 3,  
4         "completed_in": 0.015,  
5         "max_id": "536624519285583872",  
6         "since_id": "0",  
7         "next_results": "?max_id=536623674942439424&q=University%20of%20Michigan&count=3",  
8         "refresh_url": "?since_id=536624519285583872&q=University%20of%20Michigan&inclu  
9         "since_id": 0,  
10        "query": "University+of+Michigan",  
11        "max_id": 536624519285583872  
12    },  
13    "statuses": [  
14        {  
15            "contributors": None,
```

Activity: 1 -- ActiveCode (ac300_1_1)

17.7.1. Understand

At any level of the extraction process, the first task is to make sure you understand the current object you have extracted. There are few options here.

1. Print the entire object. If it's small enough, you may be able to make sense of the printout directly. If it's a little bit larger, you may find it helpful to "pretty-print" it, with indentation showing the level of nesting of the data. We don't have a way to pretty-print in our online browser-based environment, but if you're running code with a full Python interpreter, you can use the dumps function in the json module. For example:

```
import json  
print(json.dumps(res, indent=2))
```

2. If printing the entire object gives you something that's too unwieldy, you have other options for making sense of it.
 - Copy and paste it to a site like <https://jsoneditoronline.org/> which will let you explore and collapse levels
 - Print the type of the object.
 - If it's a dictionary:
 - print the keys
 - If it's a list:
 - print its length
 - print the type of the first item
 - print the first item if it's of manageable size

```
Save & Run Original - 1 of 1 Show in CodeLens
```

```
1 import json  
2 print(json.dumps(res, indent=2)[:100])  
3 print("-----")  
4 print(type(res))  
5 print(res.keys())  
6
```

```
{"search_metadata":{"count":3,"completed_in":0.015,"max_id_str":"536624519285583872","since_id_st  
r":  
-----  
<class 'dict'>  
['search_metadata', 'statuses']
```

Activity: 2 -- ActiveCode (ac300_1_2)

17.7.2. Extract

In the extraction phase, you will be diving one level deeper into the nested data.

1. If it's a dictionary, figure out which key has the value you're looking for, and get its value. For example: `res2 = res['statuses']`
2. If it's a list, you will typically be wanting to do something with each of the items (e.g., extracting something from each, and accumulating them in a list). For that you'll want a for loop, such as `for res2 in res`. During your exploration phase, however, it will be easier to debug things if you work with just one item. One trick for doing that is to iterate over a slice of the list containing just one item. For example, `for res2 in res[:1]`.

```
Save & Run Original - 1 of 1 Show in CodeLens  
1 print(type(res))  
2 print(res.keys())  
3 res2 = res['statuses']  
4
```

```
<class 'dict'>  
['search_metadata', 'statuses']
```

Activity: 3 -- ActiveCode (ac300_1_3)

17.7.3. Repeat

Now you'll repeat the Understand and Extract processes at the next level.

17.7.3.1. Level 2

First understand.

```
Save & Run Original - 1 of 1 Show in CodeLens  
1 print(type(res))  
2 print(res.keys())  
3 res2 = res['statuses']  
4 print("----Level 2----")  
5 print(type(res2)) # it's a list!  
6 print(len(res2))  
7
```

```
<class 'dict'>  
['search_metadata', 'statuses']  
---Level 2---  
<class 'list'>  
3
```

Activity: 4 -- ActiveCode (ac300_1_4)

It's a list, with three items, so it's a good guess that each item represents one tweet.

Now extract. Since it's a list, we'll want to work with each item, but to keep things manageable for now, let's use the trick for just looking at the first item. Later we'll switch to processing all the items.

```
Save & Run Original - 1 of 1 Show in CodeLens  
1 import json  
2 print(type(res))  
3 print(res.keys())
```

```

4 res2 = res['statuses']
5 print("----Level 2: a list of tweets----")
6 print(type(res2)) # it's a list!
7 print(len(res2)) # looks like one item representing each of the three tweets
8 for res3 in res2[:1]:
9     print("----Level 3: a tweet----")
10    print(json.dumps(res3, indent=2)[:30])
11

```

```

<class 'dict'>
['search_metadata', 'statuses']
---Level 2: a list of tweets---
<class 'list'>
3
---Level 3: a tweet---
{"id": "536624519285583872", "id"

```

Activity: 5 -- ActiveCode (ac300_1_5)

17.7.3.2. Level 3

First understand.

```

Save & Run Original - 1 of 1 Show in CodeLens
1 import json
2 print(type(res))
3 print(res.keys())
4 res2 = res['statuses']
5 print("----Level 2: a list of tweets----")
6 print(type(res2)) # it's a list!
7 print(len(res2)) # looks like one item representing each of the three tweets
8 for res3 in res2[:1]:
9     print("----Level 3: a tweet----")
10    print(json.dumps(res3, indent=2)[:30])
11    print(type(res3)) # it's a dictionary
12    print(res3.keys())
13

```

```

<class 'dict'>
['search_metadata', 'statuses']
---Level 2: a list of tweets---
<class 'list'>
3
---Level 3: a tweet---
{"id": "536624519285583872", "id"
<class 'dict'>
['id', 'id_str', 'entities', 'lang', 'created_at', 'contributors', 'truncated', 'text', 'in_reply_to_status_id', 'favorite_count', 'retweeted', 'source', 'coordinates', 'in_reply_to_screen_name', 'in_reply_to_user_id', 'retweet_count', 'favorited', 'user', 'geo', 'in_reply_to_user_id_str', 'in_reply_to_status_id_str', 'place', 'metadata', 'retweeted_status']

```

Activity: 6 -- ActiveCode (ac300_1_6)

Then extract. Let's pull out the information about who sent each of the tweets. Probably that's the value associated with the 'user' key.

```

Save & Run Original - 1 of 1 Show in CodeLens
1 import json
2 print(type(res))
3 print(res.keys())
4 res2 = res['statuses']
5 print("----Level 2: a list of tweets----")
6 print(type(res2)) # it's a list!
7 print(len(res2)) # looks like one item representing each of the three tweets
8 for res3 in res2[:1]:
9     print("----Level 3: a tweet----")
10    print(json.dumps(res3, indent=2)[:30])
11    res4 = res3['user']
12

```

```

<class 'dict'>
['search_metadata', 'statuses']
---Level 2: a list of tweets---
<class 'list'>
3
---Level 3: a tweet---
{"id": "536624519285583872", "id"

```

Activity: 7 -- ActiveCode (ac300_1_7)

Now repeat.

17.7.3.3. Level 4

Understand.

The screenshot shows an ActiveCode editor window. At the top, there are three buttons: 'Save & Run', 'Original - 1 of 1', and 'Show in CodeLens'. The code in the editor is as follows:

```
1 import json
2 print(type(res))
3 print(res.keys())
4 res2 = res['statuses']
5 print("----Level 2: a list of tweets----")
6 print(type(res2)) # it's a list!
7 print(len(res2)) # looks like one item representing each of the three tweets
8 for res3 in res2[1:]:
9     print("----Level 3: a tweet----")
10    print(json.dumps(res3, indent=2)[:30])
11    res4 = res3['user']
12    print("----Level 4: the user who wrote the tweet----")
13    print(type(res4)) # it's a dictionary
14    print(res4.keys())
15
```

A callout box highlights the output of the code, which is a JSON dump of a single tweet object. The output is:

```
<class 'dict'>
['search_metadata', 'statuses']
----Level 2: a list of tweets----
<class 'list'>
3
----Level 3: a tweet----
{"id": "536624519285583872", "id_str": "536624519285583872", "created_at": "Wed Apr 09 14:34:41 +0000 2014", "source": "froyoh0", "truncated": false, "text": "Extract. Let's print out the user's screen name and when their account was created.", "in_reply_to_status_id": null, "in_reply_to_status_id_str": null, "in_reply_to_user_id": null, "in_reply_to_user_id_str": null, "in_reply_to_screen_name": null, "geo": null, "coordinates": null, "place": null, "contributors": null, "is_quote_status": false, "retweet_count": 0, "favorite_count": 0, "favorited": false, "retweeted": false, "possibly_sensitive": false, "lang": "en", "display_text_range": [0, 280], "entities": {"hashtags": [], "symbols": [], "urls": [{"url": "http://t.co/31brooks_"}], "user_mentions": [{"name": "froyoh0", "screen_name": "froyoh0", "id": 213111111111111111, "id_str": "213111111111111111", "indices": [1, 11]}], "urls": [{"url": "http://t.co/31brooks_"}]}, "profile_use_background_image": true, "profile_text_color": "white", "default_profile_image": false, "profile_background_image_url_https": "https://abs.twimg.com/images/themes/theme1/bg.png", "verified": false, "profile_location": null, "profile_image_url_https": "http://pbs.twimg.com/profile_images/536624519285583872/213111111111111111_normal.jpg", "profile_sidebar_fill_color": "000000", "entities": {"user_mentions": [{"name": "froyoh0", "screen_name": "froyoh0", "id": 213111111111111111, "id_str": "213111111111111111", "indices": [1, 11]}]}, "profile_sidebar_border_color": "000000", "profile_background_color": "000000", "listed_count": 0, "is_translator": false, "profile_link_color": "000000", "friends_count": 0, "location": null, "profile_banner_url": "http://pbs.twimg.com/profile_banners/536624519285583872/1397145211", "profile_background_tile": true, "favourites_count": 0, "notifications": null, "url": null, "contributors_enabled": false, "time_zone": null, "protected": false, "default_profile": true, "is_translator": false}
```

Below the code, the status bar says 'Activity: 8 -- ActiveCode (ac300_1_8)'.

Extract. Let's print out the user's screen name and when their account was created.

The screenshot shows an ActiveCode editor window. At the top, there are three buttons: 'Save & Run', 'Original - 1 of 1', and 'Show in CodeLens'. The code in the editor is as follows:

```
1 import json
2 # print(type(res))
3 # print(res.keys())
4 res2 = res['statuses']
5 # print("----Level 2: a list of tweets----")
6 # print(type(res2)) # it's a list!
7 # print(len(res2)) # looks like one item representing each of the three tweets
8 for res3 in res2[1:]:
9     print("----Level 3: a tweet----")
10    # print(json.dumps(res3, indent=2)[:30])
11    res4 = res3['user']
12    print("----Level 4: the user who wrote the tweet----")
13    # print(type(res4)) # it's a dictionary
14    # print(res4.keys())
15    print(res4['screen_name'], res4['created_at'])
```

A callout box highlights the output of the code, which is the screen name and creation date of the user. The output is:

```
----Level 3: a tweet----
----Level 4: the user who wrote the tweet----
31brooks_ Wed Apr 09 14:34:41 +0000 2014
```

Below the code, the status bar says 'Activity: 9 -- ActiveCode (ac300_1_9)'.

Now, we may want to go back have it extract for all the items rather than only the first item in res2.

The screenshot shows an ActiveCode editor window. At the top, there are three buttons: 'Save & Run', 'Original - 1 of 1', and 'Show in CodeLens'. The code in the editor is as follows:

```
1 import json
2 # print(type(res))
3 # print(res.keys())
4 res2 = res['statuses']
5 #print("----Level 2: a list of tweets----")
6 #print(type(res2)) # it's a list!
7 #print(len(res2)) # looks like one item representing each of the three tweets
8 for res3 in res2:
9     #print("----Level 3: a tweet----")
10    #print(json.dumps(res3, indent=2)[:30])
11    res4 = res3['user']
12    #print("----Level 4: the user who wrote the tweet----")
13    #print(type(res4)) # it's a dictionary
14    #print(res4.keys())
15    print(res4['screen_name'], res4['created_at'])
```

A callout box highlights the output of the code, which is the screen name and creation date of all three users. The output is:

```
31brooks_ Wed Apr 09 14:34:41 +0000 2014
froyoh0 Thu Jan 14 21:37:54 +0000 2010
```

Activity: 10 -- ActiveCode (ac300_1_10)

17.7.3.4. Reflections

Notice that each time we descend a level in a dictionary, we have a [] picking out a key. Each time we look inside a list, we will have a for loop. If there are lists at multiple levels, we will have nested for loops.

Once you've figured out how to extract everything you want, you *may* choose to collapse things with multiple extractions in a single expression. For example, we could have this shorter version.

The screenshot shows a Jupyter Notebook cell. At the top, there are three buttons: "Save & Run", "Original - 1 of 1", and "Show in CodeLens". The code in the cell is:

```
1 for res3 in res['statuses']:
2     print(res3['user']['screen_name'], res3['user']['created_at'])
3
```

Below the code, the output is displayed in a yellow box:

```
31brooks_ Wed Apr 09 14:34:41 +0000 2014
froyoh Thu Jan 14 21:37:54 +0000 2010
MDuncan95814 Tue Sep 11 21:02:09 +0000 2012
```

At the bottom of the cell, it says "Activity: 11 -- ActiveCode (ac300_1_11)".

Even with this compact code, we can still count off how many levels of nesting we have extracted from, in this case four. `res['statuses']` says we have descended one level (in a dictionary). `for res3 in...` says we are have descended another level (in a list). `['user']` is descending one more level, and `['screen_name']` is descending one more level.

You have attempted 12 of 11 activities on this page

✓ Completed. Well Done!

17.6. Deep and Shallow Copies">

Deep and Shallow Copies">

17.8. Exercises">

17.8. Exercises">Next Section - 17.8. Exercises