

JSON examples and exercise

- get familiar with packages for dealing with JSON
- study examples with JSON strings and files
- work on exercise to be completed and submitted

- reference: <http://oandas.pydata.org/bandas-docs/stable/10.html#jo-json-reader>
- data source: <http://gonstudio.com/resources/>

```
In [4]: import pandas as pd
```

imports for Python, Pandas

```
In [5]: import json
from pandas.io.json import json_normalize
```

JSON example, with string

- demonstrates creation of normalized dataframes (tables) from nested json string
- source: <http://oandas.pydata.org/bandas-docs/stable/10.html#normalization>

```
In [6]: # define json string
data = [{"state": "Florida",
        "shortname": "FL",
        "info": {"governor": "Rick Scott"},
        "counties": [{"name": "Dade", "population": 12345},
                     {"name": "Broward", "population": 40000},
                     {"name": "Palm Beach", "population": 60000}],
        "state": "Ohio",
        "shortname": "OH",
        "info": {"governor": "John Kasich"},
        "counties": [{"name": "Summit", "population": 1234},
                     {"name": "Cuyahoga", "population": 1337}]]

In [7]: # use normalization to create tables from nested element
json_normalize(data, 'counties')
```

```
Out[7]:
```

	name	population
0	Dade	12345
1	Broward	40000
2	Palm Beach	60000
3	Summit	1234
4	Cuyahoga	1337

```
In [8]: # further populate tables created from nested element
json_normalize(data, 'counties', ['state', 'shortname', ['info', 'governor']])
```

```
Out[8]:
```

	name	population	state	shortname	info.governor
0	Dade	12345	Florida	FL	Rick Scott
1	Broward	40000	Florida	FL	Rick Scott
2	Palm Beach	60000	Florida	FL	Rick Scott
3	Summit	1234	Ohio	OH	John Kasich
4	Cuyahoga	1337	Ohio	OH	John Kasich

JSON example, with file

- demonstrates reading in a json file as a string and as a table
- uses small sample file containing data about projects funded by the World Bank
- data source: <http://gonstudio.com/resources/>

```
In [9]: # load json as string
json.load(open('data/world_bank_projects_less.json'))

-----
FileNotFoundError: [Errno 2] No such file or directory: 'data/world_bank_projects_less.json'
Traceback (most recent call last):
  <ipython-input-9-721b6769f6f5> in <module>
    1 # load json as string
----> 2 json.load(open('data/world_bank_projects_less.json'))

FileNotFoundError: [Errno 2] No such file or directory: 'data/world_bank_projects_less.json'
```

```
In [ ]: # load as Pandas dataframe
sample_json_df = pd.read_json('data/world_bank_projects_less.json')
sample_json_df
```

JSON exercise

Using data in file 'data/world_bank_projects.json' and the techniques demonstrated above.

1. Find the 10 countries with most projects
2. Find the top 10 major project themes (using column 'mjtheme_namecode')
3. In 2. above you will notice that some entries have only the code and the name is missing. Create a dataframe with the missing names filled in.

```
In [20]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

```
In [21]: df = pd.read_json('world_bank_projects.json')
df.head()
```

```
Out[21]:
```

	_id	approval	board_approval_month	boardapprovaldate	borrower	closingdate	country
0	52b213b38594d8a2be17c7807	1999	November	2013-11-12T00:00:00Z	FEDERAL DEMOCRATIC REPUBLIC OF ETHIOPIA	2018-07-07T00:00:00Z	Feder
1	52b213b38594d8a2be17c7817	2015	November	2013-11-04T00:00:00Z	GOVERNMENT OF TUNISIA	NaN	
2	52b213b38594d8a2be17c7827	2014	November	2013-11-01T00:00:00Z	MINISTRY OF FINANCE AND ECONOMIC DEVEL	NaN	
3	52b213b38594d8a2be17c7837	2014	October	2013-10-31T00:00:00Z	MIN OF PLANNING AND INTL COOPERATION	NaN	
4	52b213b38594d8a2be17c7847	2014	October	2013-10-31T00:00:00Z	MINISTRY OF FINANCE	2019-04-30T00:00:00Z	L

5 rows × 8 columns

```
In [22]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 50 columns):
 id                500 non-null object
 approval         500 non-null int64
 board_approval_month  500 non-null object
 boardapprovaldate  500 non-null object
 borrower         485 non-null object
 closingdate      370 non-null object
 country_namecode  500 non-null object
 countrycode      500 non-null object
 countryname      500 non-null object
 countryshortname  500 non-null object
 docty           446 non-null object
 envassamentcategorycode  430 non-null object
 grantant         500 non-null int64
 ibrdcommant      500 non-null int64
 id              500 non-null object
 idacommant       500 non-null int64
 impagency        412 non-null object
 lendinginstr     495 non-null object
 lendinginstrtype  495 non-null object
 lendprojectcost  500 non-null int64
 majorsector_percent  500 non-null object
 mjsector_namecode  500 non-null object
 mjtheme          491 non-null object
 mjtheme_namecode  500 non-null object
 mjthemecode      500 non-null object
 prodlinetext     500 non-null object
 productlinetype  500 non-null object
 project_abstract  362 non-null object
 project_name     500 non-null object
 projectdocs      446 non-null object
 projectfinancialtype  500 non-null object
 projectstatusdisplay  500 non-null object
 regionname      500 non-null object
 sector          500 non-null object
 sector1         500 non-null int64
 sector2         380 non-null object
 sector3         265 non-null object
 sector4         174 non-null object
 sector_namecode  500 non-null object
 sectorcode       500 non-null object
 source          500 non-null object
 status          500 non-null object
 supplementprojectflg  498 non-null object
 them1           500 non-null object
 theme_namecode   491 non-null object
 themecode        491 non-null object
 totalamt        500 non-null int64
 totalcommant     500 non-null int64
 url             500 non-null object
 dtypes: int64(7), object(43)
memory usage: 195.4+ KB
```

```
In [23]: df.dtypes

Out[23]:
```

	id	approval	board_approval_month	boardapprovaldate	borrower	closingdate	country
id	object	int64	object	object	object	object	object
approval	int64	object	object	object	object	object	object
board_approval_month	object	object	object	object	object	object	object
boardapprovaldate	object	object	object	object	object	object	object
borrower	object	object	object	object	object	object	object
closingdate	object	object	object	object	object	object	object
country_namecode	object	object	object	object	object	object	object
countrycode	object	object	object	object	object	object	object
countryname	object	object	object	object	object	object	object
countryshortname	object	object	object	object	object	object	object
docty	object	object	object	object	object	object	object
envassamentcategorycode	object	object	object	object	object	object	object
grantant	int64	object	object	object	object	object	object
ibrdcommant	int64	object	object	object	object	object	object
id	object	object	object	object	object	object	object
idacommant	int64	object	object	object	object	object	object
impagency	object	object	object	object	object	object	object
lendinginstr	object	object	object	object	object	object	object
lendinginstrtype	object	object	object	object	object	object	object
lendprojectcost	int64	object	object	object	object	object	object
majorsector_percent	object	object	object	object	object	object	object
mjsector_namecode	object	object	object	object	object	object	object
mjtheme	object	object	object	object	object	object	object
mjtheme_namecode	object	object	object	object	object	object	object
mjthemecode	object	object	object	object	object	object	object
prodlinetext	object	object	object	object	object	object	object
productlinetype	object	object	object	object	object	object	object
project_abstract	object	object	object	object	object	object	object
project_name	object	object	object	object	object	object	object
projectdocs	object	object	object	object	object	object	object
projectfinancialtype	object	object	object	object	object	object	object
projectstatusdisplay	object	object	object	object	object	object	object
regionname	object	object	object	object	object	object	object
sector	object	object	object	object	object	object	object
sector1	object	object	object	object	object	object	object
sector2	object	object	object	object	object	object	object
sector3	object	object	object	object	object	object	object
sector4	object	object	object	object	object	object	object
sector_namecode	object	object	object	object	object	object	object
sectorcode	object	object	object	object	object	object	object
source	object	object	object	object	object	object	object
status	object	object	object	object	object	object	object
supplementprojectflg	object	object	object	object	object	object	object
them1	object	object	object	object	object	object	object
theme_namecode	object	object	object	object	object	object	object
themecode	object	object	object	object	object	object	object
totalamt	int64	object	object	object	object	object	object
totalcommant	int64	object	object	object	object	object	object
url	object	object	object	object	object	object	object
dtype: object							

```
In [24]: df.describe()

Out[24]:
```

	approval	grantamt	ibrdcommant	idacommant	lendprojectcost	totalamt	totalcommant
count	500.000000	5.000000e+02	5.000000e+02	5.000000e+02	5.000000e+02	5.000000e+02	5.000000e+02
mean	2013.108000	4.432400e+06	3.286010e+07	3.542110e+07	1.547241e+08	6.828146e+07	7.271386e+07
std	0.722066	2.023307e+07	1.089197e+08	7.681431e+07	4.784211e+08	1.242862e+08	1.234705e+08
min	1999.000000	0.000000e+00	0.000000e+00	0.000000e+00	3.000000e+04	0.000000e+00	3.000000e+04
25%	2013.000000	0.000000e+00	0.000000e+00	0.000000e+00	6.472500e+06	0.000000e+00	5.000000e+06
50%	2013.000000	0.000000e+00	0.000000e+00	0.000000e+00	3.500000e+07	2.000000e+07	2.500000e+07
75%	2013.000000	1.665000e+06	0.000000e+00	3.700000e+07	1.021250e+08	8.625000e+07	9.045000e+07
max	2015.000000	3.650000e+08	1.307800e+09	6.000000e+08	5.170000e+09	1.307800e+09	1.307800e+09

```
In [25]: df.dtypes

Out[25]:
```

	id	approval	board_approval_month	boardapprovaldate	borrower	closingdate	country
id	object	int64	object	object	object	object	object
approval	int64	object	object	object	object	object	object
board_approval_month	object	object	object	object	object	object	object
boardapprovaldate	object	object	object	object	object	object	object
borrower	object	object	object	object	object	object	object
closingdate	object	object	object	object	object	object	object
country_namecode	object	object	object	object	object	object	object
countrycode	object	object	object	object	object	object	object
countryname	object	object	object	object	object	object	object
countryshortname	object	object	object	object	object	object	object
docty	object	object	object	object	object	object	object
envassamentcategorycode	object	object	object	object	object	object	object
grantant	int64	object	object	object	object	object	object
ibrdcommant	int64	object	object	object	object	object	object
id	object	object	object	object	object	object	object
idacommant	int64	object	object	object	object	object	object
impagency	object	object	object	object	object	object	object
lendinginstr	object	object	object	object	object	object	object
lendinginstrtype	object	object	object	object	object	object	object
lendprojectcost	int64	object	object	object	object	object	object
majorsector_percent	object	object	object	object	object	object	object
mjsector_namecode	object	object	object	object	object	object	object
mjtheme	object	object	object	object	object	object	object
mjtheme_namecode	object	object	object	object	object	object	object
mjthemecode	object	object	object	object	object	object	object
prodlinetext	object	object	object	object	object	object	object
productlinetype	object	object	object	object	object	object	object
project_abstract	object	object	object	object	object	object	object
project_name	object	object	object	object	object	object	object
projectdocs	object	object	object	object	object	object	object
projectfinancialtype	object	object	object	object	object	object	object
projectstatusdisplay	object	object	object	object	object	object	object
regionname	object	object	object	object	object	object	object
sector	object	object	object	object	object	object	object
sector1	object	object	object	object	object	object	object
sector2	object	object	object	object	object	object	object
sector3	object	object	object	object	object	object	object
sector4	object	object	object	object	object	object	object
sector_namecode	object	object	object	object	object	object	object
sectorcode	object	object	object	object	object	object	object
source	object	object	object	object	object	object	object
status	object	object	object	object	object	object	object
supplementprojectflg	object	object	object	object	object	object	object
them1	object	object	object	object	object	object	object
theme_namecode	object	object	object	object	object	object	object
themecode	object	object	object	object	object	object	object
totalamt	int64	object	object	object	object	object	object
totalcommant	int64	object	object	object	object	object	object
url	object	object	object	object	object	object	object
dtype: object							

10 countries with most projects

```
In [26]: df['countryname'].value_counts().head(10)
```

```
Out[26]:
```

countryname	count
People's Republic of China	19
Republic of Indonesia	19
Socialist Republic of Vietnam	17
Republic of India	16
Republic of Yemen	13
Kingdom of Morocco	12
Nepal	12
People's Republic of Bangladesh	12
Republic of Mozambique	11
Africa	11

Name: countryname, dtype: int64

Top 10 major project themes

```
In [27]: df['mjtheme_namecode'].value_counts().head(10)
```

```
Out[27]:
```

mjtheme_namecode	count
{'code': '11', 'name': 'Environment and natural resources management', ('code': '11', 'name': 'Environment and natural resources management')}	12
{'code': '8', 'name': 'Human development', ('code': '11', 'name': '')}	11
{'code': '8', 'name': 'Human development', ('code': '8', 'name': 'Human development')}	8
{'code': '4', 'name': 'Financial and private sector development', ('code': '4', 'name': 'Financial and private sector development')}	6
{'code': '2', 'name': 'Public sector governance', ('code': '2', 'name': 'Public sector governance')}	6
{'code': '2', 'name': 'Public sector governance', ('code': '2', 'name': 'Public sector governance')}	6
{'code': '8', 'name': 'Human development', ('code': '7', 'name': 'Social dev/gender/inclusion')}	5
{'code': '8', 'name': 'Human development', ('code': '8', 'name': 'Human development')}	5
{'code': '8', 'name': 'Human development', ('code': '8', 'name': 'Human development')}	5
{'code': '11', 'name': 'Environment and natural resources management', ('code': '11', 'name': '')}	5
{'code': '11', 'name': 'Environment and natural resources management', ('code': '4', 'name': '')}	5
{'code': '4', 'name': 'Financial and private sector development', ('code': '5', 'name': 'Trade and integration')}	5
Name: mjtheme_namecode, dtype: int64	

```
In [28]: import pandas as pd
import json
from pandas.io.json import json_normalize
```

```
In [29]: with open('C:\Users\User\Desktop\spring board\data_wrangling_json\data/world_bank_projects.json') as f:
    raw = json.load(f)
    df_themes = json_normalize(raw, 'mjtheme_namecode', ['id'])
```

Created DataFrame with missing values

```
In [30]: df_themes.head()
```

```
Out[30]:
```

	code	name	id
0	8	Human development	P129828
1	11		P129828
2	1	Economic management	P144674
3	6	Social protection and risk management	P144674
4	5	Trade and integration	P145310

```
In [31]: df_themes.sort_values('code', inplace = True)
df_themes
```

```
Out[31]:
```

	code	name	id
458	1	Economic management	P130925
1235	1	Economic management	P131440
1230	1	Economic management	P129465
1229	1	Economic management	P129465
1218	1	Economic management	P130824
900	1	Economic management	P143819
648	1	Economic management	P129625
647	1	Economic management	P129625
1078	1	Economic management	P131234
1206	1	Economic management	P133255
1437	1		P130412
357	1	Economic management	P130459
363	1		P144030
1010	1	Economic management	P133706
784	1	Economic management	P110836
1024	1	Economic management	P124114
88	1	Economic management	P132425
1045	1		