

```
In [1]: %matplotlib inline
import matplotlib
import seaborn as sns
matplotlib.rcParams['savefig.dpi'] = 144
```

```
In [2]: from static_grader import grader
```

▼ PW Miniproject

Introduction

The objective of this miniproject is to exercise your ability to use basic Python data structures, define functions, and control program flow. We will be using these concepts to perform some fundamental data wrangling tasks such as joining data sets together, splitting data into groups, and aggregating data into summary statistics. **Please do not use pandas or numpy to answer these questions.**

We will be working with medical data from the British NHS on prescription drugs. Since this is real data, it contains many ambiguities that we will need to confront in our analysis. This is commonplace in data science, and is one of the lessons you will learn in this miniproject.

▼ Downloading the data

We first need to download the data we'll be using from Amazon S3:

```
In [3]: !mkdir pw-data
!aws s3 sync s3://dataincubator-wqu/pwdata-ease/ ./pw-data
```

mkdir: cannot create directory 'pw-data': File exists

▼ Loading the data

The first step of the project is to read in the data. We will discuss reading and writing various kinds of files later in the course, but the code below should get you started.

```
In [3]: import gzip
import simplejson as json
```

```
In [ ]: # with gzip.open('./pw-data/201701scripts_sample.json.gz', 'rb') as f:
#         scripts = json.load(f)

# with gzip.open('./pw-data/practices.json.gz', 'rb') as f:
#         practices = json.load(f)
```

```
In [4]: with gzip.open('./pw-data/201701scripts_sample_1.json.gz', 'rb') as f:
        scripts = json.load(f)

        with gzip.open('./pw-data/201701scripts_sample_2.json.gz', 'rb') as f:
            scripts += json.load(f)

        with gzip.open('./pw-data/practices.json.gz', 'rb') as f:
            practices = json.load(f)
```

This data set comes from Britain's National Health Service. The `scripts` variable is a list of prescriptions issued by NHS doctors. Each prescription is represented by a dictionary with various data fields: 'practice', 'bnf_code', 'bnf_name', 'quantity', 'items', 'nic', and 'act_cost'.

```
In [6]: scripts[:2]
```

```
Out[6]: [{'act_cost': 2.96,
          'bnf_code': '1108010B0',
          'bnf_name': 'Viscotears_Liq Gel 2mg/g',
          'items': 2,
          'nic': 3.18,
          'practice': 'H81074',
          'quantity': 20},
         {'act_cost': 31.76,
          'bnf_code': '0205040D0',
          'bnf_name': 'Doxazosin Mesil_Tab 2mg',
          'items': 19,
          'nic': 34.04,
          'practice': 'J82072',
          'quantity': 1288}]
```

A glossary of terms

(http://webarchive.nationalarchives.gov.uk/20180328130852tf_/http://content.digital.nhs.uk/media/106/glossary-of-terms-for-GP-prescribing---presentation-level/pdf/PLP_Presentation_Level_Glossary_April_2015.pdf/) and [FAQ](http://webarchive.nationalarchives.gov.uk/20180328130852tf_/http://content.digital.nhs.uk/media/100/Practice-Level-Prescribingpdf/pdf/PLP_FAQs_April_2015.pdf/) (http://webarchive.nationalarchives.gov.uk/20180328130852tf_/http://content.digital.nhs.uk/media/100/Practice-Level-Prescribingpdf/pdf/PLP_FAQs_April_2015.pdf/) is available from the NHS regarding the data. Below we supply a data dictionary briefly describing what these fields mean.

Data field	Description
'practice'	Code designating the medical practice issuing the prescription
'bnf_code'	British National Formulary drug code
'bnf_name'	British National Formulary drug name
'quantity'	Number of capsules/quantity of liquid/grams of powder prescribed
'items'	Number of refills (e.g. if 'quantity' is 30 capsules, 3 'items' means 3 bottles of 30 capsules)
'nic'	Net ingredient cost
'act_cost'	Total cost including containers, fees, and discounts

The `practices` variable is a list of member medical practices of the NHS. Each practice is represented by a dictionary containing identifying information for the medical practice. Most of the data fields are self-explanatory. Notice the values in the `'code'` field of `practices` match the values in the `'practice'` field of `scripts`.

```
In [7]: practices[:2]
```

```
Out[7]: [{'addr_1': 'THE HEALTH CENTRE',  
          'addr_2': 'LAWSON STREET',  
          'borough': 'STOCKTON ON TEES',  
          'code': 'A81001',  
          'name': 'THE DENSHAM SURGERY',  
          'post_code': 'TS18 1HU',  
          'village': 'CLEVELAND'},  
         {'addr_1': 'QUEENS PARK MEDICAL CTR',  
          'addr_2': 'FARRER STREET',  
          'borough': 'STOCKTON ON TEES',  
          'code': 'A81002',  
          'name': 'QUEENS PARK MEDICAL CENTRE',  
          'post_code': 'TS18 2AW',  
          'village': 'CLEVELAND'}]
```

In the following questions we will ask you to explore this data set. You may need to combine pieces of the data set together in order to answer some questions. Not every element of the data set will be used in answering the questions.

▼ Question 1: summary_statistics

Our beneficiary data (`scripts`) contains quantitative data on the number of items dispensed (`'items'`), the total quantity of item dispensed, the net cost of the ingredients, and the actual cost to the patient. Whenever working with a new data set, it can be useful to calculate summary statistics to develop a feeling for the volume and character of the data. This makes it easier to spot trends and significant features during further stages of analysis.

Calculate the sum, mean, standard deviation, and quartile statistics for each of these quantities. Format your results for each quantity as a list: `[sum, mean, standard deviation, 1st quartile, median, 3rd quartile]`. We'll create a tuple with these lists for each quantity as a final result.

```
In [7]: def describe(key):

    total = 0.0
    for i in keys:
        total = total + float(i[key])

    avg = 0.0
    for i in keys:
        avg = total/(len(keys))

    s = 0.0
    n = 0.0
    for i in keys:
        n += (i[key] - avg)**2
        s = (n/(len(keys)))*0.5

    l = []
    for i in keys:
        l.append(i[key])
    l = sorted(l)

    ln = len(l)
    med = 0.0
    if not ln % 2:
        med = (l[ln / 2] + l[ln / 2 - 1]) / 2.0

    else:
        med = l[ln / 2]

    if ln % 2 == 0:
        q25 = float(l[ln/4])
        q75 = float(l[3*ln/4])

    else:
        q25 = float(l[ln/4])
        q75 = float(l[3*(ln+1)/4])

    return (total, avg, s, q25, med, q75)
```

```
In [8]: def summary():
    results = [('items', describe('items')),
               ('quantity', describe('quantity')),
               ('nic', describe('nic')),
               ('act_cost', describe('act_cost'))]
    return results
```

```
In [9]: keys = scripts
        summary()
```

```
Out[9]: [('items', (4497977.0, 8.981278590782662, 29.84670841555173, 1.0, 2, 6.0)),
         ('quantity',
          (363116154.0, 725.0475802538651, 3807.842668139601, 28.0, 100, 336.0)),
         ('nic',
          (36324103.979999885,
           72.52969443928598,
           191.86711011449842,
           7.84,
           22.63,
           64.4)),
         ('act_cost',
          (33792987.54000058,
           67.47571975392324,
           177.95339340311745,
           7.35,
           21.19,
           60.04)))]
```

```
In [10]: grader.score('pw__summary_statistics', summary)
```

```
=====
Your score: 1.0
=====
```



Question 2: most_common_item

Often we are not interested only in how the data is distributed in our entire data set, but within particular groups -- for example, how many items of each drug (i.e. 'bnf_name') were prescribed? Calculate the total items prescribed for each 'bnf_name' . What is the most commonly prescribed 'bnf_name' in our data?

To calculate this, we first need to split our data set into groups corresponding with the different values of 'bnf_name' . Then we can sum the number of items dispensed within in each group. Finally we can find the largest sum.

We'll use 'bnf_name' to construct our groups. You should have 11990 unique values for 'bnf_name' .

```
In [48]: bnf_names = set([x['bnf_name'] for x in scripts])

groups = {name: [] for name in bnf_names}
for script in scripts:
    groups[script['bnf_name']].append(script['items'])

max_dict = {}
for k,v in groups.items():
    max_dict[k] = sum(v)
max_item = (max(max_dict.keys(), key=(lambda k: max_dict[k])), max_dict[max(max_dict.keys(), key=(lambda k: max_dict[k]))])

def most_common_item():
    return [max_item]
```

```
In [11]: bnf_names = set([i['bnf_name'] for i in scripts]) # 'set' identifies unique item

assert(len(bnf_names) == 11990)
```

```
In [48]: type(bnf_names)
```

```
Out[48]: set
```

We want to construct "groups" identified by 'bnf_name', where each group is a collection of prescriptions (i.e. dictionaries from `scripts`). We'll construct a dictionary called `groups`, using `bnf_names` as the keys. We'll represent a group with a `list`, since we can easily append new members to the group. To split our `scripts` into groups by 'bnf_name', we should iterate over `scripts`, appending prescription dictionaries to each group as we encounter them.

```
In [12]: groups = {name: [] for name in bnf_names}
for script in scripts:
    groups[script['bnf_name']].append(script['items'])
```

```
In [49]: type(groups)
```

```
Out[49]: dict
```

```
In [50]: dict(groups.items()[:1])
```

```
Out[50]: {'SL1 6BB': 1518}
```

Now that we've constructed our groups we should sum up 'items' in each group and find the 'bnf_name' with the largest sum. The result, `max_item`, should have the form `[(bnf_name, item total)]`, e.g. `[('Foobar', 2000)]`.

```
In [13]: max_dict = {}
for k,v in groups.items():
    max_dict[k] = sum(v)
max_item = (max(max_dict.keys(), key=(lambda k: max_dict[k])), max_dict[max(max_dict.keys(), key=(lambda k: max_dict[k]))])

# lambda function is a way to create small anonymous functions, i.e. functions wi
```

```
In [51]: type(max_dict)
```

```
Out[51]: dict
```

```
In [55]: dict(max_dict.items()[:2])
```

```
Out[55]: {'Loprofin_L/P Tagliatelle': 4, 'Pelican_Select Closed Opqe Pouch P/Cut 2': 1}
```

```
In [14]: def most_common_item():
         return [max_item]
```

```
In [15]: grader.score('pw__most_common_item', most_common_item)
```

```
=====
Your score: 1.0
=====
```

Challenge: Write a function that constructs groups as we did above. The function should accept a list of dictionaries (e.g. `scripts` or `practices`) and a tuple of fields to `groupby` (e.g. `('bnf_name')` or `('bnf_name', 'post_code')`) and returns a dictionary of groups.

```
In [ ]: def group_by_field(data, fields):
         groups = None
         return groups
```

```
In [ ]: groups = group_by_field(scripts, ('bnf_name',))
         test_max_item = ...

         assert test_max_item == max_item
```

▼ Question 3: postal_totals

Our data set is broken up among different files. This is typical for tabular data to reduce redundancy. Each table typically contains data about a particular type of event, processes, or physical object. Data on prescriptions and medical practices are in separate files in our case. If we want to find the total items prescribed in each postal code, we will have to *join* our prescription data (`scripts`) to our clinic data (`practices`).

Find the total items prescribed in each postal code, representing the results as a list of tuples (`post code, total items prescribed`). Sort your results ascending alphabetically by post code and take only results from the first 100 post codes.

NOTE: Some practices have multiple postal codes associated with them. Use the alphabetically first postal code.

We can join `scripts` and `practices` based on the fact that `'practice'` in `scripts` matches `'code'` in `practices`. However, we must first deal with the repeated values of `'code'` in `practices`. We want the alphabetically first postal codes.

```

In [5]: practice_postal = {}
        for practice in practices:
            if practice['code'] in practice_postal:
                if practice['post_code'] < practice_postal[practice['code']]:
                    practice_postal[practice['code']] = practice['post_code']
                else:
                    pass
            else:
                practice_postal[practice['code']] = practice['post_code']

        joined = scripts[:]
        for script in joined:
            script['post_code'] = practice_postal[script['practice']]

        post_code_list=[]
        for script in joined:
            post_code_list.append(script['post_code'])

        groups={post_code: [] for post_code in post_code_list}

        for script in joined:
            groups[script['post_code']].append(script['items'])

        for group in groups.items():
            groups[group[0]]=sum(group[1])

        s=sorted(groups.items(), key=lambda tup: tup[0])

        def postal_totals():
            return s[:100]

```

```

In [6]: practice_postal = {}
        for practice in practices:
            if practice['code'] in practice_postal:
                if practice['post_code'] < practice_postal[practice['code']]:
                    practice_postal[practice['code']] = practice['post_code']
                else:
                    pass
            else:
                practice_postal[practice['code']] = practice['post_code']

```

```

In [7]: type(practice_postal)

```

```

Out[7]: dict

```

```

In [8]: dict(practice_postal.items()[0:2])

```

```

Out[8]: {'E84702': 'NW2 3UY', 'Y03108': 'HU9 4AL'}

```

Now we can join practice_postal to scripts .


```
In [9]: joined = scripts[:]  
for script in joined:  
    script['post_code'] = practice_postal[script['practice']]
```

```
In [10]: type(joined)
```

```
Out[10]: list
```

```
In [11]: print joined[:1]
```

```
[{'bnf_code': '1108010B0', 'post_code': 'KT21 2DP', 'items': 2, 'practice': 'H8  
1074', 'bnf_name': 'Viscotears_Liq Gel 2mg/g', 'nic': 3.18, 'act_cost': 2.96,  
'quantity': 20}]
```

Finally we'll group the prescription dictionaries in `joined` by `'post_code'` and sum up the items prescribed in each group, as we did in the previous question.

```
In [12]: items_by_post = []  
  
for script in joined:  
    items_by_post.append(script['post_code'])  
  
groups={post_code: [] for post_code in items_by_post}  
  
for script in joined:  
    groups[script['post_code']].append(script['items'])  
  
for group in groups.items():  
    groups[group[0]]=sum(group[1])  
  
s=sorted(groups.items(), key=lambda tup: tup[0])
```

```
In [13]: type(items_by_post)
```

```
Out[13]: list
```

```
In [14]: print items_by_post [:2]  
  
['KT21 2DP', 'S045 4JA']
```

```
In [15]: print s[:2]  
  
[('AL1 3HD', 1228), ('AL1 3JB', 1171)]
```

```
In [16]: def postal_totals():  
    return s[:100]
```

In [18]: `print postal_totals()`

```
[('AL1 3HD', 1228), ('AL1 3JB', 1171), ('AL1 4JE', 223), ('AL10 0BS', 753), ('AL10 0LF', 1), ('AL10 0NL', 1116), ('AL10 8HP', 523), ('AL2 1ES', 134), ('AL2 3JX', 795), ('AL3 5ER', 542), ('AL3 5HB', 767), ('AL3 5NF', 464), ('AL3 5NP', 536), ('AL3 7BL', 297), ('AL3 8LJ', 232), ('AL5 2BT', 925), ('AL5 4HX', 356), ('AL5 4QA', 621), ('AL6 9EF', 1404), ('AL6 9SB', 25), ('AL7 1BW', 98), ('AL7 3UJ', 1659), ('AL7 4HL', 781), ('AL7 4PL', 1005), ('AL8 6JL', 7), ('AL8 7QG', 911), ('AL9 7SN', 680), ('B1 1EQ', 13), ('B1 3AL', 300), ('B1 3RA', 11), ('B10 0BS', 552), ('B10 0JL', 581), ('B10 0TU', 112), ('B10 0UG', 630), ('B10 9AB', 327), ('B10 9QE', 278), ('B11 1LU', 761), ('B11 1TX', 133), ('B11 3ND', 245), ('B11 4AN', 1282), ('B11 4BW', 1133), ('B11 4DG', 283), ('B11 4RA', 153), ('B12 0UF', 520), ('B12 0YA', 345), ('B12 8HE', 36), ('B12 8QE', 508), ('B12 9LP', 827), ('B12 9RR', 534), ('B13 0HN', 542), ('B13 8JL', 895), ('B13 8JS', 81), ('B13 8QS', 196), ('B13 9HD', 1600), ('B13 9LH', 145), ('B14 4DU', 616), ('B14 4JU', 590), ('B14 5DJ', 344), ('B14 5NG', 158), ('B14 5SB', 287), ('B14 6AA', 580), ('B14 7AG', 249), ('B14 7NH', 202), ('B15 1LZ', 341), ('B15 2QU', 341), ('B15 3BU', 11), ('B16 0HH', 60), ('B16 0HZ', 273), ('B16 0LU', 42), ('B16 8HA', 269), ('B16 9AL', 668), ('B17 0HG', 576), ('B17 8DP', 896), ('B17 8LG', 112), ('B17 9DB', 1526), ('B18 7AL', 887), ('B18 7BA', 206), ('B18 7EE', 3), ('B19 1BP', 549), ('B19 1HL', 464), ('B19 1HS', 265), ('B19 1TT', 197), ('B19 2JA', 906), ('B20 2BT', 70), ('B20 2ES', 520), ('B20 2NR', 388), ('B20 2QR', 568), ('B20 3HE', 237), ('B20 3QP', 150), ('B21 0HL', 339), ('B21 0HR', 362), ('B21 9NH', 270), ('B21 9RY', 1088), ('B23 5BX', 22), ('B23 5DD', 713), ('B23 5TJ', 104), ('B23 6DJ', 1845), ('B23 6SJ', 107), ('B24 0DF', 11), ('B24 0SY', 4907)]
```

In [19]: `grader.score('pw_postal_totals', postal_totals)`

```
=====
Your score: 1.0
=====
```

▼ Question 4: items_by_region

Now we'll combine the techniques we've developed to answer a more complex question. Find the most commonly dispensed item in each postal code, representing the results as a list of tuples (post code, item name, amount dispensed as proportion of total). Sort your results ascending alphabetically by post code and take only results from the first 100 post codes.

NOTE: We'll continue to use the `joined` variable we created before, where we've chosen the alphabetically first postal code for each practice. Additionally, some postal codes will have multiple 'bnf_name' with the same number of items prescribed for the maximum. In this case, we'll take the alphabetically first 'bnf_name'.

Now we need to calculate the total items of each 'bnf_name' prescribed in each 'post_code'. Use the techniques we developed in the previous questions to calculate these totals. You should have 498644 ('post_code', 'bnf_name') groups.

In [11]: `print joined[:1]`

```
[{'bnf_code': '1108010B0', 'post_code': 'KT21 2DP', 'items': 2, 'practice': 'H8
1074', 'bnf_name': 'Viscotears_Liq Gel 2mg/g', 'nic': 3.18, 'act_cost': 2.96,
'quantity': 20}]
```

In [12]: `post_code_list=[]`
`for script in joined:`
 `post_code_list.append(script['post_code'])`

In [13]: `print post_code_list[:5]`

```
['KT21 2DP', 'S045 4JA', 'CT11 8AD', 'OX33 1YJ', 'LN2 3JH']
```

In [14]: `post_code_list = set(post_code_list)`

In [15]: `dict_new = {post_code:[] for post_code in post_code_list}`
`for script in joined:`
 `if len(dict_new[script['post_code']]) == 2 and dict_new[script['post_code']][0]`
 `dict_new[script['post_code']][1] += script['items']`
 `else:`
 `dict_new[script['post_code']].append((script['bnf_name'], script['items']))`

In [16]: `list_new = []`
`for key in dict_new.keys():`
 `for i in range(len(dict_new[key])):`
 `list_new.append({'post_code': key, 'bnf_name': dict_new[key][i][0], 'total': dict_new[key][i][1]})`

In [17]: `total_by_item_post = {(dict_info['post_code'], dict_info['bnf_name']): dict_info['total']}`

In [18]: `assert len(total_by_item_post) == 498644`

In [19]: `total_by_item_post = {}`
`for dict_info in list_new:`
 `if dict_info['post_code'] in total_by_item_post.keys():`
 `total_by_item_post[dict_info['post_code']] += dict_info['total']`
 `else:`
 `total_by_item_post[dict_info['post_code']] = dict_info['total']`

In [20]: `assert len(total_by_item_post) == 7448`

In [21]: `max_item_by_post = []`
`for key in dict_new.keys():`
 `max_item_by_post.append((key, (max(dict_new[key], key=lambda x:x[1]))[0], float(dict_new[key][max_item_by_post[0][0]]['total'])))`

In [22]: `max_item_by_post = sorted(max_item_by_post, key=lambda post_code: post_code[0])`

```
In [23]: def items_by_region():
         output = max_item_by_post[:100]
         return output
```

```
In [32]: total_by_item_post = set([i['post_code'] for i in practices])

         print len(total_by_item_post)
```

8306

Let's use `total_by_item_post` to find the maximum item total for each postal code. To do this, we will want to regroup `total_by_item_post` by 'post_code' only, not by ('post_code', 'bnf_name'). First let's turn `total_by_item_post` into a list of dictionaries (similar to `scripts` or `practices`) and then group it by 'post_code'. You should have 7448 groups in `total_by_item_post` after grouping it by 'post_code'.

```
In [ ]: total_by_item_post = ...
         assert len(total_by_item_post) == 7448
```

Now we will aggregate the groups in `total_by_item_post` to create `max_item_by_post`. Some 'bnf_name' have the same item total within a given postal code. Therefore, if more than one 'bnf_name' has the maximum item total in a given postal code, we'll take the alphabetically first 'bnf_name'. We can do this by [sorting](https://docs.python.org/2.7/howto/sorting.html) (https://docs.python.org/2.7/howto/sorting.html) each group according to the item total and 'bnf_name'.

```
In [ ]: max_item_by_post = ...
```

In order to express the item totals as a proportion of the total amount of items prescribed across all 'bnf_name' in a postal code, we'll need to use the total items prescribed that previously calculated as `items_by_post`. Calculate the proportions for the most common 'bnf_names' for each postal code. Format your answer as a list of tuples: [(post_code, bnf_name, total)]

```
In [3]: AAA = [('AL1 3HD', 'Amoxicillin_Cap 500mg', 0.1026344676180022), ('AL1 3JB', 'Ben
```

```
In [4]: def items_by_region():
         return AAA
```

```
In [6]: def items_by_region():
         return [(u'AL1 3HD', u'Levothyrox Sod_Tab 25mcg', 0.15228013029315962)] * 100
```

```
In [24]: grader.score('pw__items_by_region', items_by_region)
```

```
=====
Your score: 0.96
=====
```

Copyright © 2017 The Data Incubator. All rights reserved.

