

BeeChecklist_solution

December 31, 2018

0.1 Solution of 5.9.1, Bee Checklist

First of all, we import the two modules we'll need to read the csv file, and to use regular expressions:

```
In [1]: import csv
import re
```

Then, we read the file, and store the columns Scientific Name and Taxon Author in two lists:

```
In [2]: with open('../data/bee_list.txt') as f:
        csvr = csv.DictReader(f, delimiter = '\t')
        species = []
        authors = []
        for r in csvr:
            species.append(r['Scientific Name'])
            authors.append(r['Taxon Author'])
```

How many species?

```
In [3]: len(species)
```

```
Out[3]: 19508
```

```
In [4]: len(authors)
```

```
Out[4]: 19508
```

Pick one of the authors element to use for testing. Choose one that is quite complicated, such as the 38th element:

```
In [5]: au = authors[37]
```

```
In [6]: au
```

```
Out[6]: 'Tadauchi, Hirashima & Matsumura, 1987'
```

Now we need to build a regular expression. After some twiddling, you should end up with something like this, which captures the authors in one group, and the year in another group:

```
In [7]: my_reg = re.compile(r'\((?([\w\s,\.\- \&]*),\s(\d{4}))\)?')
# Translation
# \(? -> open parenthesis (or not)
# ([\w\s,\.\- \&]+) -> the first group is the list of authors
#                               which can contain \w (word character)
#                               \s (space) \. (dot) \- (dash) \& (ampersand)
# ,\s -> followed by comma and space
# (\d{4}) -> the second group is the year, 4 digits
# \)? -> potentially, close parenthesis
```

Test the expression

```
In [8]: re.findall(my_reg, au)
```

```
Out[8]: [('Tadauchi, Hirashima & Matsumura', '1987')]
```

Now we write a function that uses the regular expression to extract an author list (useful when there are multiple authors), and the year

```
In [9]: def extract_list_au_year(au):
    tmp = re.match(my_reg, au)
    authorlist = tmp.group(1)
    year = tmp.group(2)
    # split authors into a list using re.split
    authorlist = re.split(', | \& ', authorlist)
    # Translation: either separate using ', ' or '& '
    return [authorlist, year]
```

Let's see the output of this function:

```
In [10]: extract_list_au_year(au)
```

```
Out[10]: [['Tadauchi', 'Hirashima', 'Matsumura'], '1987']
```

Finally, let's build two dictionaries: - one tracking the number of times each year is mentioned in the database; - one tracking the number of times each author is mentioned

```
In [11]: dict_years = {}
    dict_authors = {}
    for au in authors:
        tmp = extract_list_au_year(au)
        for aunum in tmp[0]:
            if aunum in dict_authors.keys():
                dict_authors[aunum] = dict_authors[aunum] + 1
            else:
                dict_authors[aunum] = 1
        if tmp[1] in dict_years.keys():
            dict_years[tmp[1]] = dict_years[tmp[1]] + 1
        else:
            dict_years[tmp[1]] = 1
```

For example, these are all the authors:

```
In [12]: dict_authors
```

```
Out[12]: {'Michener': 455,  
          'Cockerell': 3394,  
          'Timberlake': 864,  
          'Zavortink': 7,  
          'Dours': 62,  
          'Eversmann': 41,  
          'Pérez': 309,  
          'Warncke': 507,  
          'Hirashima': 210,  
          'Viereck': 130,  
          'Morawitz': 513,  
          'LaBerge': 192,  
          'Frieze': 1330,  
          'Smith': 943,  
          'Scopoli': 8,  
          'Thorp': 19,  
          'Osytschnjuk': 74,  
          'Tadauchi': 76,  
          'Matsumura': 19,  
          'Popov': 86,  
          'Alfken': 149,  
          'Radoszkowski': 113,  
          'Rossi': 11,  
          'Perkins': 65,  
          'Robertson': 135,  
          'Benoist': 164,  
          'Miyanaaga': 2,  
          'Dawut': 1,  
          'Lebedev': 9,  
          'Stoeckhert': 10,  
          'Linsley': 59,  
          'MacSwain': 16,  
          'Cresson': 437,  
          'Kirby': 74,  
          'Ribble': 31,  
          'Cameron': 158,  
          'Blüthgen': 315,  
          'Kreichbaumer': 1,  
          'Scheuchl': 10,  
          'Gusenleitner': 39,  
          'Schmiedeknecht': 50,  
          'Mitchell': 270,  
          'Hedicke': 24,  
          'Xu': 37,
```

'Donovan': 24,
'Nurse': 53,
'Malloch': 3,
'Bouseman': 8,
'Panzer': 29,
'Schulthess': 1,
'Fabricius': 138,
'Vachal': 557,
'Larkin': 6,
'Stephens': 1,
'Bingham': 47,
'Wu': 213,
'Dalla Torre': 65,
'Strand': 249,
'Casad': 7,
'Pittioni': 7,
'Brullé': 33,
'Neff': 5,
'Yasumatsu': 53,
'Dubitzky': 12,
'Baker': 44,
'Schenck': 41,
'Linnaeus': 39,
'Christ': 4,
'Lepeletier': 185,
'Schwarz': 170,
'Theunert': 1,
'Morice': 11,
'Thomson': 9,
'Gribodo': 75,
'Mavromoustakis': 100,
'Erichson': 23,
'Provancher': 25,
'Lanham': 4,
'Dufour': 4,
'Kohl': 18,
'Haneda': 1,
'Giraud': 12,
'Grünwaldt': 7,
'Guiglia': 2,
'Ashmead': 30,
'Magretti': 18,
'Müller': 4,
'Patiny': 24,
'van der Vecht': 12,
'Imhoff': 8,
'Dunning': 3,
'Jaeger': 1,

'Pohl': 1,
'Atwood': 6,
'Latreille': 35,
'Zetterstedt': 5,
'Schönitzer': 2,
'Kim': 2,
'Fox': 25,
'Saunders': 22,
'Graenicher': 9,
'Kriechbaumer': 7,
'Noskiewicz': 75,
'Nylander': 37,
'Uchida': 3,
'De Stefani': 1,
'Schwenninger': 1,
'Engel': 89,
'Spinola': 108,
'Hazir': 1,
'Illiger': 9,
'Tamasana': 1,
'Meunier': 5,
'Ducke': 92,
'Hurd': 36,
'Schuberth': 1,
'Rozen': 42,
'Moure': 417,
'Urban': 282,
'Gerstäcker': 55,
'Ascher': 1,
'Griswold': 50,
'Klug': 35,
'Sichel': 22,
'Graf': 2,
'Seabra': 12,
'Compagnucci': 10,
'Shinn': 25,
'Gonzalez': 19,
'Holmberg': 181,
'Fowler': 7,
'Jørgensen': 38,
'Porter': 4,
'Swenk': 62,
'Crawford': 88,
'Ruz': 45,
'Brèthes': 42,
'Toro': 140,
'Herrera': 3,
'Cure': 11,

'Wittman': 1,
'Roig-Alsina': 71,
'Rodríguez': 3,
'Walker': 112,
'Eardley': 93,
'Brauns': 35,
'Schwammberger': 8,
'Mocsáry': 67,
'Lucas': 9,
'Aurivillius': 1,
'Richards': 7,
'Ortiz-Sánchez': 1,
'Michez': 14,
'Snelling': 72,
'Danforth': 4,
'Parker': 15,
'Say': 35,
'Stevens': 4,
'Schrottky': 246,
'Melo': 23,
'Tapia': 4,
'Vivallo': 3,
'Reed': 5,
'Ramos': 1,
'Lucas de Oliveira': 1,
'Schlindwein': 15,
'Chiappa': 3,
'Rayment': 209,
'Saussure': 28,
'Schulz': 13,
'Lieftinck': 108,
'Gussakovsky': 14,
'Priesner': 20,
'Brooks': 103,
'Pauly': 175,
'Buysson': 6,
'Meade-Waldo': 45,
'Fedtschenko': 29,
'W. F. Kirby': 13,
'de Villers': 1,
'Mariskovskaya': 1,
'Olivier': 14,
'Huard': 2,
'Pallas': 4,
'Banaszak': 1,
'Tkalcu': 86,
'Laboulbene': 1,
'Packard': 8,

'Westrich': 3,
'Romankova': 4,
'Stadelmann': 5,
'Enderlein': 19,
'Skorikov': 22,
'Franklin': 7,
'Dahlbom': 2,
'Vogt': 5,
'Seidl': 1,
'Handlirsch': 8,
'Frison': 13,
'Wahlberg': 1,
'Wang': 1,
'Guérin-Méneville': 19,
'Milliron': 1,
'DeGeer': 6,
'Bischoff': 10,
'Greene': 3,
'Schönherr': 1,
'Labougle': 1,
'Ayala': 22,
'Swederus': 2,
'Sladen': 3,
'Sparre-Schneider': 1,
'De Geer': 1,
'Schulthess-Rechberg': 1,
'Curtis': 4,
'Vollenhoven': 1,
'Geoffroy': 2,
'Herbst': 11,
'Blanchard': 4,
'Jensen-Haarup': 2,
'Azevedo': 2,
'Silveira': 7,
'Perty': 5,
'Burmeister': 5,
'Westwood': 11,
'Ruiz': 16,
'Oliveira': 8,
'Viana': 2,
'Zanella': 4,
'Castro': 1,
'Audinet-Serville': 11,
'Rodriguez': 2,
'Haliday': 7,
'Thiele': 1,
'Romand': 1,
'Bertoni': 7,

'Sitdikov': 5,
'Risch': 23,
'Dusmet y Alonso': 33,
'Fonscolombe': 3,
'Bradley': 1,
'Bär': 1,
'Dover': 3,
'de Gaulle': 1,
'Ogloblin': 7,
'Lovell': 17,
'Krombein': 11,
'Tucker': 1,
'Genaro': 16,
'Germar': 2,
'Kimsey': 10,
'Neves': 1,
'Nemésio': 5,
'González': 3,
'Gaiani': 3,
'Hoffmannsegg': 1,
'Dressler': 36,
'Sakagami': 38,
'Cheesman': 30,
'Hinojosa-Díaz': 8,
'Rebêlo': 5,
'Roubik': 13,
'Bembé': 1,
'Ramírez': 2,
'Rasmussen': 1,
'Skov': 1,
'Ospina-Torres': 2,
'Sandino-Franco': 1,
'Anjos-Silva': 1,
'González-Vaquero': 3,
'Almeida': 6,
'Rojas': 11,
'Taschenberg': 2,
'Forster': 1,
'Rightmyer': 40,
'Meyer': 31,
'Mackie': 2,
'de Beaumont': 1,
'Sivik': 1,
'Darchen': 2,
'Camargo': 75,
'Inoue': 3,
'Sepúlveda': 1,
'Marchi': 9,

'Jobiraj': 1,
'Narendran': 1,
'Bennett': 1,
'Kerr': 3,
'Lobo Segura': 1,
'Vélez': 1,
'Pedro': 25,
'Wille': 3,
'Puls': 1,
'Silvestri': 2,
'Franck': 1,
'Boongird': 1,
'Albuquerque': 8,
'Shanks': 15,
'Shrottky': 1,
'Dominique': 2,
'Harter-Marques': 2,
'Cunha': 1,
'Truylio': 1,
'Heinrich': 3,
'Jurine': 4,
'Kocourek': 4,
'Hicks': 2,
'Cooper': 3,
'Proshchalykin': 1,
'Lelej': 2,
'Ehrenfeld': 3,
'Sandhouse': 78,
'Schilling': 1,
'Giordani Soika': 1,
'Brues': 2,
'Tsuneki': 74,
'Lozinski': 2,
'Herrich-Schäffer': 8,
'Mazzucco': 1,
'Standfuss': 4,
'Benzi': 1,
'Moalif': 4,
'Broemeling': 5,
'Mitai': 3,
'Linné': 3,
'Evans': 4,
'Noble': 1,
'Rodeck': 1,
'Gmelin': 1,
'Arnold': 1,
'Ikudome': 8,
'Reyes': 21,

'Borges': 1,
'Syed': 2,
'Masi': 1,
'Houston': 50,
'Tierney': 1,
'H. S. Smith': 7,
'Kokujev': 1,
'Daly': 31,
'Chevrier': 1,
'Terzo': 12,
'Rasmont': 4,
'S. Lee': 1,
'Shiokawa': 4,
'Sickmann': 1,
'Sonan': 1,
'Maa': 22,
'Ritsema': 16,
'Maidl': 16,
'Vecht': 1,
'LeVeque': 10,
'Drury': 2,
'Leys': 2,
'Wiedemann': 1,
'Ponomareva': 2,
'Guilding': 1,
'Patton': 5,
'Trucco Aleman': 5,
'Maynard': 16,
'Kuhlmann': 65,
'Stephen': 12,
'Janvier': 2,
'Metz': 16,
'Ortiz': 1,
'Ornosa': 1,
'Fourcroy': 1,
'Verhoeff': 1,
'Schmidt': 1,
'Raw': 7,
'Gistel': 3,
'Titus': 10,
'Frey-Gessner': 1,
'Dubitzki': 1,
'Packer': 29,
'Cabezas': 15,
'Davies': 11,
'Vergara': 1,
'Deyrup': 1,
'Roberts': 32,

'Exley': 231,
'Cardale': 1,
'Bridwell': 16,
'Dathe': 38,
'Magnacca': 10,
'Erlandsson': 1,
'Förster': 21,
'Blackburn': 3,
'Hensen': 1,
'Motschulski': 1,
'Gorski': 1,
'Sumner': 1,
'Abe': 1,
'Moldenke': 37,
'Gibbs': 1,
'Philippi': 2,
'Willis': 3,
'Smith-Pardo': 10,
'Eickwort': 7,
'Coelho': 5,
'Ordway': 3,
'Doering': 1,
'Yanega': 6,
'Klein': 2,
'E. A. B. Almeida': 5,
'Gonçalves': 4,
'Wolcott': 1,
'M. C. de Almeida': 1,
'Laroca': 1,
'Ebmer': 215,
'Walckenaer': 1,
'Fan': 13,
'Bramson': 1,
'Pesenko': 56,
'Niu': 1,
'Huang': 1,
'Janjic': 1,
'McGinley': 24,
'Munakata': 1,
'Svensson': 1,
'Ellis': 23,
'Murao': 4,
'Munzinger': 2,
'Knerer': 5,
'Yáñez-Ordóñez': 3,
'Bytinski-Salz': 1,
'Wcislo': 3,
'Timmermann': 1,

'Takahashi': 1,
'Schrack': 2,
'Maeta': 9,
'Herrmann': 1,
'Davydova': 1,
'Godínez-García': 1,
'Stage': 4,
'Nobile': 7,
'Turrisi': 7,
'LaRoche': 3,
'Hagens': 10,
'R. P. Urban': 5,
'Wesmael': 1,
'Kato': 1,
'Itino': 1,
'Cross': 2,
'Wickwar': 1,
'Schletterer': 23,
'He': 2,
'Costa': 6,
'Astafurova': 2,
'Bohart': 47,
'Batra': 1,
'Whitehead': 16,
'Gupta': 16,
'Pasteels': 197,
'Simlote': 2,
'Perris': 4,
'Jaycox': 1,
'Rohwer': 2,
'Grigarick': 2,
'Stange': 3,
'Sharma': 6,
'Peters': 16,
'van der Zanden': 53,
'Rahman': 1,
'Tewari': 1,
'Sielfeld': 1,
'Schummel': 1,
'Fritz': 15,
'Hill': 1,
'Nagase': 3,
'Proschchalykin': 1,
'Dewitz': 1,
'King': 9,
'Rebmann': 18,
'Durante': 4,
'Schwimmer': 1,

```

'Schulten': 7,
'Ferton': 6,
'Abramovich': 2,
'Lucia': 2,
'Silveira et al.': 1,
'Stanek': 5,
'Le Goff': 1,
'Teunissen': 1,
'van Achterberg': 1,
'Rust': 3,
'Atanassov': 1,
'White': 4,
'Haeseler': 1,
'Robinaeau': 1,
'Dumeril': 1,
'von Schulthess': 1,
'Harris': 1,
'Rudow': 2,
'Quilis': 2,
'Steiner': 13,
'Vogel': 1}

```

0.1.1 What is the name of the author with most entries in the database?

We use the following strategy: - we find the maximum value in the dictionary - we use the function `index` to find to which entry is it associated - we find the corresponding author

```

In [13]: max_value_author = max(dict_authors.values())
         max_value_author

```

```

Out[13]: 3394

```

```

In [14]: which_index = list(dict_authors.values()).index(max_value_author)
         which_index

```

```

Out[14]: 1

```

An the winner is:

```

In [15]: list(dict_authors.keys())[which_index]

```

```

Out[15]: 'Cockerell'

```

0.1.2 Which year of publication is most represented in the database?

We use the same strategy to find that the golden year of bee publication is:

```

In [16]: max_value_year = max(dict_years.values())
         which_index = list(dict_years.values()).index(max_value_year)
         list(dict_years.keys())[which_index]

```

```

Out[16]: '1903'

```