

Big Data; 02 Practical - Python Basics 2

Ryan Greenup

July 23, 2020

Contents

Calculate Data Statistics	1
Sum Values	1
.1 Improvements	1
Minimum Value	2
Maximum Value	3
Vector Norm, Inner Product and Distance	3
Piping	3
Using Built ins	4
Inner Product	5
Distance	5
Vote Counting	6
Word Capitaliser	6
Parse File	7
Set up the Text File	7
Parse the Text File	8
.1 Read the Text File	8
.2 Return only Matching Data	9
Wrap it into a function	11
Parse Dictionary	13
grep	14
Top 10 Words	16
Create the Dictionary	17
:HTML: line-style when	

Calculate Data Statistics

Sum Values

So what we need to do is go through each value and tally it up, it is important however to return the value and then print it:

```
1  def mysum(x):
2      total = float(0)
3      for i in x:
4          total = total + float(i)
5      return total
6
7  value = mysum([1,2,3])
8  print(value)
```

6.0

Improvements

We could however improve this by using a try / except test, in the event that a non-numerical list is provided:

```
1  def mysum(x):
2      total = float(0)
3      for i in x:
4          try:
5              total = total + float(i)
6          except:
7              print("The Values of the list must be numeric")
8              print("Discarding Value")
9      return total
10
11 value = mysum([1,2,3, "apple"])
12 print(value)
```

The Values of the list must be numeric

6.0

Minimum Value

Take the first item of the list as a candidate, for every item in the list, compare it to the candidate, if the next value is bigger that will become the new candidate, finally the candidate will be the maximum value.

Just like above we use a try / except to prevent issues.

```

1  def mymax(list_of_vals):
2      candidate = list_of_vals[0]  # Unlike R/Mathematica/Julia, python
    ↪ starts from 0.
3      for i in list_of_vals:
4          try:
5              if i > candidate:
6                  candidate = float(i)
7          except:
8              print("The list items must be numeric")
9              print("Discarding Value")
10     return candidate
11
12 print(mymax([4, 6, 2, 5, 7, 3, 8, "john", -9]))

```

The list items must be numeric
Discarding Value
8.0

Maximum Value

Same as above, just remember to:

- wrap in float() as appropriate
- print the function call.

```

1  def mymin(thelist):
2      candidate = thelist[0]
3      for i in thelist:
4          try:
5              if float(i) > float(candidate):
6                  candidate = i
7          except:
8              print("list items must be numeric, discarding value")
9      return candidate
10
11 value = mymin([1, 5, 3, "apple", 8, 2, -9])
12 print(value)

```

list items must be numeric, discarding value
8

Vector Norm, Inner Product and Distance

Piping

The *Toolz* Module gives something very similar to piping in bash / julia / R

Using Built ins

```
1  import math as mt
2  import copy
3
4  ## because python counts from 0 indexing is confusing,
5  ## the count will come back as 4, but the indexes will be 0, 1, 2 and 3.
6
7  def getNorm(x):
8      total = 0
9      for i in range(len(x)):
10         total=x[i]**2+total
11     return mt.sqrt(total)
12
13 print(len([1,2,3,4]))
14
15
16 xvec = [0, 1, 2, 3, 4]
17 yvec = [4, 3, 2, 1, 0]
18
19 norm = getNorm(xvec);          print(norm)
```

4
5.477225575051661

Inner Product

```
1 def getInnerProd(x, y):
2     z = copy.deepcopy(x)    ## Careful, you need to copy, not just
    ↪ assign
3     if len(x) == len(y):
4         for i in range(len(x)):
5             z[i] = x[i]*y[i]
6         return sum(z)
7     else:
8         print("The vectors must have the same dimension")
9 xvec = [0, 1, 2, 3, 4]
10 yvec = [4, 3, 2, 1, 0]
11
12 norm = getNorm(xvec);        print(norm)
13 norm = getNorm(yvec);        print(norm)
14 prod = getInnerProd(xvec, yvec); print(prod)
```

5.477225575051661

5.477225575051661

10

6.324555320336759

Distance

```
1 def getDist(x, y):
2     if len(x) == len(y):
3         z = mt.sqrt(getNorm(x)**2 + getNorm(y)**2 - 2 * getInnerProd(x,
    ↪ y))
4         return z
5     else:
6         print("The vectors must have the same dimension")
7
8
9 xvec = [0, 1, 2, 3, 4]
10 yvec = [4, 3, 2, 1, 0]
11
12 norm = getNorm(xvec);        print(norm)
13 norm = getNorm(yvec);        print(norm)
14 prod = getInnerProd(xvec, yvec); print(prod)
15 dist = getDist(xvec, yvec);   print(dist)
```

Vote Counting

```
1 votes = "N , Y, Y,N,n , N , N N ,n ,y, n,N,Y, y,Y,N , N , n ,y,N"
2 def countVotes(ballot):
3     ballot = ballot.replace(",","").replace(" ", "").upper()
4     neg_ballots = ballot.count("N")
5     pos_ballots = ballot.count("Y")
6
7     # Could also have used a loop
8     print(str(pos_ballots) + " Yes votes and " + str(neg_ballots) + "
9         ↪ No votes")
10
11 countVotes(votes)
12 votes = ",,yyyyn,,y,y,nn,y,,nn,y"
```

7 Yes votes and 13 No votes

Word Capitaliser

```
1 def capitalise(sentence):
2     ## Split the words into a list
3     wordsList = sentence.split()
4     ## These are escape words
5     EscWords = ["am", "a", "an", "the", "am", "is", "are", "and", "of",
6         ↪ "in", "on", "with", "from", "to"]
7     ## The number of words starting from 0
8     for i in range(len(wordsList)-1):
9         ## if not in the escape words
10        if i not in EscWords:
11            ## replace the ith word for a capitalized one
12            wordsList[i] = wordsList[i].capitalize()
13        ## Take a space and use it to join the list together
14        sentence_Capitalized = " ".join(wordsList)
15        ## Print the output
16        print(sentence_Capitalized)
17        return sentence_Capitalized
18
19 capitalise("The quick brown fox jumped over the lazy dogs")
```

The Quick Brown Fox Jumped Over The Lazy dogs

Parse File

Set up the Text File

Take the following text:

Unit ID, unit name, course name 301046, Big Data, MICT 300581, Programming Techniques, BICT 300144, Object Oriented Analysis, BICT 300103, Data Structure, BCS 300147, Object Oriented Programming, BCS 300569, Computer Security, BIS 301044, Data Science, MICT 300582, Technologies for Web Applications, BICT

Let's write it to a file:

```
1  pwd
2  ls
```

```
1  scemunits = """Unit ID, unit name, course name
2  301046, Big Data, MICT
3  300581, Programming Techniques, BICT
4  300144, Object Oriented Analysis, BICT
5  300103, Data Structure, BCS
6  300147, Object Oriented Programming, BCS
7  300569, Computer Security, BIS
8  301044, Data Science, MICT
9  300582, Technologies for Web Applications, BICT"""
10
11 def writeTextFile(text, filename) :
12     f = open(filename, 'w')
13     for i in text.split('\n'):
14         f.writelines(i+'\n')
15     f.close()
16
17 writeTextFile(scemunits, 'scemunits.txt')
```

In order to check that worked we can run cat from *Python*:

```
1  import subprocess
2  MyCommand = "cat scemunits.txt"
3  scemunits_txt = subprocess.run(MyCommand.split(), capture_output = True)
4  print(scemunits_txt)
```

CompletedProcess(args=['cat', 'scemunits.txt'], returncode=0, stdout=b'Unit ID, unit name, course name\n301046, Big Data, MICT\n300581, Programming Techniques, BICT\n300144, Object Oriented Analysis, BICT\n300103, Data Structure, BCS\n300147, Object Oriented Programming, BCS\n300569, Computer Security, BIS\n301044, Data Science, MICT\n300582, Technologies for Web Applications, BICT')

Observe that:

1. The `"""` are necessary for new line strings
2. The `open(file, w)` will write over any pre-existing file (like `>` in bash)
(a) using `open(file, a)` would append to a file (like `»` in bash)
3. Nothing is written to disk until after `f.close()`, that's when the changes go from memory to disk.

Parse the Text File

Read the Text File

```
1  ## Open the File
2  scemunits_fid = open('./scemunits.txt')
3
4  ## Dispense with the first line
5  header = scemunits_fid.readline()
6
7  ## Read the remaining Lines into a var
8  scemunits_txt = scemunits_fid.read()
9
10 ## Print what we have
11 print(scemunits_txt)
12
13 ## Close the file
14 scemunits_fid.close()
```

301046, Big Data, MICT
300581, Programming Techniques, BICT
300144, Object Oriented Analysis, BICT
300103, Data Structure, BCS
300147, Object Oriented Programming, BCS
300569, Computer Security, BIS
301044, Data Science, MICT
300582, Technologies for Web Applications, BICT

Return only Matching Data

```
1  ## Split each line into a list element
2  obs = scemunits_txt.split('\n')
3
4  ## Throw away the empty line
5  obs = list(filter(None, obs))
6
7  ## Get the Course Names
8      ## Use replace so whitespace is not required after ,
9  courses = [ obs[i].replace(', ', ',').split(',')[2] for i in
    ↪ range(len(obs)) ]
10 units = [ obs[i].replace(', ', ',').split(',')[1] for i in
    ↪ range(len(obs)) ]
11
12 ## Enumerate the obs so that they
13 obs = list(obs)
14
15 ## Make an empty list for the matches
16 matches = []
17
18 ##
19 for i in range(len(obs)):
20     ## Don't Require whitespace after comma
21     if courses[i] == "MICT":
22         matches.append(obs[i])
23
24
25 #print([header] + join(matches).insert(header))
26 #print([header].append(matches))
27 print(matches)
28 matches.insert(0, header.replace('\n', ''))
29 print("\n".join(matches))
30
31 out_fid = open('outfile.txt', "w")
32 # out_fid.write("\n".join(matches))
33 for i in matches:
34     out_fid.write(i+'\n')
35     print(i)
36
37 out_fid.close()
```

```
['301046, Big Data, MICT', '301044, Data Science, MICT']
```

Unit ID, unit name, course name

301046, Big Data, MICT

301044, Data Science, MICT

```
Unit ID, unit name, course name
301046, Big Data, MICT
301044, Data Science, MICT
```

We can now inspect the contents of that file:

Wrap it into a function

```
1  #!/usr/bin/env python3
2
3  # * Create the Text File
4  scemunits = """Unit ID, unit name, course name
5  301046, Big Data, MICT
6  300581, Programming Techniques, BICT
7  300144, Object Oriented Analysis, BICT
8  300103, Data Structure, BCS
9  300147, Object Oriented Programming, BCS
10 300569, Computer Security, BIS
11 301044, Data Science, MICT
12 300582, Technologies for Web Applications, BICT"""
13
14 def writeTextFile(text, filename) :
15     f = open(filename, 'w')
16     for i in text.split('\n'):
17         f.writelines(i+'\n')
18     f.close()
19
20 writeTextFile(scemunits, 'scemunits.txt')
21
22
23 # * Main Functions
24 def readWriteFile(infile, outfile):
25     readTheTextFile(infile)
26     listOfLines = returnMatchingData(outfile)
27     print(listOfLines)
28     writeToFile(listOfLines, outfile)
29
30
31 # ** Sub Functions
32 # *** Input
33 def readTheTextFile(infile):
34     ## Open the File
35     scemunits_fid = open(infile)
36     ## ////////////////////////////////// File Open //////////////////////////////////
37
38     ## Dispense with the first line
39     readTheTextFile.header = scemunits_fid.readline().replace('\n', '')
40
41     ## Read the remaining lines into an attribute
42     readTheTextFile.scemunits_txt = scemunits_fid.read()
43
44     ## Close the File
45     ## ////////////////////////////////// File Closed //////////////////////////////////
46     scemunits_fid.close()
47
48
49 # ** Output
```

```
['Unit ID, unit name, course name', '301046, Big Data, MICT', '301044, Data Science, MICT']  
Unit ID, unit name, course name  
301046, Big Data, MICT  
301044, Data Science, MICT
```

and to confirm that it has written to the file:

```
1 cat outfile.txt
```

Parse Dictionary

```
1  #!/usr/bin/env python3
2
3  # * Create the Dictionary
4
5  units = {('301046', 'Big Data'): 'MICT',
6           ('300581', 'Programming Techniques'): 'BICT',
7           ('300144', 'OOA'): 'BICT',
8           ('300103', 'Data Structures'): 'BCS',
9           ('300147', 'OOP'): 'BCS',
10          ('300569', 'Computer Security'): 'BIS',
11          ('301044', 'Data Science'): 'MICT',
12          ('300582', 'TWA'): 'BICT'}
13
14
15  def displayUnits(unitsDict, keyword):
16      # Should Return Gracefully if the input is wrong
17      # Could have used Try/Except
18      if type(unitsDict) != dict:
19          print("ERROR; Require Dictionary of Unit Values")
20          return
21      # Make an empty List to fill
22      matches = []
23      # For each dictionary item if it corresponds to the keyword
24      # append it to the list
25      for i in unitsDict:
26          if units[i] == keyword:
27              matches.append(i)
28      # Use to get back matches[][1]
29      matching_units = [matches[i][1] for i in range(len(matches))]
30      # Return the Value
31      return matching_units
32
33
34  # To Print the Values join the list together with new line characters.
35  # The function should return data in a list not a string
36  # (python => data, bash => string)
37
38  print("Match MICT \n -----")
39  print("\n".join(displayUnits(units, 'MICT')))
40
41  print("Match BCS \n -----")
42  print("\n".join(displayUnits(units, 'BCS')))
```

Big Data
Data Science
Data Structures
OOP

grep

This is easy, just loop through the lines and print if the word is in the line.

```
1  #!/bin/python
2
3  def pygrep(filename, expr):
4      try:
5          inputfile_fid = open(filename)
6      except:
7          print("ERROR: Could not open file")
8      for line in inputfile_fid:
9          if expr in line:
10             print(line)
11
12  pygrep("./bigdata.txt", 'Big')
13  pygrep("./bigdata.txt", 'technology')
```

Big data is a broad term for data sets so large or complex that they are difficult to process with their tools, and expanding capabilities make Big Data a moving target. Thus, what is considered Big data usually includes data sets with sizes beyond the ability of commonly used software tools to capture, manage, and process data within a tolerable elapsed time. Big data "size" is a constantly moving target, as of 2012 ranging from a few dozen terabytes to many petabytes of data. Big data is a definition as follows: "Big data is high volume, high velocity, and/or high variety information." Big data uses inductive statistics and concepts from nonlinear system identification to infer patterns. Big data can also be defined as "Big data is a large volume unstructured data which cannot be handled by traditional data processing applications." Big data can be described by the following characteristics: Volume The quantity of data that is under consideration and whether it can actually be considered as Big Data or not. The name is Volume. Variety - The next aspect of Big Data is its variety. This means that the category to which Big Data belongs is upholding the importance of the Big Data.

complexity of Big Data.

Big data analytics consists of 6 Cs in the integrated industry 4.0 and Cyber Physical Systems (CPS) and its implications in an article titled "Big Data Solution Offering". The methodology addresses handling of Big Data.

Big Data Analytics for Manufacturing Applications can be based on a 5C architecture (connection, computation, communication, context, and control).

Big Data Lake - With the changing face of business and IT sector, capturing and storage of data is becoming a challenge.

Big data requires exceptional technologies to efficiently process large quantities of data with high accuracy.

Bus wrapped with SAP Big data parked outside IDF13.

Big data has increased the demand of information management specialists in that Software AG, Oracle, SAP, etc.

While many vendors offer off-the-shelf solutions for Big Data, experts recommend the development of custom solutions.

The use and adoption of Big Data, within governmental processes, is beneficial and allows efficient decision making.

Governmental Big Data space.

In 2012, the Obama administration announced the Big Data Research and Development Initiative, which aims to advance the use of big data in government.

different big data programs spread across six departments. Big data analysis played a large role in the success of the initiative.

Big data analysis was, in parts, responsible for the BJP and its allies to win a highly successful election in 2014.

benefit of big data for manufacturing. Big data provides an infrastructure for transparency in manufacturing.

In order to hone into the manner in which the media utilises Big Data, it is first necessary to understand the media landscape.

Practitioners in Advertising and Media approach Big Data as many actionable points of information.

The media industries process Big Data in a dual, interconnected manner:

Big Data and the IoT work in conjunction. From a media perspective, Data is the key derivative of the IoT.

far-reaching impacts on media efficiency. The wealth of data generated by this industry (i.e. the IoT) is vast.

Engineering Education. Gautam Siwach engaged at Tackling the challenges of Big Data by MIT Computer Science and Technology.

In March 2012, The White House announced a national "Big Data Initiative" that consisted of six pillars.

The White House Big Data Initiative also included a commitment by the Department of Energy to promote the use of big data.

The U.S. state of Massachusetts announced the Massachusetts Big Data Initiative in May 2012, which aims to advance the use of big data in government.

Massachusetts Institute of Technology hosts the Intel Science and Technology Center for Big Data. The European Commission is funding the 2-year-long Big Data Public Private Forum through their presenters from various industrial companies discussed their concerns, issues and future goals. Computational social sciences Anyone can use Application Programming Interfaces (APIs) provided the emergence of the typical network characteristics of Big Data". In their critique, Snijders, Marston. Big data has been called a "fad" in scientific research and its use was even made fun of as an "illusion". Questions for Big Data", the authors title big data a part of mythology: "large data sets offer a seductive but often "lost in the sheer volume of numbers", and "working with Big Data is still subjective, and often such as pro-active reporting especially target improvements in usability of Big Data, through a Big data analysis is often shallow compared to analysis of smaller data sets. In many big data cases Big data is a buzzword and a "vague term", but at the same time an "obsession" with entrepreneurs, consultants, scientists and the media. Big data showcases such as Google Flu Trends failed to correctly predict election predictions solely based on Twitter were more often off than on target. Big data often went public with the launch of a company called Ayasdi.

ICT4D) suggests that big data technology can make important contributions but also present unique challenges in the later utilization stage. Finally, with ubiquitous connectivity offered by cloud computing technology as well as queries from more than half a million third-party sellers. The core technology that keeps the data in form at cloud interface by providing the raw definitions and real time examples within the technology experiments (i.e. process a big amount of scientific data; although not with big data technology).

Top 10 Words

ope

So the idea here is to first try and open the file, use try/catch so that errors relating to missing files are descriptive and exit gracefully.

- If the word is in the list but not the dictionary
 - Put it in the dictionary with a value of 1
- If the word is in the dictionary then increment its value

```
{'unit': 2, 'id.': 1, 'name.': 1, 'course': 1, 'name\n301046.': 1, 'big': 1, 'data.': 1, '}
```

17

```
1  from operator import itemgetter
2  myDict = createDict('bigdata.txt')
3
4  sortedList = sorted(myDict.items(), key = itemgetter(1), reverse = True)
5
6  for key, value in sortedList[:10]:
7      print(key, value)
```

the 302
of 213
and 201
data 178
to 170
in 122
a 99
big 80
is 75
as 59