

TextProcessing-Solutions

February 21, 2015

1 Workign with text

One of Python's strengths is the ease of working with text. Here are some examples.

1.1 String methods

```
In [1]: # multi-line strings use triple quotes
```

```
s = """
it was the best of times,
it was the worst of times,
it was the age of wisdom,
it was the age of foolishness,
it was the epoch of belief,
it was the epoch of incredulity,
it was the season of Light,
it was the season of Darkness,
it was the spring of hope,
it was the winter of despair,
"""

print s.count('of')
print s.find('wisdom')
print s.find('foolsihness')
```

```
10
72
-1
```

```
In [2]: print s.upper()
```

```
IT WAS THE BEST OF TIMES,
IT WAS THE WORST OF TIMES,
IT WAS THE AGE OF WISDOM,
IT WAS THE AGE OF FOOLISHNESS,
IT WAS THE EPOCH OF BELIEF,
IT WAS THE EPOCH OF INCREULITY,
IT WAS THE SEASON OF LIGHT,
IT WAS THE SEASON OF DARKNESS,
IT WAS THE SPRING OF HOPE,
IT WAS THE WINTER OF DESPAIR,
```

```
In [3]: print s.replace('was', 'might have been')
```

```

it might have been the best of times,
it might have been the worst of times,
it might have been the age of wisdom,
it might have been the age of foolishness,
it might have been the epoch of belief,
it might have been the epoch of incredulity,
it might have been the season of Light,
it might have been the season of Darkness,
it might have been the spring of hope,
it might have been the winter of despair,

```

1.2 Splitting and joining strings

```

In [4]: paths = ['echo $PATH'
               print paths[0]

```

```

/bin/sh: /usr/local/bin:/Users/cliburn/git/julia:/Developer/NVIDIA/CUDA-6.5/bin:/Users/cliburn/anaconda

```

```

In [5]: for path in paths[0].split(':'):
         print '=> '.join(path.strip().split('/'))

```

```

=> bin=> sh
=> usr=> local=> bin
=> Users=> cliburn=> git=> julia=>
=> Developer=> NVIDIA=> CUDA-6.5=> bin
=> Users=> cliburn=> anaconda=> bin
=> usr=> bin
=> bin
=> usr=> sbin
=> sbin
=> usr=> local=> bin
=> opt=> X11=> bin
=> usr=> texbin
No such file or directory

```

2 The string module

The string module provides a very useful maketrans function. It is easier to show than to explain what this does.

```

In [6]: from string import maketrans

         dna_to_rna = maketrans('ACTG', 'ACUG')

         dna = 'gattaca'
         print dna.upper().translate(dna_to_rna).lower()

```

```

gauuaca

```

```

In [7]: # Incidentally the translate function is useful for getting rid of unwanted characters in a string

```

```

         from string import punctuation
         print punctuation

!"#$%&'()*+,-./:;<=>?@[\\]^_`{|}~

```

```
In [8]: import os

        # Alice in Wonderland from Project Gutenberg

        if not os.path.exists('alice.txt'):
            ! wget http://www.gutenberg.org/cache/epub/11/pg11.txt -O alice.txt

In [9]: from collections import Counter

        # Remove
        alice = open('alice.txt').read()
        words = alice.translate(None, punctuation).lower().split()
        word_counts = Counter(words)
        for item in word_counts.most_common(10):
            print item
        print 'alice', word_counts['alice']

('the', 1804)
('and', 912)
('to', 801)
('a', 684)
('of', 625)
('it', 541)
('she', 538)
('said', 462)
('you', 429)
('in', 428)
alice 385
```

2.1 Regular expressions

Regular expressions are a domain specific language for flexible text processing. It is a useful tool, but can be hard to decipher unless you use it often. Where possible, use string methods in preference to regular expressions. Sometimes, however, regular expressions are extremely useful. We will illustrate its use for motif finding in DNA sequences.

See [Regular Expression HOWTO](#) and the [re documentation](#) for details.

```
In [10]: # Here is the E Coli DNA sequence for the beta-D-galactosidase enzyme.

gene = """
>ENA|BAE76126|BAE76126.1 Escherichia coli str. K-12 substr. W3110 beta-D-galactosidase
ATGACCATGATTACGGATTCACTGGCCGTCGTTTTACAACGTCGTGACTGGGAAAACCT
GGCGTTACCCAACTTAATCGCCTTGCAGCACATCCCCCTTTCGCCAGCTGGCGTAATAGC
GAAGAGGCCCGACCGATCGCCCTTCCCAACAGTTGCGCAGCCTGAATGGCGAATGGCGC
TTTGCCTGGTTTCCGGCACCAGAAGCGGTGCCGAAAAGCTGGCTGGAGTGCATCTTCT
GAGGCCGATACTGTCTGCTGCCCTCAAAGTGGCAGATGCACGGTTACGATGCGCCATC
TACACCAACGTGACCTATCCCATACGGTCAATCCGCCGTTTGTTCACGGAGAATCCG
ACGGGTTGTTACTCGCTCACATTTAATGTTGATGAAAGCTGGCTACAGGAAGGCCAGACG
CGAATTATTTTGTATGGCGTTAACTCGGCGTTTCTATCTGTGGTGCAACGGCGCTGGGTC
GGTTACGGCCAGGACAGTCGTTTGCCGCTCTGAATTTGACCTGAGCGCATTTTACGCGCC
GGAGAAAACCGCCTCGCGGTGATGGTGCTGCGCTGGAGTGACGGCAGTTATCTGGAAGAT
CAGGATATGTGGCGGATGAGCGGCATTTTCCGTGACGTCTCGTTGCTGCATAAACCGACT
ACACAAATCAGCGATTTCCATGTTGCCACTCGCTTTAATGATGATTTACGCCGCGCTGTA
CTGGAGGCTGAAGTTCAGATGTGCGGCGAGTTGCGTGACTACCTACGGGTAACAGTTTCT
TTATGGCAGGTTGAAACGCAGGTCGCCAGCGGCACCGCGCCTTTCGGCGGTGAAATTATC
```

```

GATGAGCGTGGTGGTTATGCCGATCGCGTCACACTACGTCTGAACGTCGAAAACCCGAAA
CTGTGGAGCGCCGAAATCCCGAATCTCTATCGTGCGGTGGTTGAACTGCACACGCGCGAC
GGCAGCGTGATTGAAGCAGAAGCCTGCGATGTCGGTTTCCGCGAGGTGCGGATTGAAAAT
GGTCTGCTGCTGCTGAACGGCAAGCCGTTGCTGATTGAGGCGTTAACCGTCACGAGCAT
CATCTCTGCATGGTCAGGTCATGGATGAGCAGACGATGGTGCAGGATATCCTGCTGATG
AAGCAGAACAACCTTTAACGCCGTGCGCTGTTTCGATTATCCGAACCATCCGCTGTGGTAC
ACGCTGTGCGACCGCTACGGCCTGTATGTGGTGGATGAAGCCAATATTGAAACCCACGGC
ATGGTGCCAATGAATCGTCTGACCGATGATCCGCGCTGGCTACCGGCGATGAGCGAACGC
GTAACGCGAATGGTGCAGCGCGATCGTAATCACCCGAGTGTGATCATCTGGTCGCTGGGG
AATGAATCAGGCCACGGCGCTAATCAGCAGCGCTGTATCGCTGGATCAAATCTGTGAT
CCTTCCCGCCCGGTGTCAGTATGAAGGCGGCGAGCCGACACCACGGCCACCGATATTATT
TGCCCGATGTACGCGCGGTGGATGAAGACCAGCCCTTCCCGGTGTGCCGAAATGGTCC
ATCAAAAAATGGCTTTTCGCTACCTGGAGAGACGCGCCCGCTGATCCTTTGCGAATACGCC
CACGCGATGGGTAACAGTCTTGGCGTTTCGCTAAATACTGGCAGGCGTTTCGTCAGTAT
CCCCGTTTACAGGGCGGCTTCGTCTGGGACTGGGTGGATCAGTCGCTGATTAAATATGAT
GAAAACGGCAACCCGTGGTCGGCTTACGGCGGTGATTTTGGCGATACGCCGAACGATCGC
CAGTTCTGTATGAACGGTCTGGTCTTTGCCGACCGCACGCGCATCCAGCGCTGACGGAA
GCAAAACACCAGCAGCAGTTTTTCCAGTTCGTTTATCCGGGCAAACCATCGAAGTGACC
AGCGAATACCTGTTCCGTCATAGCGATAACGAGCTCCTGCACTGGATGGTGGCGCTGGAT
GGTAAGCCGTGGCAAGCGGTGAAGTGCCTCTGGATGTGCTCCACAAGGTAACAGTTG
ATTGAACTGCCTGAACTACCGCAGCCGAGAGCGCCGGGCAACTCTGGCTCACAGTACGC
GTAGTGCAACCGAACGCGACCGCATGGTCAGAAGCCGGGCACATCAGCGCTGGCAGCAG
TGGCGTCTGGCGGAAAACCTCAGTGTGACGCTCCCGCGCGTCCACGCCATCCCGCAT
CTGACCACCAGCGAAATGGATTTTTGCATCGAGCTGGGTAATAAGCGTTGGCAATTTAAC
CGCCAGTCAGGCTTTCTTTCACAGATGTGGATTGGCGATAAAAAACAATGCTGACGCCG
CTGCGCGATCAGTTACCCGTGCACCGCTGGATAACGACATTGGCGTAAGTGAAGCGACC
CGCATTGACCCTAACGCCTGGGTGCAACGCTGGAAGGCGCGGGCCATTACCAGGCCGAA
GCAGCGTTGTTGCAGTGCACGGCAGATACACTTGCTGATGCGGTGCTGATTACGACCGCT
CACGCGTGGCAGCATCAGGGGAAAACCTTATTTATCAGCCGAAAACCTACCGGATTGAT
GGTAGTGGTCAAATGGCGATTACCGTTGATGTTGAAGTGGCGAGCGATACACCGCATCCG
GCGCGGATTGGCCTGAACTGCCAGCTGGCGCAGGTAGCAGAGCGGGTAACTGGCTCGGA
TTAGGGCCGCAAGAAAACCTATCCCGACCGCCTTACTGCCGCTGTTTTGACCGCTGGGAT
CTGCCATTGTCAGACATGTATAACCCGTACGTCTTCCCGAGCGAAAACGGTCTGCGCTGC
GGGACGCGCGAATTGAATTATGGCCACACCAAGTGGCGCGGCGACTTCCAGTTCAACATC
AGCCGCTACAGTCAACAGCAACTGATGGAACACAGCCATCGCCATCTGCTGCACGCGGAA
GAAGGCACATGGCTGAATATCGACGTTTCCATATGGGGATTGGTGGCGACGACTCCTGG
AGCCCGTCAGTATCGGCGGAATCCAGCTGAGCGCCGGTTCGCTACCATTACCAGTTGGTC
TGGTGTCAAAAATAA
"""

```

```

In [11]: # Suppose we want to replace motifs that start with 'ATA',
# followed by between 1 and 4 of any nucleotide, followed by 'CG'
# with a blank string of the same length

```

```

import re
from toolz import partition

def replace(match):
    return ' ' * len(match.group(0))

# convert FASTA into single DNA sequence
dna = ''.join(line for line in gene.strip().split('\n'))
if not line.startswith('>')
pattern = 'ATA.{1,4}CG'

```

```

modified_dna = re.sub(pattern, replace, dna)

# pretty print modified sequence
linewidth = 60
print '\n'.join([''.join(line) for line
                  in partition(linewidth, modified_dna)])

ATGACCATGATTACGGATTCACTGGCCGTCGTTTTACAACGTCGTGACTGGGAAAACCCCT
GGCGTTACCCAACCTTAATCGCCTTGACGACATCCCCCTTTTCGCCAGCTGGCGTA
AAGAGGCCCGACCGATCGCCCTTCCCAACAGTTGCGCAGCCTGAATGGCGAATGGCGC
TTTGCTGTTTTCCGGCACCAGAAGCGGTGCCGAAAGCTGGCTGGAGTGCATCTTCCT
GAGGCCG          TCGTCCCCTCAAACCTGGCAGATGCACGGTTACGATGCGCCCATC
TACACCAACGTGACCTATCCATTACGGTCAATCCGCCGTTTGTTCACGGAAGAATCCG
ACGGGTTGTTACTCGCTCACATTTAATGTTGATGAAAGCTGGCTACAGGAAGGCCAGACG
CGAATTATTTTTGATGGCGTTAACTCGGCGTTTCATCTGTGGTGCAACGGGCGCTGGGTC
GGTTACGGCCAGGACAGTCGTTTGCCGTCTGAATTTGACCTGAGCGCATTTTTACGCGCC
GGAGAAAACCGCCTCGCGTGATGGTGCTGCGCTGGAGTGACGGCAGTTATCTGGAAGAT
CAGGATATGTGGCGGATGAGCGGCATTTCCGTGACGTCTCGTTGCTGC          ACT
ACACAAATCAGCGATTTCCATGTTGCCACTCGCTTTAATGATGATTTACGCCGCGCTGTA
CTGGAGGCTGAAGTTCAGATGTGCGCGAGTTGCGTGACTACCTACGGGTAACAGTTTCT
TTATGGCAGGGTGAAACGCAGGTCGCCAGCGGCACCGCGCTTTCGGCGGTGAAATTATC
GATGAGCGTGGTGGTTATGCCGATCGCGTCACACTACGTCTGAACGTCGAAAACCCGAAA
CTGTGGAGCGCCGAAATCCCGAATCTCTATCGTGCGGTGGTTGAACTGCACACCGCCGAC
GGCAGCGTGATTGAAGCAGAAGCCTGCGATGTCGTTTCCGCGAGGTGCGGATTGAAAAT
GGTCTGCTGCTGCTGAACGCAAGCCGTTGCTGATTTCGAGGCGTTAACCGTCACGAGCAT
CATCTCTGCATGGTCAGGTCATGGATGAGCAGACGATGGTGCAAGGATATCCTGCTGATG
AAGCAGAACAACCTTTAACGCCGTGCGCTGTTTCGCATTATCCGAACCATCCGCTGTGGTAC
ACGCTGTGCGACCGCTACGGCCTGTATGTGGTGGATGAAGCCAATATTGAAACCCACGGC
ATGGTGCCAATGAATCGTCTGACCGATGATCCGCGCTGGCTACCGGCGATGAGCGAACGC
GTAACGCGAATGGTGACGCGCATCGTAATCACCCGAGTGATCATCTGGTCGCTGGGG
AATGAATCAGGCCACGGCGCTAATCACGACGCGCTGTATCGCTGGATCAAATCTGTGAT
CCTTCCCGCCCGGTGACGATATGAAGGCGCGGAGCCGACACCACGGCCACCGATATTATT
TGCCCGATGTACGCGCGCTGGATGAAGACCAGCCCTTCCCGGCTGTGCCGAAATGGTCC
ATCAAAAAATGGCTTTCGCTACCTGGAGAGACGCGCCCGCTGATCCTTTGCGAATACGCC
CACGCGATGGGTAACAGTCTTGCGGCTTTCGCTAAATACTGGGAGGCGTTTCGTCAGTAT
CCCCGTTTACAGGGCGGCTTCGTCTGGGACTGGGTGGATCAGTCGCTGATTAAATATGAT
GAAAACGGCAACCCGTGGTTCGGCTTACGGCGGTGATTTTGGCG          AACGATCGC
CAGTTCTGTATGAACGCTCTGGTCTTTGCCGACCGCACGCCGATCCAGCGCTGACGGAA
GCAAAACACCAGCAGCAGTTTTTCCAGTTCCGTTTATCCGGGCAAACCATCGAAGTGACC
AGCGAATACCTGTTCCGTC          AGCTCCTGCACTGGATGGTGGCGCTGGAT
GGTAAGCCGCTGGCAAGCGGTGAAGTGCTCTGGATGTGCTCCACAAGGTAAACAGTTG
ATTGAACTGCCTGAACTACCGCAGCCGAGAGCGCCGGGCAACTCTGGCTCACAGTACGC
GTAGTGCAACCGAACCGCAGCCGATGGTCAGAAGCCGGGCACATCAGCGCTGGCAGCAG
TGGCGTCTGGCGGAAAACCTCAGTGTGACGCTCCCGCGCGCTCCACGCCATCCCGCAT
CTGACCACAGCGAAATGGATTTTTGCATCGAGCTGGGTA          TTGGCAATTTAAC
CGCCAGTCAGGCTTTCTTTACAGATGTGGATTGGCGATAAAAAACAACCTGCTGACGCCG
CTGCGCGATCAGTTCACCCGTGCACCGCTGG          ACATTGGCGTAAGTGAAGCGACC
CGCATTGACCCTAACGCCTGGGTGCAACGCTGGAAGGCGGCGGGCCATTACCAGGCCGAA
GCAGCGTTGTTGCAGTGCACGGCAGATACACTTGTGATGCGGTGCTGATTACGACCGCT
CACGCGTGGCAGCATCAGGGGAAAACCTTATTTATCAGCCGGAACCTACCGGATTGAT
GGTAGTGGTCAAATGGCGATTACCGTTGATGTTGAAGTGGCGAGCG          CATCCG
GCGCGGATTGGCCTGAACTGCCAGCTGGCGCAGGTAGCAGAGCGGGTAAACTGGCTCGGA
TTAGGGCCGCAAGAAAACCTATCCCGACCGCCTTACTGCCGCTGTTTTGACCGCTGGGAT
CTGCCATTGTCAGACATGT          TACGTCTTCCCGAGCGAAAACGGTCTGCGCTGC

```

```
GGGACGCGCGAATTGAATTATGGCCACACACAGTGGCGCGGCGACTTCCAGTTCAACATC
AGCCGCTACAGTCAACAGCAACTGATGGAAACCAGCCATCGCCATCTGCTGCACGCGGAA
GAAGGCACATGGCTGA          GTTCCATATGGGGATTGGTGGCGACGACTCCTGG
AGCCCGTCAGTATCGGCGGAATTCCAGCTGAGCGCCGGTCGCTACCATTACCAGTTGGTC
```

2.2 The NLTK toolkit

If you will be doing statistical natural language processing or significant amounts of machine learning on natural text, check out the [Natural Language Toolkit](#).

2.3 Exercises

1. Write a function to find the complementary strand given a DNA sequence. For example
Given ATCGTTA Return TAGCAAT
Note: The following are complementary bases A|T, C|G.

In [12]: # YOUR CODE HERE

```
def complement(dna):
    """Return complementary strand given DNA sequence."""
    import string
    table = string.maketrans('actgACTG', 'tgacTGAC')
    return dna.translate(table)

print complement('ATCGTTA')
```

TAGCAAT

2. Write a regular expression that matches the following:

- Phone numbers with the format: (919)-1234567 (i.e. (123)-9876543 should match but not 234-1234567 or (123)-666666)
- Email addresses john.doe@duke.edu (i.e. steve@gmail.com should match but not steve@gmail)
- DNA sequences with the motif A-C-T-G where - indicates 0 or 1 other nucleotide (any of A,C,T or G)

In [13]: # YOUR CODE HERE

```
phone_pat = re.compile(r'\(\d{3}\)-\d{7}')

for s in ['(123)-9876543', '234-1234567', '(123)-666666']:
    m = phone_pat.match(s)
    if m:
        print 'Matched', s
    else:
        print 'Not matched', s
```

Matched (123)-9876543
Not matched 234-1234567
Not matched 123)-666666)

Note: This is just for practice - actual email validators should not be using regular expressions because the rules for a valid email are insanely [complex](#), and should probably be checked with a *parser*.

In [14]: email_pat = re.compile(r'[\w]+[\.\[\w]+\]?@([\w]+\.)+[\w]+')

```
for s in ['john@', 'john.doe@duke.edu', 'steve@gmail.com', 'steve@gmail']:
```

```

m = email_pat.match(s)
if m:
    print 'Mathced', s
else:
    print 'Not matched', s

```

```

Not matched johm@
Mathced john.doe@duke.edu
Mathced steve@gmail.com
Not matched steve@gmail

```

In [15]: `motif_pat = re.compile(r'A.?C.?T.?G')`

```

for s in ['GATTACA', 'ACTG', 'AACCTTGG', 'AAACCCTTTGGG']:
    m = motif_pat.match(s)
    if m:
        print 'Mathced', s
    else:
        print 'Not matched', s

```

```

Not matched GATTACA
Mathced ACTG
Mathced AACCTTGG
Not matched AAACCCTTTGGG

```

3. Download ‘Pride and Prejudice’ by Jane Austen from Project Gutenbrg.

- Remove all punctuation and covert to lower case
- Count how many times the word ‘married’ appears
- Count how often the word ‘daughter’ and ‘married’ appear in the same 10-word window

In [16]: *# YOUR CODE HERE*

```

if not os.path.exists('pride_and_prejudice.txt'):
    ! curl 'http://www.gutenberg.org/cache/epub/1342/pg1342.txt' > 'pride_and_prejudice.txt'

```

In [17]: `import string`

```

with open('pride_and_prejudice.txt') as f:
    s = f.read()
    s = s.lower().translate(None, string.punctuation)

    words = s.split()
    size = 10
    windows = list(partition(size, words))
    print "'daughter' and 'married' appera %d times in the same 10-word window" % \
        sum('daughter' in window and 'married' in window for window in windows)
    print "The word 'married' appears %d times" % s.count('married')

```

```

'daughter' and 'married' appera 5 times in the same 10-word window
The word 'married' appears 61 times

```

4. Download “The Gutenberg Webster’s Unabridged Dictionary” from Project Gutenbrg

- First extract all defined words (109561 words) - oops I cannot replicate this number
- Count the number of *defined* English words containing 3 or more vowels (aeiou)

- Find all longest palindrome (a palindrome is a word that is spelt the same forwards as backwards - e.g. 'deified')

In [18]: # YOUR CODE HERE

```
# If you look at the plain text file,
# it is quite hard to figure out how to extract a defined word.
# We have more luck with the HTML file.

if not os.path.exists('websters.html'):
    ! curl 'www.gutenberg.org/cache/epub/29765/pg29765.html' > 'websters.html'
```

In [19]: ! head -n 400 websters.html | tail -n 30

In [20]: # Notice that in the HTML, word definitions have the structure <p id="xxxxxx">WORD</br> or <p

```
text = open('websters.html').read()
word = re.compile(r'<p id="id\d+">([A-Z]+)[<br/>|\r\n+]')
```

words = word.findall(text)

count = 0

```
for word in words:
    if word.count('A') + word.count('E') + word.count('I') + word.count('O') + word.count('U')
        count += 1
```

```
print "Number of words is %d" % len(words)
print "Number of words with 3 or more vowels is %d" % count
```

palindromes = [word for word in words if word == word[::-1]]

lengths = map(len, palindromes)

max_len = max(lengths)

```
print "Longest palindromes are", [p for p in palindromes if len(p) == max_len]
```

Number of words is 103020

Number of words with 3 or more vowels is 69210

Longest palindromes are ['MALAYALAM']

In [20]: