

Chapter 19: Gematria: Numeric encoding of text using ASCII values

Gematria is a system for assigning a number to a word by summing the numeric values of each of the characters ^[1]. In the standard encoding (*Mispar hechrechi*), each character of the Hebrew alphabet is assigned a numeric value ranging from 1 to 400, but there are more than a dozen other methods for calculating the numeric value for the letters. To encode a word, these values are added together. Revelation 13:18 from the Christian Bible says "Let the one who has insight calculate the number of the wild beast, for it is a man's number, and its number is 666." Some scholars believe that number is derived from the encoding of the characters representing Nero Caesar's name and title and was used as a way of writing about the Roman emperor without naming him.



We will write a program called `gematria.py` that will numerically encode each word in a given text. There are many ways we could assign a numeric value to the characters of the English alphabet. For instance, we could start by giving "a" the value 1, "b" the value 2, and so forth. Instead, we can use the ASCII table ^[2] to derive a numeric for English alphabet characters. For non-English characters, we could consider using a Unicode value, but this exercise will stick to ASCII letters.

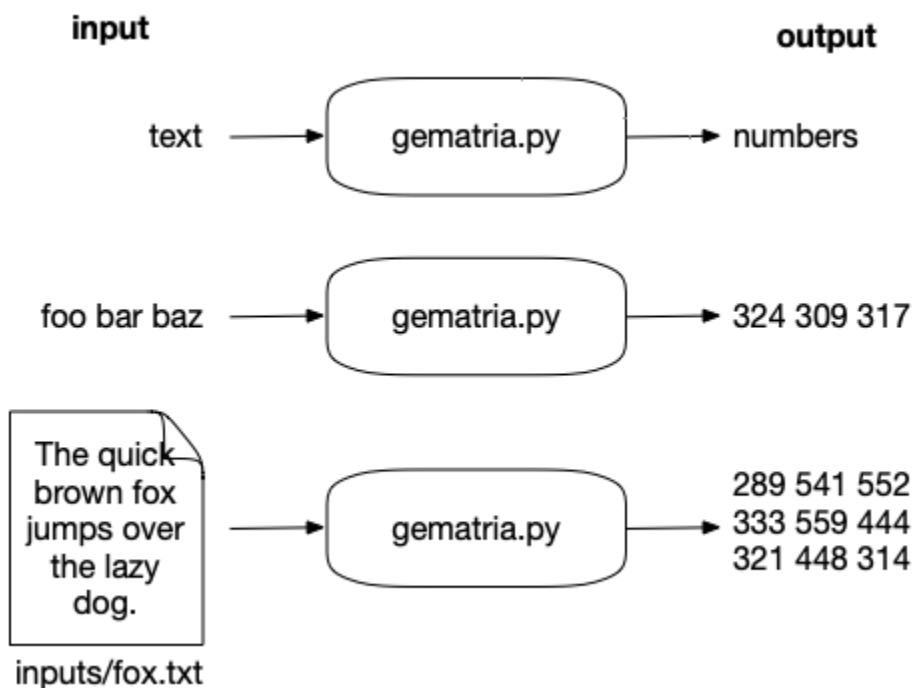
The input text may be given on the command line:

```
$ ./gematria.py 'foo bar baz'
324 309 317
```

Or in a file:

```
$ ./gematria.py ../inputs/fox.txt
289 541 552 333 559 444 321 448 314
```

Here is a string diagram showing how the program should work:



In this exercise, you will:

- Learn about the `ord` and `chr` functions
- Explore how characters are organized in the ASCII table
- Understand character ranges used in regular expression
- Use the `re.sub` function
- Learn how `map` can be written without `lambda`
- Use the `sum` function
- Learn how to perform case-insensitive string sorting

Writing `gematria.py`

I will always recommend you start your programs in some way that avoids having to type all the boilerplate text. Either copy `template/template.py` to `gematria/gematria.py` or use `new.py gematria.py` in the `gematria` directory to create a starting point. Modify the program until it prints the following usage if given no arguments or the `-h` or `--help` flag:

```
$ ./gematria.py
usage: gematria.py [-h] str
gematria.py: error: the following arguments are required: str
$ ./gematria.py -h
usage: gematria.py [-h] str
```

Gematria

positional arguments:

str Input text or file

optional arguments:

-h, --help show this help message and exit

As in previous exercises, the input may come from the command line or from a file. I suggest you copy the code you used in "Howler" to handle this, then modify your `main` to the following:

```
def main():
    args = get_args()
    print(args.text.rstrip())
```

Verify that your program will print text from the command line:

```
$ ./gematria.py 'Death smiles at us all, but all a man can do is smile back.'
Death smiles at us all, but all a man can do is smile back.
```

Or from a file:

```
$ ./gematria.py ../inputs/spiders.txt
Don't worry, spiders,
I keep house
casually.
```

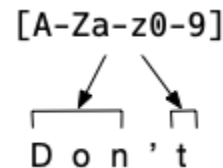
Cleaning a word

Let's discuss how a single word will be encoded as it will affect how we will break the text in the next section. In order to be absolutely sure we are only dealing with ASCII values, let's remove anything that is not an upper- or lowercase English alphabet character or any of the Arabic numerals 0-9. We can define that class of characters using the regular expression `[A-Za-z0-9]`. We can use the `re.findall` function we used in "Mad Libs" to find all the characters in `word` that match this class. In the word "Don't," we should expect to find everything except the apostrophe:

```
>>> re.findall('[A-Za-z0-9]', "Don't")
['D', 'o', 'n', 't']
```

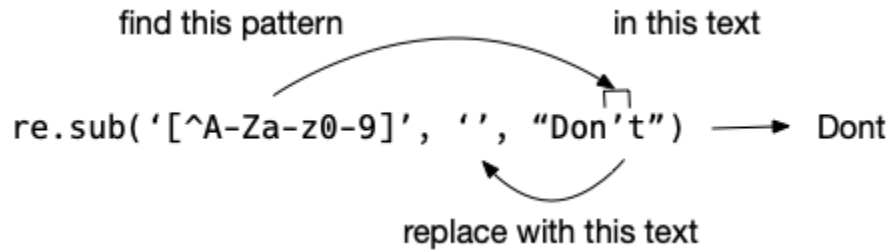
Anything *not* in that class can be defined by placing a caret (^ or "hat") as the first character inside the class like `[^A-Za-z0-9]`. Now we would expect to match *only* the apostrophe:

```
>>> re.findall('[^A-Za-z0-9]', "Don't")
['']
```



We can use the `re.sub` function to replace any characters in that second class with the empty string. As you learned in "Mad Libs," this will replace *all* occurrences of the pattern unless we use the `count=n` option:

```
>>> word =
re.sub('[^A-Za-z0-9]', '', "Don't")
>>> word
'Dont'
```



We will want use this operation to clean each word that we'll encode.

Ordinal character values and ranges

We will encode a string like "Dont" by converting *each character* to a numeric value and then adding them together, so let's first figure out how to encode a single character. Python has a function called `ord` that will convert a character to its "ordinal" value (its order in the ASCII table):

```
>>> ord('D')
68
>>> ord('o')
111
```

The `chr` function works in reverse to convert a number to a character:

```
>>> chr(68)
'D'
>>> chr(111)
'o'
```

Here is a table showing the lower-order ASCII values 0-127 ^[3]. For simplicity's sake, I show "NA" for "not available" for the ones up to index 31 as they are not printable:

\$./asciitbl.py							
0 NA	1 NA	2 NA	3 NA	4 NA	5 NA	6 NA	7 NA
8 NA	9 NA	10 NA	11 NA	12 NA	13 NA	14 NA	15 NA
16 NA	17 NA	18 NA	19 NA	20 NA	21 NA	22 NA	23 NA
24 NA	25 NA	26 NA	27 NA	28 NA	29 NA	30 NA	31 NA
32 SPACE	33 !	34 "	35 #	36 \$	37 %	38 &	39 '
40 (41)	42 *	43 +	44 ,	45 -	46 .	47 /
48 0	49 1	50 2	51 3	52 4	53 5	54 6	55 7
56 8	57 9	58 :	59 ;	60 <	61 =	62 >	63 ?
64 @	65 A	66 B	67 C	68 D	69 E	70 F	71 G
72 H	73 I	74 J	75 K	76 L	77 M	78 N	79 O
80 P	81 Q	82 R	83 S	84 T	85 U	86 V	87 W
88 X	89 Y	90 Z	91 [92 \	93]	94 ^	95 _
96 `	97 a	98 b	99 c	100 d	101 e	102 f	103 g
104 h	105 i	106 j	107 k	108 l	109 m	110 n	111 o
112 p	113 q	114 r	115 s	116 t	117 u	118 v	119 w
120 x	121 y	122 z	123 {	124	125 }	126 ~	127 DEL

We can use a `for` loop to cycle through all the characters in a string:

```
>>> word = "Dont"
>>> for char in word:
...     print(char, ord(char))
...
D 68
o 111
n 110
t 116
```

Note that upper- and lowercase letters have different `ord` values. It makes sense because they are two different letters:

```
>>> ord('D')
68
>>> ord('d')
100
```

We can iterate over the values from "a" to "z" by finding their `ord` values:

```
>>> [chr(n) for n in range(ord('a'), ord('z') + 1)]
['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm',
 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
```

The letters "a" through "z" lies contiguously in the ASCII table. You can do the same for "A" to "Z" and "0" to "9" which is why we can use `[A-Za-z0-9]` as a regex.

Note that the uppercase letters have lower ordinal values than their lowercase versions, which is why you cannot use the range `[a-Z]`. Try this in the REPL and note the error you get:

```
>>> re.findall('[a-Z]', word)
```

The last line I see is this:

```
re.error: bad character range a-Z at position 1
```

You *can* use the range `[A-z]`:

```
>>> re.findall('[A-z]', word)
['D', 'o', 'n', 't']
```

But see that that "Z" and "a" are not contiguous:

```
>>> ord('Z'), ord('a')
(90, 97)
```

There are other characters in between them:

```
>>> [chr(n) for n in range(ord('Z') + 1, ord('a'))]
['[', '\\', ']', '^', '_', ``']
```

If we try to use that range on all the printable characters, you'll see that it matches characters that are not letters:

```
>>> import string
>>> re.findall('[A-z]', string.printable)
['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm',
 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z',
 'A', 'B', 'C', 'D', 'E', 'F', 'G', 'H', 'I', 'J', 'K', 'L', 'M',
 'N', 'O', 'P', 'Q', 'R', 'S', 'T', 'U', 'V', 'W', 'X', 'Y', 'Z',
 '[', '\\', ']', '^', '_', ``']
```



That is why it is safest to specify the characters we want as the three ranges, [A-Za-z0-9], which you may sometimes hear pronounced as "a to z, A to Z, zero to nine" as it assumes you understand that there are two "a to z" ranges which are distinct according to their case.

Encoding a word

Let's keep reminding ourselves what the goal is here: convert all the characters in a word and then sum those values. There is a handy Python function called `sum` that will add a `list` of numbers:

```
>>> sum([1, 2, 3])
6
```

We can manually encode the string "Dont" by calling `ord` on each letter and passing the results as a `list` to `sum`:

```
>>> sum([ord('D'), ord('o'), ord('n'), ord('t')])
405
```

So the question is how to apply the function `ord` to all the elements of a `str` and give the result to `sum`. You've seen this pattern many times now. What's the first tool you'll reach for? We can always start with our handy `for` loop:

```
>>> word = 'Dont'
>>> vals = []
>>> for char in word:
...     vals.append(ord(char))
...
>>> vals
[68, 111, 110, 116]
```

Can you see how to make that into a single line using a list comprehension?

```
>>> vals = [ord(char) for char in word]
>>> vals
[68, 111, 110, 116]
```

And from there, we can move to a `map`:

```
>>> vals = map(lambda char: ord(char), word)
>>> list(vals)
[68, 111, 110, 116]
```

Here I'd like to show that the `map` version here doesn't even need the `lambda` declaration because the `ord` function expects a single value which is exactly what it will get from `map`. Here is a nicer way to write it:

```
>>> vals = map(ord, word)
>>> list(vals)
[68, 111, 110, 116]
```

To my eye, that is a really beautiful piece of code! Now we can `sum` that to get a final value for our `word`:

```
>>> sum(map(ord, word))
405
```

Which is correct:

```
>>> sum([68, 111, 110, 116])
405
```

We can create a function to encapsulate all this. I called mine `word2num`, and here is my test:

```
def test_word2num():
    """Test word2num"""

    assert word2num("a") == "97"
    assert word2num("abc") == "294"
    assert word2num("ab'c") == "294"
    assert word2num("4a-b'c,") == "346"
```

Notice that my function returns a `str` value, not an `int`. This is because I want to use the result with the `str.join` function that only accepts `str` values. So `'405'` instead of `405`:

```
>>> from gematria import word2num
>>> word2num("Don't")
'405'
```

To summarize, the `word2num` function accepts a word, removes unwanted characters, converts the remaining characters to `ord` values, and returns a `str` representation of the `sum` of those values.

Breaking the text

The tests expect you to maintain the same line breaks as the original text, so I recommend you use `str.splitlines` as in other exercises. In "Friar" and "Scrambler," we used different regexes to split each line into "words," a process sometimes called "tokenization" in programs that deal with Natural Language Processing (NLP). If you write a `word2num` that passes the tests I provide, then you can use `str.split` to break a line on spaces because the function will ignore anything that is not a character or number. You are, of course, welcome to break the line into words using whatever means you like.

This code will maintain the line breaks and reconstruct the text. Can you modify it to add the `word2num` function so that it instead prints out encoded words?

```
def main():
    args = get_args()
    for line in args.text.splitlines():
        for word in line.split():
            # what goes here?
            print(' '.join(line.split()))
```

The output will be one number for each word:

```
$ ./gematria.py ../inputs/fox.txt
289 541 552 333 559 444 321 448 314
```

The	quick	brown	fox	jumps	over	the	lazy	dog.
↓	↓	↓	↓	↓	↓	↓	↓	↓
289	541	552	333	559	444	321	448	314

Time to finish writing the solution. Be sure to use the tests! See you on the flip side.

Solution


```

1  #!/usr/bin/env python3
2  """Gematria"""
3
4  import argparse
5  import os
6  import re
7
8
9  # -----
10 def get_args():
11     """Get command-line arguments"""
12
13     parser = argparse.ArgumentParser(
14         description='Gematria',
15         formatter_class=argparse.ArgumentDefaultsHelpFormatter)
16
17     parser.add_argument('text', metavar='str', help='Input text or file') 1
18
19     args = parser.parse_args() 2
20
21     if os.path.isfile(args.text): 3
22         args.text = open(args.text).read() 4
23
24     return args 5
25
26
27 # -----
28 def main():
29     """Make a jazz noise here"""
30
31     args = get_args() 6
32
33     for line in args.text.splitlines(): 7
34         print(' '.join(map(word2num, line.split()))) 8
35
36
37 # -----
38 def word2num(word): 9
39     """Sum the ordinal values of all the characters"""
40
41     return str(sum(map(ord, re.sub('[^a-zA-Z0-9]', '', word)))) 10
42
43
44 # -----
45 def test_word2num(): 11
46     """Test word2num"""
47
48     assert word2num("a") == "97"
49     assert word2num("abc") == "294"
50     assert word2num("ab'c") == "294"
51     assert word2num("4a-b'c,") == "346"
52
53
54 # -----
55 if __name__ == '__main__':
56     main()

```

- 1 The `text` argument is a string which might be a file name.
- 2 Get the parsed command-line arguments.
- 3 Check if the `text` argument names an existing file.

- 4 If it does, overwrite the `args.text` with the contents of the file.
- 5 Return the fixed up arguments.
- 6 Get the parsed arguments.
- 7 Split `args.text` on newlines to retain line breaks.
- 8 Split the `line` on spaces, map the result through `word2num`, then join that result on spaces.
- 9 Define a function to convert a word to a number.
- 10 Use `re.sub` to remove anything not an alpha-numeric character, map the resulting string through the `ord` function, sum the ordinal values of the characters, and return a `str` representation of the number.
- 11 Define a function to test the `word2num` function.

Discussion

I trust you understand the `get_args` as we've used this exact code several times now. Let's jump to the `word2num` function.

Writing `word2num`

I could have written the function like this:

```
def word2num(word):
    vals = []
    for char in re.sub('[^a-zA-Z0-9]', '', word):
        vals.append(ord(char))

    return str(sum(vals))
```

- 1 Initialize an empty `list` to hold the values.
- 2 Iterate all the characters returned from `re.sub`.
- 3 Convert the character to an ordinal value and append that to the values.
- 4 Sum the values and return a string representation.

That's four lines of code instead of one I wrote. I would at least rather use a list comprehension which collapses three lines of code into one:

```
def word2num(word):
    vals = [ord(char) for char in re.sub('[^a-zA-Z0-9]', '', word)]
    return str(sum(vals))
```

Which can then be moved into one line:

```
def word2num(word):
    return str(sum([ord(char) for char in re.sub('[^a-zA-Z0-9]', '', word)]))
```

I still think the `map` version is the most readable and concise:

```
def word2num(word):
    return str(sum(map(ord, re.sub('[^a-zA-Z0-9]', '', word))))
```

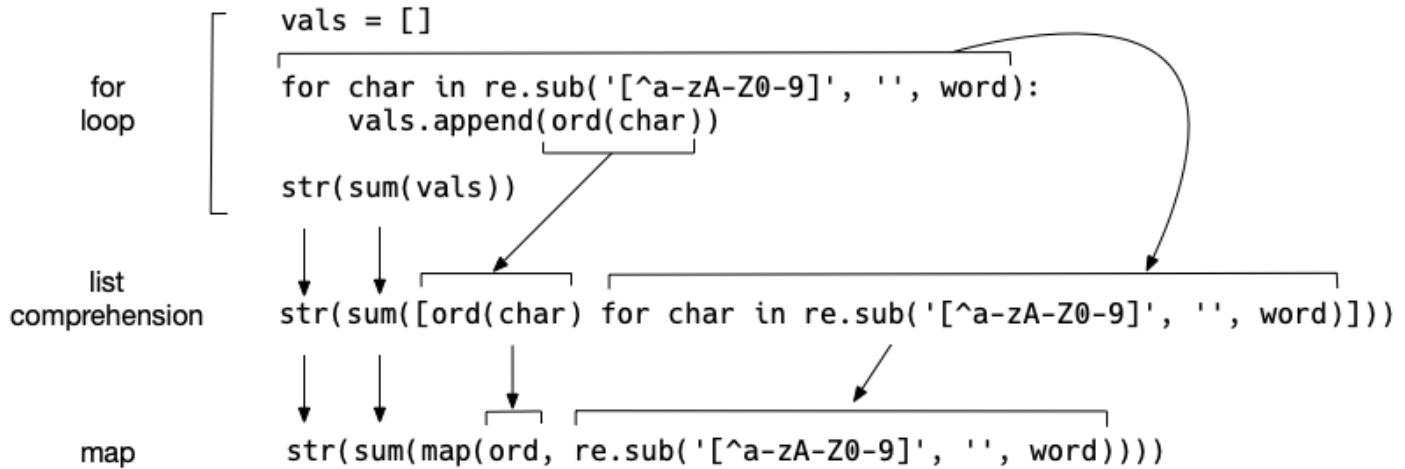
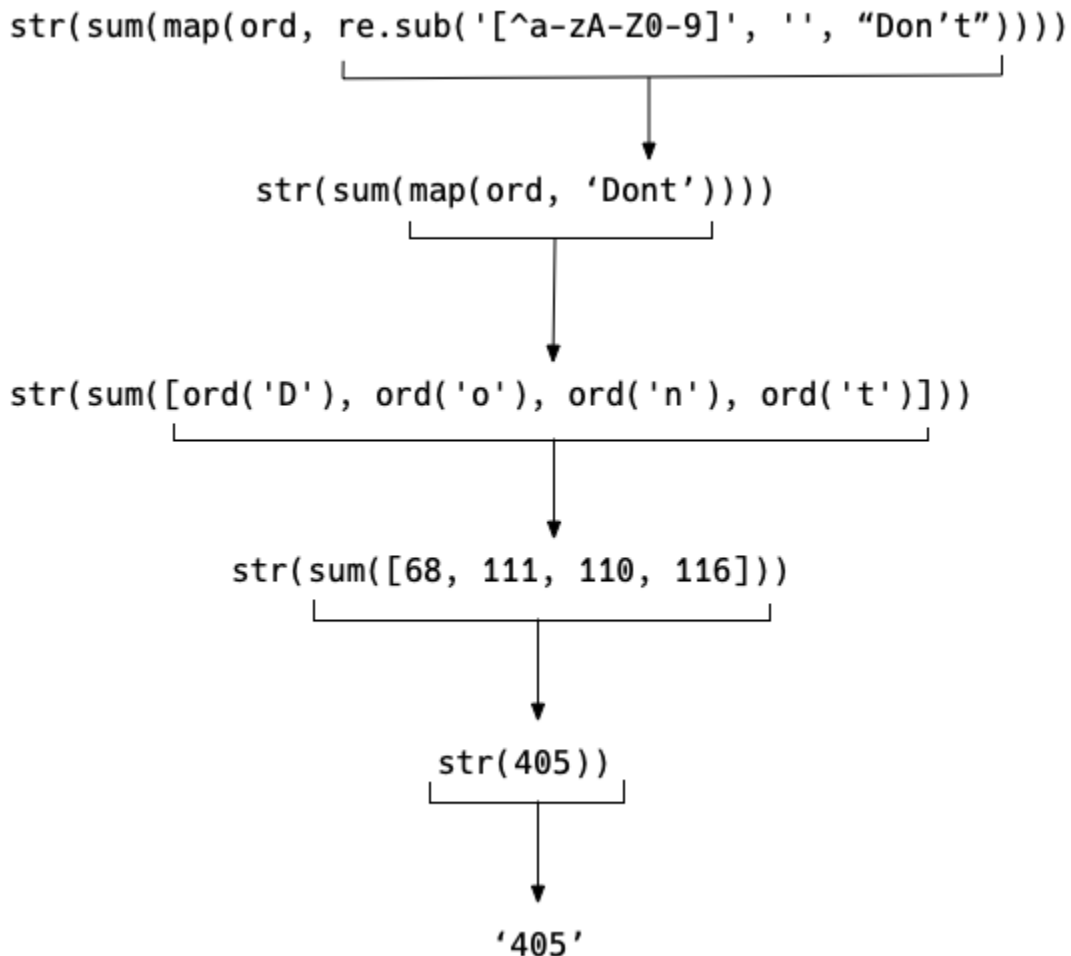
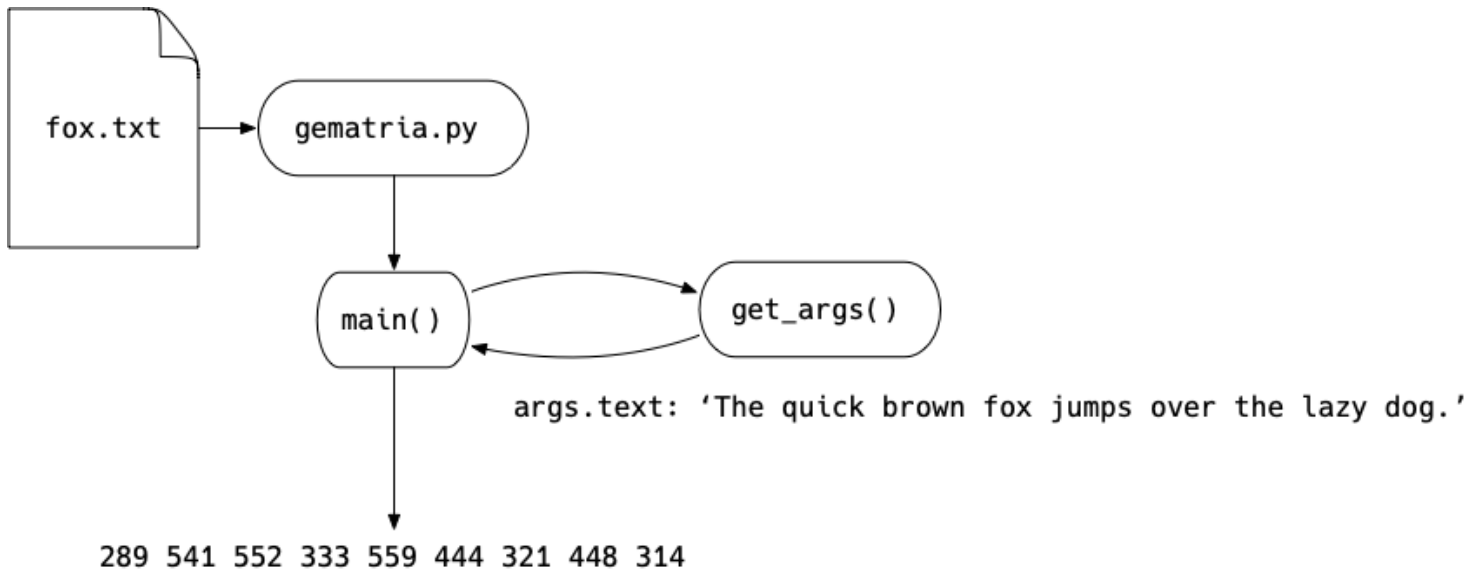


Figure 1. How the *for* loop, a list comprehension, and a *map* relate to each other.

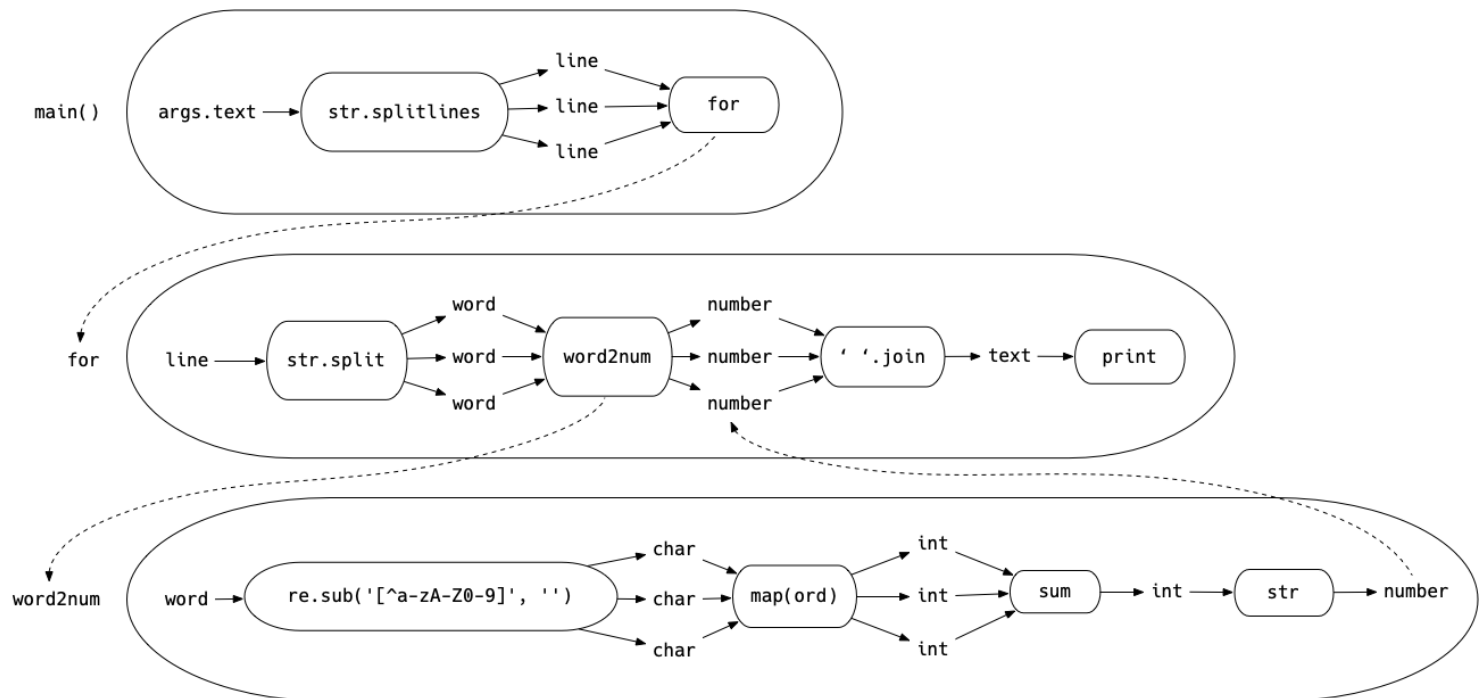
Here's a diagram to help you see how the data moves through the `map` version with the string "Don't":



From a higher level, here is how the data moves through the program as a whole:



And here is a string diagram showing how the data moves through the functions:



Sorting

The point of this exercise was less about the `ord` and `chr` functions and more about exploring regular expressions, function application, and how characters are represented inside programming languages like Python. For instance, sorting of strings is case-sensitive. Assume these words:

```
>>> words = 'banana Apple Cherry anchovies cabbage Beets'
```

Note that all the uppercase letters get sorted as a group and then the lowercase:

```
>>> sorted(words)
['Apple', 'Beets', 'Cherry', 'anchovies', 'banana', 'cabbage']
```

This is because all the uppercase ordinal values are lower than those of the lowercase letters. In order to perform a case-insensitive sorting of strings, you can use the `str.casefold` to get "a version of the string suitable for caseless comparisons" as the value to the `key` option:

```
>>> sorted(words, key=str.casefold)
['anchovies', 'Apple', 'banana', 'Beets', 'cabbage', 'Cherry']
```

The option is the same with `list.sort` if you prefer to sort the list in-place:

```
>>> words.sort(key=str.casefold)
>>> words
['anchovies', 'Apple', 'banana', 'Beets', 'cabbage', 'Cherry']
```

Unix command-line tools like the `sort` program behave in the same way due to the same representation of characters. Given a file of these same words:

```
$ cat words
banana
Apple
Cherry
anchovies
cabbage
Beets
```

The `sort` program will first sort the uppercase words and then the lowercase:

```
$ sort words
Apple
Beets
Cherry
anchovies
banana
cabbage
```

You have to read the `sort` manual page (via `man sort`) to find the `-f` flag to perform a case-insensitive sort:

```
$ sort -f words
anchovies
Apple
banana
Beets
cabbage
Cherry
```

Testing

I would like to take a moment to point out how often I use my own tests. Everytime I write an alternate version of a function or program, I run my own tests to verify that I'm not accidentally showing you buggy code. Having a test suite gives me the freedom and confidence to extensively refactor my programs because I know I can check my work. If I ever find a bug in my code, I add a test to verify that the bug exists. Then I fix the bug and verify that it's handled. I know if I accidentally reintroduce that bug, my tests will catch it!

Review

- The `ord` function will return the ordinal value of a character which is its position in the ASCII table.
- The `chr` function will return the character at a given position in the ASCII table.
- We can use character ranges like `a-z` in regular expressions because the characters lie contiguously in the ASCII table.
- The `re.sub` function will replace matching patterns of text in a string with new values such as replacing all non-characters with the empty string to remove punctuation and whitespace.
- A `map` can be written with a direct function reference instead of a `lambda` if the function expects a single positional argument.
- The `sum` function combines a list of numbers to a single sum. We could also say it "reduces" the list. Think of how that relates to the "map-reduce" concept we discussed before.
- To perform a case-insensitive sort of string values, use `key=str.casefold` option with both the `sorted` and `list.sort` functions.

Going Further

- Create a version that will handle non-English characters by using Unicode values for characters. You'll find that the idea of a "character" is quite different in Unicode. Something that appears to be one entity may actually be composed of more than one Unicode character, hence ideas like "grapheme clusters" ^[4] and "code points" ^[5].
- Analyze text files to find other words that sum to the value 666. Are these particularly scary words?

-
1. <https://en.wikipedia.org/wiki/Gematria>
 2. <https://en.wikipedia.org/wiki/ASCII>
 3. I included the `asciitbl.py` program I used to create this
 4. <https://unicode.org/reports/tr29/>
 5. https://en.wikipedia.org/wiki/Code_point

Last updated 2020-01-14 20:38:17 -0700