

# PROGRAMMING FOR HUMANISTS

STUART M. SHIEBER

## CONTENTS

<b>Part 1. Programming by imitation</b>	2
1. Where we're headed	4
2. Installing Python	5
3. The synoptic gospels	6
<b>Part 2. Programming from first principles</b>	12
4. Python documentation	13
5. The Python interpreter	14
6. Variables and the naming of values	15
7. Sequence data types	16
8. Functions	20
9. Words, types, and tokens	22
10. Files	23
11. Special characters	25
12. Splitting and joining strings	26
13. List comprehensions	27
14. Sets	29
15. Calculating with $n$ -grams	30
16. Defining your own functions	31
17. Dictionaries	33
18. Loops and conditionals	35
19. A concordance	38
Appendix A. A Concordance With Description	40
Appendix B. Statistics	46
References	47
Index	48

## Part 1. Programming by imitation

These notes are intended to provide an introduction to programming in the programming language `Python` for an audience of techno-savvy humanities scholars who are primarily interested in the use of computers for performing simple analyses of text. I originally prepared them for an audience of historians and philologists of premodern Europe, and the notes may reflect that audience, but should be appropriate for scholars from other disciplines as well.

There are two ways to learn a new language: by imitation and from first principles. This holds for both natural languages and programming languages. Under the `IMITATION` approach, learners see some examples and generate new examples by replacing parts of expressions they've seen. This approach has the benefit of allowing learners to use the language in interesting ways from early on, but they may do so without a full understanding of why the things they are saying work the way they do. Under the `FIRST PRINCIPLES` approach, learners study the elementary units of the language and how they are composed – the lexicon, grammar, and semantics of the language – and construct new examples from these first principles. This approach has the benefit that at every step the learner understands why the expressions work the way they do, but it may take a while to get to the point of being able to use the language to do much that is worthwhile.

For natural languages, the imitation approach is undoubtedly the preferred method. The lexicons and grammars of natural languages are large and complex and not well understood. Further, human beings have an ability to learn natural languages through immersion that allows even very young children to acquire a natural language with no explicit training in the first principles. Finally, the agents that understand natural languages are quite forgiving in their behavior. Fluent speakers can understand disfluent speech. So imperfections in the imitations don't have to hold up communication too much.

For programming languages, the case is somewhat different. Programming languages are artificial languages, and thus we cannot rely on innate language learning abilities. Furthermore, the agents that understand programming languages, computers, are quite unforgiving in their behavior. Even the most trivial variance from the well-formedness principles of the language may be met with utter failure to communicate the programmer's intent to the computer. On the other hand, the lexicons, grammars, and semantics of programming languages are much better understood than those of natural languages, because they have been explicitly designed and sometimes even specified with mathematical rigor. It is thus more practical to learn these first principles and apply them.

In these notes, I use both approaches, starting in this first part with the imitation method to get started and build some intuition and sense of what can be done, and then moving in the second part to the first principles that underly the language.

During this part, the idea is to merely get you used to the idea of commanding the computer to carry out calculations. Don't worry about the details of the language. Just let the code waft over you, like a pleasant sea breeze. Type the examples in and marvel at the results even if you can't fully understand yet why they work. Learn the following important lessons from the exercise:

- (1) *There's nothing to fear here.* You won't damage your computer by typing the wrong thing. You can experiment. If you wonder "what would happen if", just try it.
- (2) *First principles are important.* To really understand what's going on, the zen-like approach of Part 1 is insufficient. If you're motivated, move on to Part 2. Then go back to Part 1 afterwards and you'll see how much better you understand what's going on.

## 1. WHERE WE'RE HEADED

The coverage of these notes is not sufficient to make you a proficient Python programmer. They do not even provide a basic understanding of the full language. But the notes should get you to the point of writing simple programs to do basic text analyses. To get a sense of what can be achieved, by the end of working through these notes you'll have written code to generate a concordance of the text in Figure 1 (page 22) as found in Appendix A.

You'll also have enough familiarity with Python programming that it should be a simpler transition to learning about and working with the [Natural Language Toolkit](#) (NLTK), a free and open source Python toolkit for language processing that comes with its own book *Natural Language Processing with Python*.

Like all skills, programming requires practice. You don't get it by reading about it but by doing it. I recommend that you do *all* of the exercises and problems in these notes in order, even the ones that feel trivial, as well as playing around with small problems and tasks of your own devising.

### KEY CONCEPTS

**1.1. Conventions used in the notes.** First mentions of KEY CONCEPTS are shown in small caps and marked in the margins. You'll find them in the index at the end of the notes as well.

### CLICKABLE LINKS

The URLs provided in these notes, and some other items are CLICKABLE. Clickable links are shown in dark blue.

### EXERCISES PROBLEMS

There are EXERCISES and PROBLEMS interspersed throughout. The problems are more difficult than the exercises.

 Advanced material that can be skipped on first reading is marked as here.

**1.2. Disclaimer.** I apologize ahead of time for the rather breathless nature of these notes. They go through things quickly, and may be incomplete in various ways. You may (in fact likely will) have to augment them with reading in the Python documentation. On the other hand, I'll be available in class to answer questions, so there's that.

If you find errors or disfluencies in the notes, please let me know so that I can correct them.

## 2. INSTALLING PYTHON

Go no further without getting access to a Python interpreter. You'll want to try out the samples of Python code as they are presented and do your own experimentation as well.

**Mac OS:** Python is available natively on Mac OS. From a window in the Terminal application, type "python". The interpreter will be launched.

**Windows OS:** Python executables for Windows can be downloaded from <https://www.python.org/downloads/windows/>. Good luck with that. In case of failure, see the section below on web-based Python interpreters.

**Linux:** If you're running Linux, you're not going to need these notes.

**Web-based:** On any operating system with a browser, you can set up an account at [Pythonanywhere](#) and run a Python interpreter from within your browser. This will get you started for now.

The material in these notes is sufficiently straightforward that it probably makes little difference which version of Python you are running. However, for concreteness, all the examples below were run with Python 2.7.2.

**Exercise 1.** *Obtain access to a Python interpreter via one of the methods above.* □

**Exercise 2.** *Test that the Python interpreter is working by running it and typing in a simple command for the interpreter to execute. You should see something like this:*

```
% python
Python 2.7.2 (default, Oct 11 2012, 20:14:37)
[GCC 4.2.1 Compatible Apple Clang 4.0 (tags/Apple/clang-418.0.60)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
Finished loading pythonrc file
>>> 1+1
2
>>>
```

□

## 3. THE SYNOPTIC GOSPELS

We'll be looking primarily at text processing. Suppose we're interested in the synoptic gospels (and who isn't?). Each gospel is a text, which we can think of as a sequence of characters. Here, for instance, are the first four verses of the Gospel of Mark, generated using Python by *opening* a file named `Mark.txt` containing the Clementine Vulgate version of the Gospel of Mark, *reading* all of its lines into a list of lines, and then *extracting* the first four items in that list:

```
>>> open('Mark.txt').readlines()[:4]
['1:1 Initium Evangelii Jesu Christi, Filii Dei.\r\n', '1:2 Sicut scriptum est in Isaia propheta :
```

Let's give that list of lines a name. We'll call it `markLines`.

```
>>> markLines = open('Mark.txt').readlines()
>>> markLines[:4]
['1:1 Initium Evangelii Jesu Christi, Filii Dei.\r\n', '1:2 Sicut scriptum est in Isaia propheta :
```

☞ Notice that the expression `markLines[:4]` has exactly the same value as the previous expression `open('Mark.txt').readlines()[:4]`. This fact can be seen as an instance of Leibniz's law of the indiscernability of identicals. The command `markLines = open('Mark.txt').readlines()` has the effect of making `markLines` identical to `open('Mark.txt').readlines()`. Leibniz's law means that we can "substitute equals for equals". By substituting `markLines` for `open('Mark.txt').readlines()` in `open('Mark.txt').readlines()[:4]`, we get the equivalent `markLines[:4]`.

Viewing the list of lines that way isn't too readable. Here's a nicer presentation:

```
>>> for verse in markLines[:4]:
...     print verse,
...
1:1 Initium Evangelii Jesu Christi, Filii Dei.
1:2 Sicut scriptum est in Isaia propheta : [Ecce ego mitto angelum meum ante faciem tuam,/ qui pra
1:3 Vox clamantis in deserto :/ Parate viam Domini, rectas facite semitas ejus.]
1:4 Fuit Joannes in deserto baptizans, et praedicans baptismum poenitentiae in remissionem peccato
```

**Exercise 3.** *If the text of Matthew is in the file named `Matthew.txt`, how would you print out the first four verses of Matthew? The first six verses?* □

Instead of a list of lines (verses), it might be useful to extract a list of words. We'll start by joining all of the lines together, separated by, say, a colon.

```
>>> markString = ':'.join(markLines)
```

We can then take a look at the first few characters of this string. (Restricting to the first few avoids the whole giant string running off the end of the page.)

```
>>> markString[:60]
'1:1 Initium Evangelii Jesu Christi, Filii Dei.\r\n:1:2 Sicut s'
```

**Exercise 4.** What do you think the `[:60]` at the end does? Try substituting different numbers, like `[:5]` or `[:100]` and see what happens. □

**Exercise 5.** Suppose instead that you wanted to join the lines together with a space instead of a colon. How would you do that? □

Let's simplify and normalize the text a bit, by making it all lowercase.

```
>>> markLower = markString.lower()
>>> markLower[:60]
'1:1 initium evangelii jesu christi, filii dei.\r\n:1:2 sicut s'
```

**Exercise 6.** How would you assign the name `markUpper` to the uppercased text of `Mark`? □

The next step in extracting the words is to get rid of a bunch of characters that we aren't interested in – the chapter and verse markers for instance.

```
>>> markSimple = markLower.translate(None, '0123456789:')
>>> markSimple[:60]
' initium evangelii jesu christi, filii dei.\r\n sicut scriptum'
```

There are other characters we may want to remove, punctuation and newlines and such, so let's redo the process with a broader set of characters to exclude.

```
>>> markSimple = markLower.translate(None, '\n\r,.;\|/')(?0123456789:')
>>> markSimple[:60]
' initium evangelii jesu christi filii dei sicut scriptum est'
```

Finally, let's get rid of any extraneous `WHITESPACE` – the nonprinting layout characters like spaces, tabs, and newlines – at the start and end of the string. WHITESPACE

```
>>> markSimple = markSimple.strip()
>>> markSimple[:60]
'initium evangelii jesu christi filii dei sicut scriptum est '
```

Now, we can split the string into the component words at the whitespace that separate the words.

```
>>> markWords = markSimple.split()
>>> markWords[:7]
['initium', 'evangelii', 'jesu', 'christi', 'filii', 'dei', 'sicut']
```

Let's encapsulate this whole process of turning a file into the list of words by defining a *function* that carries out that process.

```
>>> def wordsNormed(fileName):
...     return ' '.join(open(fileName).readlines())
...     .lower()
```

```

...         .translate(None, '\n\r,.;\|/[]()?'0123456789') \
...         .strip() \
...         .split()
...

```

Now we can do that for several different documents.

```

>>> matthew = wordsNormed('Matthew.txt')
>>> mark = wordsNormed('Mark.txt')
>>> luke = wordsNormed('Luke.txt')
>>> john = wordsNormed('John.txt')

```

To make sure it worked, let's look at the first few words of each.

```

>>> matthew[:7]
['liber', 'generationis', 'jesu', 'christi', 'filii', 'david', 'filii']
>>> mark[:7]
['initium', 'evangelii', 'jesu', 'christi', 'filii', 'dei', 'sicut']
>>> luke[:7]
['quoniam', 'quidem', 'multi', 'conati', 'sunt', 'ordinare', 'narrationem']
>>> john[:7]
['in', 'principio', 'erat', 'verbum', 'et', 'verbum', 'erat']

```

Let's look at some contiguous word sequences from the gospels. Here's the third through fifth words in Mark.

```

>>> mark[2:5]
['jesu', 'christi', 'filii']

```

(Even though we want the third through fifth words, we use the numeric indices 2 and 5. You'll see why later in Section 7.1.)

How about generating a whole series of such three word sequences? Contiguous sequences of  $n$  words in a document are called  $n$ -grams; in the case where  $n$  is 3, they are called trigrams. Here are the first ten trigrams in Mark.

```

>>> mark10trigrams = [mark[i:i+3] for i in range(10)]
>>> for trigram in mark10trigrams:
...     print trigram
...
['initium', 'evangelii', 'jesu']
['evangelii', 'jesu', 'christi']
['jesu', 'christi', 'filii']
['christi', 'filii', 'dei']
['filii', 'dei', 'sicut']
['dei', 'sicut', 'scriptum']
['sicut', 'scriptum', 'est']
['scriptum', 'est', 'in']

```

```
['est', 'in', 'isaia']
['in', 'isaia', 'propheta']
```

We can define a process to generate a list of *all* of the trigrams in a list of words.

```
>>> def ngrams(aList, N=3):
...     return [aList[i:i+N] for i in range(len(aList)-N+1)]
... 
```

 The argument specification `N=3` means that the second argument named `N` is *OPTIONAL*, and if it is not provided, a default value of 3 will be used as its value. Thus `ngrams` by default computes trigrams, but can also be used to compute *n*-grams for other values of *n* if desired.

OPTIONAL ARGUMENTS

**Exercise 7.** *Why is the range limit `len(aList)-N+1` rather than just `len(aList)`? What is the point of the extra arithmetic? Hint: Try it with just `len(aList)` and see what happens.* □

Let's test it on Mark again, printing the first few trigrams found to verify that it worked.

```
>>> markTrigrams = ngrams(mark)
>>> for trigram in markTrigrams[:5]:
...     print trigram
... 
```

```
['initium', 'evangelii', 'jesu']
['evangelii', 'jesu', 'christi']
['jesu', 'christi', 'filii']
['christi', 'filii', 'dei']
['filii', 'dei', 'sicut']
```

For completeness, we can generate the trigrams in the other gospels as well.

```
>>> matthewTrigrams = ngrams(matthew)
>>> lukeTrigrams = ngrams(luke)
>>> johnTrigrams = ngrams(john)
```

One way to measure the similarity of two documents is to examine what trigrams (or other *n*-grams) they have in common. We start by defining the intersection of two lists, that is, the items they have in common:

```
>>> def intersect(list1, list2):
...     return [item
...             for item in list1
...             if item in list2]
... 
```

Now we can find all of the trigrams in common between Matthew and Mark:

```
>>> commonMatthewMark = intersect(matthewTrigrams, markTrigrams)
>>> for common in commonMatthewMark[:5]:
```

```

...     print common
...
['jesu', 'christi', 'filii']
['quod', 'est', 'interpretatum']
['cum', 'illo', 'et']
['principes', 'sacerdotum', 'et']
['at', 'illi', 'dixerunt']

```

How many such common trigrams are there?

```

>>> len(commonMatthewMark)
1906

```

That's about 18 percent of the Mark trigrams.

**Exercise 8.** *Knowing the raw count of common  $n$ -grams may not be as useful as knowing the proportion of common  $n$ -grams. What is the proportion of the Mark trigrams that are also found in Matthew?* □

Is that a lot? We can compare it against the proportion of trigrams found in some other more or less unrelated Latin document. Let's use the *Vita Sancti Germani*.

```

>>> vsgTrigrams = ngrams(wordsNormed('vsg.txt'))
>>> len(intersect(markTrigrams, vsgTrigrams))
13

```

The 13 common trigrams accounts for only 0.13 percent. So (unsurprisingly) Mark looks to be extremely similar to Matthew.

Let's make a table that shows how similar the gospels are to each other (at least as measured by common trigrams).

```

>>> gospels = {'Matthew': matthewTrigrams,
...           'Mark': markTrigrams,
...           'Luke': lukeTrigrams,
...           'John': johnTrigrams}
>>> N = 3
>>> for (g1, w1) in gospels.items():
...     for (g2, w2) in gospels.items():
...         print "{:10s} {:10s} {:10.3%}"\
...             .format(g1, g2,
...                     float(len(intersect(w1, w2)))
...                     / (len(w1) - N + 1))
...
Matthew  Matthew  100.012%
Matthew  Luke         12.985%
Matthew  John         2.586%
Matthew  Mark         11.517%

```

Luke	Matthew	11.351%
Luke	Luke	100.011%
Luke	John	2.206%
Luke	Mark	7.553%
John	Matthew	3.286%
John	Luke	3.485%
John	John	100.014%
John	Mark	3.116%
Mark	Matthew	17.127%
Mark	Luke	12.220%
Mark	John	3.065%
Mark	Mark	100.019%

**Exercise 9.** Which of the gospels is the outlier? That is, which is the most different from all the others?

**Exercise 10.** What about common 5-grams? Generate the same table but for 5-grams.

**Part 2. Programming from first principles**

The first part of these notes should give you an idea of how even a few lines of Python code can accomplish some serious textual analysis. But to really understand how to program, so that you can generate effective code directly and not merely program by analogy, you need to understand the first principles of the programming language. In this part, we present some of these first principles for Python in a graded manner with interspersed exercises.

#### 4. PYTHON DOCUMENTATION

These notes are not self-contained – on purpose. Python is a large language, with many built-in functions and add-on modules for doing all kinds of things. All are well documented at the [python.org](https://python.org) web site. You'll want to get in the habit of heading there to look up aspects of the language that you need help with.

Here are some especially important bits:

- There is a tutorial on the language at <https://docs.python.org/2/tutorial/index.html>, which you may find complementary to these notes. It does assume a bit of programming background.
- The language reference manual is at <https://docs.python.org/2/reference/index.html>.
- The Python standard library and modules are described at <https://docs.python.org/2/library/index.html>. We use some of these below, for instance, standard functions like `sorted` and the `pprint` module.

## 5. THE PYTHON INTERPRETER

## INTERPRETER

A Python INTERPRETER allows you to specify calculations as Python expressions or programs and calculates the result of those specifications. You type Python commands and expressions into the interpreter, and the interpreter executes the commands and calculates the values of the expressions printing a representation of the calculated values.

COMMANDS  
EXPRESSIONS

 We distinguish COMMANDS and EXPRESSIONS. Commands are executed for their side effects. Expressions are executed for their values (though they may have side effects as well). The difference is revealed by the interpreter: after entering a command, no output is printed by the interpreter; after entering an expression, an output is printed, namely, the expression's value.

Here is a simple example of using a Python interpreter. The user's input is on the lines beginning '>>>' (or '.' for lines continuing a single input) and the interpreter's output immediately follows.

```
>>> 3 + 4 * 5
23
```

We've used the + symbol for addition and \* for multiplication. You can find a larger listing of arithmetic operators at [https://en.wikibooks.org/wiki/Python\\_Programming/Basic\\_Math](https://en.wikibooks.org/wiki/Python_Programming/Basic_Math).

**Exercise 11.** Enter the expression  $3 + 4 * 5$  into the Python interpreter and verify that it works like it should. □

**Exercise 12.** Use the Python interpreter to determine the values of the following arithmetic expressions:

- (1)  $4/4 - 4/4$
- (2)  $\frac{4+4}{4+4}$
- (3)  $\frac{4-4}{4+4}$

This exercise is inspired by the "four fours" puzzle, which involves constructing arithmetic expressions for each positive integer using four fours combined however you want. Feel free to generate more examples and use Python to verify them for you. □

## 6. VARIABLES AND THE NAMING OF VALUES

We can name the results of computations for later use. These names are called **VARIABLES**. Variables are alphanumeric names that start with an alphabetic character (or underscore “\_”). Here’s an example of the use of a variable to name a value and then using that value in later computations.

```
>>> largeSquare = 128 ** 2
>>> largeSquare / 2
8192
```

The first line constitutes an **ASSIGNMENT**; it assigns the name given on the left side of the = operator to the value specified by the expression on the right side. Thus the variable `largeSquare` names the value 16384. Assignments are executed for their *effect*, not their *value*. For that reason, the interpreter doesn’t print anything after this line. (Don’t be confused. The = does not mean “is equal to”, as it does in standard mathematical notation. It’s a kind of command, not a statement of fact.)

The second line then uses that variable by dividing its value by 2. The interpreter prints the value specified by that last expression.

## 7. SEQUENCE DATA TYPES

It is conventional in defining programming languages to carefully distinguish the different types of data that programs can manipulate. We've seen one `DATA TYPE` already – numbers.

 In actuality, Python treats numbers as falling into a set of different data subtypes: integers, real numbers, complex numbers, each of which operates slightly differently.

Our primary application in these notes is analysis of text. We will therefore move quickly to look at the data type most useful for representing text, namely, strings. Strings are a kind of sequence data type; a string is essentially a sequence of characters. In fact, Python provides several different data types for sequences: strings of course, but also lists and tuples. These sequence data types share many properties, so we introduce them together.

`LIST` 7.1. **Lists.** The Python `LIST` data type is used to represent sequences of other data objects, sequences that can be adjusted in various ways, for instance, by adding or removing elements. The notation for lists is to place the individual listed objects, separated by commas and surrounded by brackets.

```
>>> [1, 2, 3]
[1, 2, 3]
>>> aList = [1, 4, 1, 5, 9, 2, 6]
>>> aList
[1, 4, 1, 5, 9, 2, 6]
```

`POSITION` Each item in a list has its own `POSITION` in the list. The individual items within a `INDEXING` list can be extracted by `INDEXING` them based on their respective positions. We use the indexing notation `·[·]`. For instance, to retrieve the fifth item from `aList`, we use the notation `aList[4]`.

```
>>> aList[4]
9
```

Notice that the value of this expression is indeed the fifth item in the list, the number 9.

Why use the index 4 for the fifth item? Because we think of the positions as being numbered *starting from index zero*. Alternatively, you can think of the indices as numbering the points *between* the items, starting with zero, like in this picture.

1	4	1	5	9	2	6	
0	1	2	3	4	5	6	7

Under this conception, the indexing `aList[4]` extracts the item *following* position 4, that is, the fifth item.

7.2. **Sequence lengths.** We may want to know how many items there are in one of these kinds of sequences. We use the `len` function to calculate the length of a list. (We'll have much more to say about functions shortly, starting in Section 8.)

LEN FUNCTION

```
>>> len(aList)
7
```

Since the length of a list is a number, you can operate on it as you would any other number, applying arithmetic operations to it for instance.

```
>>> len(aList) * 2
14
```

7.3. **Strings.** We'll use the `STRING` data type for representing text. Strings in Python are specified by enclosing a sequence of characters within matching string `DELIMITERS`, such as single quotes.

STRING

DELIMITERS

```
>>> 'sanctus Germanus'
'sanctus Germanus'
```

Strings can be specified with other delimiters, such as double quotes, or triple double or single quotes.

```
>>> "This example uses double quotes"
'This example uses double quotes'
>>> """Triple quotes are
... often used for
... multi-line strings."""
'Triple quotes are\noften used for\nmulti-line strings.'
```

Note that Python always prints out the strings using the single quote delimiter.

 This last string has some `NEWLINE` characters in it. They're specified with the `'\n'` characters. See Section 11 below.

NEWLINE

Strings can be concatenated using the `+` operator.

```
>>> "This" + ' that'
'This that'
```

(We can freely combine strings specified with the different delimiters.)

Like all data values, strings can be named by variables.

```
>>> aString = " be as it were as it"
>>> "Let it" + aString * 2 + " were"
'Let it be as it were as it be as it were as it were'
```

Interesting how Python uses the “multiplication” operator `*` for repeating strings, no? This “arithmetic” on strings works for lists as well.

```
>>> motto = [ "nihil", "agere", "delectat" ]
>>> motto
['nihil', 'agere', 'delectat']
```

```

>>> len(motto)
3
>>> motto + motto
['nihil', 'agere', 'delectat', 'nihil', 'agere', 'delectat']
>>> len(motto * 2) - len(motto) * 2
0
>>> aString
' be as it were as it'
>>> len(aString)
20

```

**Exercise 13.** *What will Python print in response to each of the following inputs?*

```

aList = [ "agere", "delectat", "nihil" ]
aList[2] + aList[0] + aList[1]
aList[2] + " " + aList[0] + " " + aList[1]
aList[1][1] + aList[2][2]
len(aList * 2) - len(aList) * 2

```

□

**7.4. Substrings.** Strings, like lists, are sequences – in particular, sequences of characters. We can do many of the same operations on strings that we can on lists. For instance, we can extract a character from a string using the same indexing notation `[·]`. To retrieve the fifth character from `aString`, we use the notation `aString[4]`.

```

>>> aString = "sanctus Germanus"
>>> aString[4]
't'

```

As before we think of the indices as numbering the points *between* the characters, starting with zero, like in this picture.

```

  s   a   n   c   t   u   s   _   G   e   r   m   a   n   u   s
0   1   2   3   4   5   6   y   8   9   10  11  12  13  14  15  16

```

Under this conception, the indexing `aString[4]` extracts the character *following* string position 4, that is, the fifth character.

## SLICING

Substrings can be specified by a SLICING notation, similar to the indexing notation but providing both starting and ending positions within the full string, separated by a colon. For instance, to extract the substring between string positions 2 and 6 (that is, the second through fifth characters):

```

>>> aString[2:6]
'nctu'

```

**Exercise 14.** *What strings are specified by the following Python expressions? Recall the value of `aString` defined above.*

- (1) `aString[0:3]`
- (2) `aString[3]`
- (3) `aString[3:4]`
- (4) `aString[3:3]`
- (5) `aString[3:2]`
- (6) `aString[:4]`
- (7) `aString[4:]`
- (8) `aString[4:-3]`
- (9) `aString[3:100]`
- (10) `aString[8:0:-1]`
- (11) `aString[::-1]`

We really haven't given enough detail about how the indexing notation works to determine all of these, so you'll have to experiment to figure them out. □

**Exercise 15.** Based on your experiments with the previous exercise, how would you reverse a string in Python, that is, generate a string with the characters in the reverse order? □

**Exercise 16.** This method that allows extracting substrings from strings also allows extracting sublists from lists. Suppose the variable `vsgList` names the value `['sanctus', 'Germanus', 'abba', 'et', 'martyr']`. How would you extract all but the first and last elements from the list? □

**Exercise 17.** How would you extract the final two elements from `vsgList`, without recourse to prior knowledge of the number of items in the list? □

7.5. **Tuples.** The final sequence data type we'll cover is the `TUPLE`. The name `TUPLE` derives from the suffix seen in quintuple, sextuple, septuple, and the like.

[[[To be written.]]]

## 8. FUNCTIONS

FUNCTIONS  
ARGUMENTS  
RESULT

Data – numeric or string values, and all the other types of data that Python makes available – are manipulated through the application of **FUNCTIONS**, engines that take inputs, called **ARGUMENTS**, and transform them into an output, the **RESULT**. We’ve seen examples of such functions already: the arithmetic and string operators like `+` and `*`, indexing operators like `[:]`. These are special built-in functions that are invoked via special “idiomatic” notations. The arithmetic operators, for instance, are written infix, as, e.g., `1 + 2`, and the indexing operator is written with brackets.

But in general, Python uses two notations that are more uniform for applying a function to its arguments.

MATHEMATICAL NOTATION

- (1) *Mathematical notation*: Mimicking a traditional **MATHEMATICAL NOTATION** due originally to Leibniz, a function, say  $f$ , applied to its arguments is notated by placing the comma-separated arguments after the function in parentheses, viz.,

$$f(\langle arg1 \rangle, \langle arg2 \rangle, \dots) \quad .$$

OBJECT NOTATION

- (2) *Object notation*: A second notation, **OBJECT NOTATION**, derived from conventions used in so-called object-oriented programming languages, places the function *after* its first argument separated by a dot, with all other arguments following as in the mathematical notation, viz.,

$$\langle arg1 \rangle.f(\langle arg2 \rangle, \dots) \quad .$$

☞ The latter notation makes more sense once Python’s status as an object-oriented language is understood, but in the interest of introducing the least language for our purposes, we introduce it as just a fixed idiom.

Any given function uses either the first or second notation, in much the same way that any given Latin verb inflects as per one of a small set of conjugations. You might think of functions that use the mathematical notation “first conjugation” functions and those using object notation “second conjugation”.

LEN FUNCTION

An example of the mathematical notation is the built-in `len` function, which takes a single argument and returns its length. It can be applied to any kind of list, and in particular, to strings, for instance,

```
>>> len(aString)
16
```

**Exercise 18.** *What do the following Python expressions return?*

- (1) `aString[0:len(aString)]`
- (2) `aString[1:len(aString)]`
- (3) `aString[0:len(aString)-1]`

Can you find simpler ways of getting the same values? □

Another useful function is the built-in `sorted` function, which takes a single argument representing a sequence (such as a list or string) and returns a corresponding object representing the elements of its argument in sorted order. SORTED FUNCTION

```
>>> sorted([3, 1, 4, 1, 5])
[1, 1, 3, 4, 5]
```

**Exercise 19.** Recall the value of `motto`, which is `['nihil', 'agere', 'delectat']`. What do the following Python expressions return?

- (1) `sorted(motto)`
- (2) `sorted(motto[0])`

□

It is often useful to generate a list of sequential numbers. We'll see use examples later. The `range` function serves that purpose. Its two arguments specify the start and end of the range; the included numbers are obtained by starting with the first, and incrementing repeatedly until the second number is reached (or surpassed). If the first argument is left off, it is assumed to be 0. If a third argument is added, it is taken to be the increment used between numbers. RANGE FUNCTION

```
>>> range(5, 10)
[5, 6, 7, 8, 9]
>>> range(10)
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
>>> range(1, 10, 2)
[1, 3, 5, 7, 9]
>>> range(10, 0, -1)
[10, 9, 8, 7, 6, 5, 4, 3, 2, 1]
```

As a final example, we consider the `count` function (which uses the object notation), which counts the number of occurrences of its second argument as elements of its first list argument. COUNT FUNCTION

```
>>> motto.count('nihil')
1
>>> motto.count('ipsum')
0
>>> motto[0].count('i')
2
```

## 9. WORDS, TYPES, AND TOKENS

As we turn to processing of text, some standard terminology about words is useful, starting with the word “word” itself. The question of what is a word is itself somewhat fraught. For the time being, we’ll just consider the words in a text to be the maximal sequences of alphabetic characters separated by whitespace. (As it turns out, this is an exceptionally poor definition, but sufficient for the time being.)

TYPES  
TOKENS

We distinguish word `TYPES` from word `TOKENS`. A text is made up of a series of word tokens. Each word token belongs to a word type. Consider the text corpus in Figure 1, a sentence from Gertrude Stein’s 1929 poem “An Acquaintance With Description” (Stein, 1929). This corpus has 225 word tokens (ignoring punctuation), which are instances of just eight word types (if we conflate upper and lower case). The eight types, in decreasing order of frequency, are: “be”, “to”, “it”, “sure”, “let”, “mine”, “when”, “is”. Each of these word types has several occurrences as tokens in the poem.

---

Let it be when it is mine to be sure let it be when it is mine when it is mine  
 let it be to be sure when it is mine to be sure let it be let it be let it be to be  
 sure let it be to be sure when it is mine to be sure let it to be sure when it  
 is mine let it be to be sure let it be to be sure to be sure let it be to be sure  
 let it be to be sure to be sure let it be to be sure let it be to be sure let it be  
 to be sure let it be mine to be sure let it be to be sure to be mine to be sure  
 to be mine to be sure to be mine let it be to be mine let it be to be sure to  
 be mine to be sure let it be to be mine let it be to be sure let it be to be sure  
 to be sure let it to be sure mine to be sure let it be mine to let it be to be  
 sure to let it be mine when to be sure when to be sure to let it to be sure to  
 be mine.

---

FIGURE 1. A sentence from Stein’s “An Acquaintance with Description” (1929).

## 10. FILES

10.1. **Strings from files.** Typing in the kinds of long strings we'll be analyzing, entire books in some cases, is painful. Better to store the text in a text file and load that file into Python. Let's imagine that we have a file called "stein.txt" that contains the line from Figure 1. We want to read that file into Python so that we can operate with it.

We'll use an idiom to get the contents of a text file into a variable. The idiom is this:

```
contents = open('filename.txt').readlines()
```

We are using two different functions in this idiom, the `open` function, invoked using the mathematical function notation, and the `readlines` function, invoked using the object notation. The `open` function takes a single string as an argument, and returns as value an object that designates the file with that name. The `readlines` function's first argument is a file designator (as returned by `open`), and since it takes no further arguments, the parentheses for the remaining arguments are empty. The function returns a list, each component of which is a string containing a line of the file that was read in.

To read the Stein poem in, we can therefore use:

```
>>> steinLines = open('stein.txt').readlines()
```

Now, let's examine what we've read in.

☞ Instead of just evaluating (and having the interpreter print the value of) `stein`, here we are "importing" a special "pretty-printing" facility, the `pprint` function, to print the value of `stein` in a more attractive manner.

PPRINT FUNCTION

```
>>> from pprint import pprint
>>> pprint(steinLines)
['Let it be when it is mine to be sure let\n',
 'it be when it is mine when it is mine\n',
 'let it be to be sure when it is mine to\n',
 'be sure let it be let it be let it be to\n',
 'be sure let it be to be sure when it is\n',
 'mine to be sure let it to be sure when\n',
 'it is mine let it be to be sure let it\n',
 'be to be sure to be sure let it be to be\n',
 'sure let it be to be sure to be sure let\n',
 'it be to be sure let it be to be sure\n',
 'let it be to be sure let it be mine to\n',
 'be sure let it be to be sure to be mine\n',
 'to be sure to be mine to be sure to be\n',
 'mine let it be to be mine let it be to\n',
```

```
'be sure to be mine to be sure let it be\n',  
'to be mine let it be to be sure let it\n',  
'be to be sure to be sure let it to be\n',  
'sure mine to be sure let it be mine to\n',  
'let it be to be sure to let it be mine\n',  
'when to be sure when to be sure to let\n',  
'it to be sure to be mine.\n']
```

**Exercise 20.** *Read into Python the contents of a text file for some document you are interested in. The Vita Sancti Germani comes to mind.* □

## 11. SPECIAL CHARACTERS

Let's examine the first line of the poem.

```
>>> steinLines[0]
'Let it be when it is mine to be sure let\n'
```

It's a string of 41 characters.

**Exercise 21.** *How could you verify that length? Do it.*

**Exercise 22.** *Use Python to extract the last character from the first line of the poem.*

The last character of the first line is the newline character, which unlike all the “normal” characters, is notated with an `ESCAPE SEQUENCE`, a backslash followed by an `n`: `'\n'`. There are other escape sequences, used for characters that are otherwise hard to make clear in a printed representation, such as `'\t'` for the tab character or `'\''` for the single quote character (which is otherwise hard to put in a single-quoted string without prematurely terminating the string).

[ESCAPE SEQUENCE](#)

**Exercise 23.** *How would you notate the single-quoted string containing the possessive form of your first name?*

## 12. SPLITTING AND JOINING STRINGS

We introduce some useful string manipulation functions. To concatenate together a list of strings to form a single string, use the `join` function that takes a separator string and a list of strings to join and combines the strings in the list together separated by the separator string.

JOIN FUNCTION

```
>>> ' '.join(['sanctus', 'Germanus'])
'sanctus Germanus'
```

**Exercise 24.** Use `join` to generate the following strings from the list of number strings `['1', '2', '3']`.

- (1) `'1-2-3'`
- (2) `'1, 2, 3'`
- (3) `'123'`
- (4) `'3, 2, 1'`

For the last problem, recall Exercise 15. For further extra credit, start from the list of numbers themselves `[1, 2, 3]`. Check out the functions `map` and `str`. □

SPLIT FUNCTION

The converse of the `join` function is the `split` function. Again, `split` takes two arguments in object notation. The first is the string to be split up into substrings and the second is a string that specifies where to split. Each occurrence of the second string in the first string generates a split point. To split at the spaces in the string, then, the second argument would be the string `' '`:

```
>>> line = "He told me you had been to her and mentioned me to him"
>>> line.split(' ')[0:5]
['He', 'told', 'me', 'you', 'had']
```

The splitting can occur at any substring we want:

```
>>> line.split(' me ')
['He told', 'you had been to her and mentioned', 'to him']
```

**Exercise 25.** Extra credit: What is this line from? □

LOWER FUNCTION

**Exercise 26.** Use Python's `lower` function (*inter alia*) to generate the list of word tokens in `line` but with all words in lower case. Step one: Click on the link in this exercise to go to the Python documentation on the `lower` function. While you're there, look around at the range of *other string-processing functions* that may come in handy some day. □

**Exercise 27.** Use Python to split the first line of Stein's poem into its separate word tokens, storing the resulting list of tokens in the variable `steinWords1`. □

## 13. LIST COMPREHENSIONS

We've seen the notation for specifying a list **EXTENSIONALLY**, that is, by enumerating its elements explicitly. Here for instance are the first letters of the first few words (the first eight, say) in the first line of the Stein poem, enumerated explicitly:

```
>>> firstLetters = ['L', 'i', 'b', 'w', 'i', 'i', 'm', 't']
>>> firstLetters
['L', 'i', 'b', 'w', 'i', 'i', 'm', 't']
```

EXTENSIONAL

It's much more elegant and less error-prone to let Python do the work for you. We use **LIST COMPREHENSIONS** for the task. List comprehensions allow specifying a single generic list element computation that captures all of the elements of the list. It allows defining lists **INTENSIONALLY** rather than extensionally. The list comprehension notation is

LIST COMPREHENSIONS

INTENSIONAL

```
[ <generic element> for <variable> in <list> ]
```

☞ For the mathematically inclined, it may be useful to think of this notation as analogous to the familiar mathematical notation for defining sets intensionally, for example,

$$\{x^2 \mid 0 \leq x < 10\}$$

which defines the set containing the first 10 squares. The braces become brackets in Python, and the vertical bar becomes the word **for**, which separates the generic element  $x^2$  on its left from the specification of the possible values of  $x$  on its right.

For the current example, each element of the list can be calculated as `word[0]` where `word` is one of the first few words in the first line of the poem. (Recall that the words in the first few lines in the poem are named by the variable `steinWords` from Exercise 27.)

```
>>> firstLetters = [word[0] for word in steinWords1[0:8]]
>>> firstLetters
['L', 'i', 'b', 'w', 'i', 'i', 'm', 't']
```

Here, the variable `word` takes on each element of the list `steinWords[0:8]`, and for each one, an element of the list is computed as `word[0]`.

**Exercise 28.** Generate a list each element of which is a list of all of the word tokens in a line of the Stein poem. □

**Exercise 29.** Generate a list named `steinWords` of all the word tokens in the Stein poem. Make sure that all the words are lower case. You may find the `strip` function to be useful. You should be able to get the following behavior:

STRIP FUNCTION

```
>>> steinWords[6:12]
['mine', 'to', 'be', 'sure', 'let', 'it']
```

□

**Exercise 30.** *Generate a list of the first 10 squares (0, 1, 4, 9, etc.). Hint: You'll want to recall the `range` function.*

□

## 14. SETS

Time to introduce another data type, the `SET`. A set is a compound data type; like the list, each set contains elements. But the elements of a set are unique. A set does not contain multiple tokens of the same value. You can create a set from a list with the `set` function.

SET

SET FUNCTION

```
>>> set([1, 2, 3])
set([1, 2, 3])
>>> set([1, 2, 3, 2, 1])
set([1, 2, 3])
>>> set('how much would would a woodchuck chuck'.split(' '))
set(['a', 'would', 'chuck', 'how', 'much', 'woodchuck'])
```

As you can see, the printed representation for a set shows a list of the elements but still marks it as a set.

Many of the same functions that apply to lists apply to sets as well: `len` for counting the number of elements, `+` for combining two sets (taking their UNION), etc.

UNION

**Exercise 31.** Use Python to calculate how many word types (not tokens) there are in the Stein poem. Ignore case distinctions. Hint: The answer is 8. The hint is to emphasize that the point of the exercise is the code, not the answer. □

📖 The elements of a set can be of many types – numbers, strings, and tuples, in particular – but unfortunately not lists or sets. Only HASHABLE data types are allowed.

HASHABLE

15. CALCULATING WITH  $n$ -GRAMS $n$ -GRAMS

We'll spend some time looking at  $n$ -GRAMS, contiguous sequences of  $n$  words. When  $n$  is 1, 2, or 3, we call them unigrams, bigrams, and trigrams, respectively. Here are some examples of trigrams built from the vocabulary seen in Gertrude Stein's poem:

- (1) let it be
- (2) it is mine
- (3) it is sure
- (4) to be sure

**Problem 32.** *Generate a list of all of the trigram tokens in the Stein poem. You'll want to use the word list you generated in Exercise 29.*

**Problem 33.** *How many times do each of the four sample trigrams above occur in the poem? If you resort to counting them yourself, go back to the beginning of these notes and start over.*

**Exercise 34.** *How many unique trigrams are there in the Stein poem? (You may want to look at the earlier discussion about hashable data types.)*

## 16. DEFINING YOUR OWN FUNCTIONS

Functions like `len`, `sorted`, `count`, and the like can be fabulously useful. If there's a function that does just what you need, a single line of code can accomplish your purposes.

Sadly, there often is not a function tailor-made for your purposes. But you can write your own. Indeed, writing functions is the heart of computer programming (in spite of the fact that it took until page 31 to get to the topic).

In Python, you can define your own function of zero or more arguments using the **def** command. The notation is as follows:

```
def <function name>(<arguments>):
    <function body>
```

DEF COMMAND

Within the body of the function, the **return** command generates the value to return as the result of the function.

RETURN COMMAND

For example, here we define a function to calculate the first letter of a string.

```
>>> def firstLetter(aString):
...     return aString[0]
... 
```

We can use this function by CALLING it just as we would a built-in function using mathematical notation:

FUNCTION CALL

```
>>> firstLetter('nihil')
'n'
>>> firstLetter(steinLines[0])
'L'
```

**Exercise 35.** Define and test a function that returns the last letter of the first word in a string. □

**Exercise 36.** Define and test a function that returns the reversal of a string or list. □

**Exercise 37.** Define and test a function that returns the alphabetically first word in a string. □

**Exercise 38.** Define and test a function that returns the middle element of a list, that is, the element that has the same number of elements before and after it. (If the list has an even number of elements, the chosen element should have one more element before than after.) □

**Exercise 39.** Define and test a function that returns a list of all the trigram tokens in a list of tokens. For instance, it should have the following behavior:

```
>>> pprint(ngrams(motto * 2))
[('nihil', 'agere', 'delectat'),
 ('agere', 'delectat', 'nihil'),
```

```
('delectat', 'nihil', 'agere'),  
( 'nihil', 'agere', 'delectat')]
```

*Test it on the Stein poem.*

□

**Exercise 40.** Define and test a function that returns a set of all trigram types in a list of tokens. For instance, it should have the following behavior:

```
>>> pprint(ngramSet(motto * 2))  
set([('agere', 'delectat', 'nihil'),  
     ('delectat', 'nihil', 'agere'),  
     ('nihil', 'agere', 'delectat')])
```

*Test it on the first line of the Stein poem.*

□

## 17. DICTIONARIES

A **DICTIONARY** is a data structure for associating one kind of data with another. We might want to associate words with their locations in a document, or  $n$ -grams with their number of occurrences, or any of a variety of other associations.

In Python, a dictionary can be specified extensionally using a notation with braces. Here, we build a dictionary that associates a few words with their length.

```
>>> lengths = { 'the': 3, 'a': 1, 'is': 2, 'an': 3 }
>>> lengths
{'a': 1, 'the': 3, 'is': 2, 'an': 3}
```

Notice that when the dictionary is printed, the association between **KEYS** (the words) and their **VALUES** (the lengths) is preserved, but the order of presentation is not. Dictionaries are important for the association, not the ordering. (That's what lists are for.)

The value for a given key can be recovered using the indexing notation we've already used, but now we're indexing not by numeric positions but by keys to retrieve the corresponding values.

```
>>> lengths['the']
3
>>> lengths['an']
3
>>> lengths['some']
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
KeyError: 'some'
```

You may have noticed a problem with the `lengths` list: One of the values is wrong. That's what you get when building things extensionally. Better to build the dictionary intensionally. First, we build a list of pairs of words and their lengths using a list comprehension.

```
>>> lenList = [ (word, len(word)) for word in ['the', 'a', 'is', 'an'] ]
>>> lenList
[('the', 3), ('a', 1), ('is', 2), ('an', 2)]
```

Then we convert this list of pairs into a dictionary using the `dict` function.

```
>>> lenDict = dict(lenList)
>>> lenDict['the']
3
>>> lenDict['an']
2
>>> lenDict['some']
Traceback (most recent call last):
```

```
File "<stdin>", line 1, in <module>
KeyError: 'some'
```

**Exercise 41.** Build a dictionary named `firstLetters` of words and their first letters. The word types should be taken from the Stein poem. The result should look like this:

```
>>> firstLetters
{'be': 'b', 'sure': 's', 'is': 'i', 'when': 'w', 'it': 'i', 'mine': 'm', 'to': 't', 'let': 'l'}
```

□

There are a few additional functions for manipulating dictionaries that may prove useful. The `keys` function returns a list of all of the keys defined in a dictionary

KEYS FUNCTION

```
>>> firstLetters.keys()
['be', 'sure', 'is', 'when', 'it', 'mine', 'to', 'let']
```

VALUES FUNCTION

and the `values` function returns a list of all of the values in a dictionary.

```
>>> firstLetters.values()
['b', 's', 'i', 'w', 'i', 'm', 't', 'l']
```

ITEMS FUNCTION

Finally, the `items` function returns a list of key-value pairs from the dictionary.

```
>>> firstLetters.items()
[('be', 'b'), ('sure', 's'), ('is', 'i'), ('when', 'w'), ('it', 'i'), ('mine', 'm'), ('to', 't'),
```

## 18. LOOPS AND CONDITIONALS

It's now page 35, and I've postponed as long as possible a discussion of the kind of control structures that many people think of as the hallmark of computer programming, such constructs as loops and conditionals. The style of programming I've been implicitly using – a kind of functional programming over compound data structures – eschews these kinds of structures. But for the next steps, we'll need to use them a bit.

The `FOR LOOP` allows executing a block of code several times, once *for* each value that a certain variable takes on. The notation is as follows: `FOR LOOP`

```
for <variable> in <list or set or other iterable data>:
    <body>
```

For example,

```
>>> for letter in 'sanctus':
...     print letter
...
s
a
n
c
t
u
s
```

Note the `INDENTATION`. It is crucial. Python uses indentation to convey the structure of the program. What constitutes the body of a for loop, for instance, is exactly the sequence of textual lines that follow the first line and that are *indented more deeply*. Similarly for other constructs in the language. Indentation is important; pay attention to it. `INDENTATION`

The `print` command (it's not a function) used above, when executed, has the side effect of presenting the printed representation of the comma-separated items following it (they're not really arguments) to the screen. I've used it inside the loop so that we can see what's happening inside the loop. `PRINT COMMAND`

The `CONDITIONAL` allows different code to be executed depending on whether a particular condition holds or not. We test the condition, and if it holds execute one branch of the conditional, otherwise executing the other branch. `CONDITIONAL`

```
if <condition>:
    <true branch>
else:
    <false branch>
```

☞ The **else:** and *(else branch)* can be dropped if nothing needs to be done in case the condition is false.

Here's an (admittedly artificial) example:

```
>>> occurs = {}
>>> for letter in 'sanctus':
...     if letter in 'Germanus':
...         occurs[letter] = True
...     else:
...         occurs[letter] = False
...
>>> occurs
{'a': True, 'c': False, 'n': True, 's': True, 'u': True, 't': False}
```

**Exercise 42.** *What does this snippet of code do?* □

What kinds of expressions can be in the test part of a conditional? Any expression whose value is a truth value, or **BOOLEAN**. The Boolean data type contains just two values: **True** and **False**. (In the above snippet, the values in the dictionary were also Booleans.) There are several functions that return Boolean values. Here are just a few:

- IN FUNCTION** • ***x in y***: Returns **True** just in case the value *x* is one of the values in the list, set, or other iterable data object *y*. Otherwise, it returns **False**.
- == FUNCTION** • ***x == y***: Returns **True** just in case *x* and *y* are the same value.
- < FUNCTION** • ***x < y***: Returns **True** just in case the value *x* is less than the value *y* under whatever ordering is appropriate for their data type (numerically for numbers, lexicographically for strings).
- AND FUNCTION** • ***x and y***: Returns **True** just in case both *x* and *y* have the value **True**.

There are many other built-in functions that return Booleans, and of course you can define your own.

**Exercise 43.** *Define and test a function `isPalindrome` that returns a Boolean: **True** if its argument is a palindromic string, and **False** otherwise.* □

**Exercise 44.** *Define and test a function `printPalindromes` that prints all of the words in its list argument that are palindromes, one palindrome per line.* □

**Exercise 45.** *Define and test a function `commonLetters` that takes two string arguments and returns a string containing all of the letters that its two arguments have in common. Demonstrate it on the two strings `'disproportionableness'` and `'absolutism'`. Hint: The answer is `'isotabl'`.* □

A useful idiom is to loop over all of the key-value pairs in a dictionary by taking advantage of the fact that the `items` function returns an iterable list:

```
>>> for (key, value) in firstLetters.items():
...     print key, "has first letter", value
...
be has first letter b
sure has first letter s
is has first letter i
when has first letter w
it has first letter i
mine has first letter m
to has first letter t
let has first letter l
```

## 19. A CONCORDANCE

In this section, you'll put together code to generate a simple keyword-in-context (KWIC) concordance, which lists for each word in a text all of the contexts in which it occurs.

Recall the dictionary you built in Exercise 41. This dictionary associates each word with its first letter. Of course, in a traditional dictionary (in the nontechnical sense of the word 'dictionary'), the association is the other way around: Each letter is associated with a list of the words that it is the first letter of. We could generate such a dictionary from the one we already built if we had a way of "inverting" dictionaries. Such a dictionary inverter will turn out to be useful for other tasks as well.

**Problem 46.** Write a function that takes a dictionary as its argument and returns a new dictionary that is the "inversion" of its argument. The keys in the new dictionary are the values in the original, and the values for a key  $x$  is the list of all keys in the original whose value in the original was  $x$ . □

If you've done this problem properly, you should get the following behavior:

```
>>> pprint(invertDictionary(firstLetters))
{'b': ['be'],
 'i': ['is', 'it'],
 'l': ['let'],
 'm': ['mine'],
 's': ['sure'],
 't': ['to'],
 'w': ['when']}
```

We've turned our `firstLetters` dictionary into a dictionary in the conventional sense, a mapping from letters to words they start with.

Now a slightly more sophisticated case.

**Problem 47.** Generate a dictionary that for a given list of words (the word in the Stein poem, say) associates each position or index with the word at that index. The dictionary should associate the number 0 with 'let' (because the Stein poem has the word 'let' at index 0), the number 1 with 'it', and so forth. Then invert the dictionary. The inverted dictionary will map words to a list of positions where that word occurs – a concordance! □

Finally, we can keep track not only of the index of each word, but also its context, the few words surrounding it.

**Problem 48.** Choose an appropriate dictionary structure that, when inverted, associates with each word a list of pairs. Each pair has an index and a surrounding  $n$ -gram at that position. Create such a dictionary and invert it. Write some code to print out the contents

*of that dictionary in a nice format. The output of such a concordance generator operating on the Stein poem can be found in Appendix A.* □

## APPENDIX A. A CONCORDANCE WITH DESCRIPTION

```
>>> printConcordance(concordance)
```

```
be:
```

```
91 -- let it be to be
139 -- mine to be sure to
210 -- when to be sure when
136 -- sure to be mine to
57 -- mine to be sure let
72 -- be to be sure let
201 -- be to be sure to
2 -- let it be when it
62 -- it to be sure when
114 -- be to be sure let
121 -- mine to be sure let
108 -- be to be sure let
125 -- let it be to be
93 -- be to be sure to
25 -- be to be sure when
152 -- let it be to be
81 -- sure to be sure let
100 -- let it be to be
160 -- mine to be sure let
194 -- let it be mine to
190 -- mine to be sure let
70 -- let it be to be
148 -- be to be mine let
50 -- be to be sure when
76 -- let it be to be
170 -- let it be to be
78 -- be to be sure to
87 -- be to be sure let
214 -- when to be sure to
32 -- mine to be sure let
176 -- let it be to be
199 -- let it be to be
172 -- be to be sure let
118 -- let it be mine to
39 -- let it be let it
146 -- let it be to be
133 -- mine to be sure to
96 -- sure to be sure let
130 -- sure to be mine to
```

```
42 -- let it be to be
  8 -- mine to be sure let
12 -- let it be when it
154 -- be to be sure to
157 -- sure to be mine to
127 -- be to be sure to
178 -- be to be sure to
 48 -- let it be to be
102 -- be to be sure let
142 -- sure to be mine let
106 -- let it be to be
 44 -- be to be sure let
186 -- it to be sure mine
166 -- be to be mine let
 36 -- let it be let it
181 -- sure to be sure let
164 -- let it be to be
206 -- let it be mine when
112 -- let it be to be
 85 -- let it be to be
 23 -- let it be to be
```

sure:

```
211 -- to be sure when to
 33 -- to be sure let it
115 -- to be sure let it
 73 -- to be sure let it
  9 -- to be sure let it
128 -- to be sure to be
103 -- to be sure let it
 97 -- to be sure let it
 82 -- to be sure let it
134 -- to be sure to be
179 -- to be sure to be
122 -- to be sure let it
182 -- to be sure let it
 58 -- to be sure let it
215 -- to be sure to let
 88 -- to be sure let it
187 -- to be sure mine to
 94 -- to be sure to be
 63 -- to be sure when it
140 -- to be sure to be
```

79 -- to be sure to be  
161 -- to be sure let it  
26 -- to be sure when it  
109 -- to be sure let it  
173 -- to be sure let it  
51 -- to be sure when it  
155 -- to be sure to be  
191 -- to be sure let it  
45 -- to be sure let it  
202 -- to be sure to let

is:

66 -- when it is mine let  
19 -- when it is mine let  
54 -- when it is mine to  
5 -- when it is mine to  
29 -- when it is mine to  
15 -- when it is mine when

when:

208 -- be mine when to be  
212 -- be sure when to be  
13 -- it be when it is  
17 -- is mine when it is  
3 -- it be when it is  
64 -- be sure when it is  
27 -- be sure when it is  
52 -- be sure when it is

it:

14 -- be when it is mine  
47 -- sure let it be to  
99 -- sure let it be to  
69 -- mine let it be to  
60 -- sure let it to be  
184 -- sure let it to be  
90 -- sure let it be to  
175 -- sure let it be to  
193 -- sure let it be mine  
38 -- be let it be let  
11 -- sure let it be when  
205 -- to let it be mine  
18 -- mine when it is mine  
75 -- sure let it be to  
105 -- sure let it be to

```
111 -- sure let it be to
163 -- sure let it be to
169 -- mine let it be to
 35 -- sure let it be let
 84 -- sure let it be to
  4 -- be when it is mine
117 -- sure let it be mine
198 -- to let it be to
 28 -- sure when it is mine
 41 -- be let it be to
124 -- sure let it be to
 65 -- sure when it is mine
151 -- mine let it be to
 22 -- mine let it be to
218 -- to let it to be
 53 -- sure when it is mine
145 -- mine let it be to
mine:
 67 -- it is mine let it
137 -- to be mine to be
188 -- be sure mine to be
119 -- it be mine to be
 30 -- it is mine to be
131 -- to be mine to be
  6 -- it is mine to be
167 -- to be mine let it
195 -- it be mine to let
158 -- to be mine to be
207 -- it be mine when to
 55 -- it is mine to be
143 -- to be mine let it
 16 -- it is mine when it
 20 -- it is mine let it
149 -- to be mine let it
to:
141 -- be sure to be mine
159 -- be mine to be sure
135 -- be sure to be mine
 31 -- is mine to be sure
171 -- it be to be sure
 61 -- let it to be sure
129 -- be sure to be mine
```

189 -- sure mine to be sure  
95 -- be sure to be sure  
138 -- be mine to be sure  
56 -- is mine to be sure  
24 -- it be to be sure  
209 -- mine when to be sure  
203 -- be sure to let it  
49 -- it be to be sure  
156 -- be sure to be mine  
86 -- it be to be sure  
113 -- it be to be sure  
71 -- it be to be sure  
165 -- it be to be mine  
7 -- is mine to be sure  
147 -- it be to be mine  
177 -- it be to be sure  
92 -- it be to be sure  
185 -- let it to be sure  
200 -- it be to be sure  
180 -- be sure to be sure  
132 -- be mine to be sure  
216 -- be sure to let it  
126 -- it be to be sure  
101 -- it be to be sure  
213 -- sure when to be sure  
43 -- it be to be sure  
80 -- be sure to be sure  
120 -- be mine to be sure  
77 -- it be to be sure  
196 -- be mine to let it  
153 -- it be to be sure  
107 -- it be to be sure

let:

10 -- be sure let it be  
168 -- be mine let it be  
40 -- it be let it be  
150 -- be mine let it be  
34 -- be sure let it be  
110 -- be sure let it be  
197 -- mine to let it be  
37 -- it be let it be  
144 -- be mine let it be

```
83 -- be sure let it be
116 -- be sure let it be
74 -- be sure let it be
46 -- be sure let it be
174 -- be sure let it be
68 -- is mine let it be
123 -- be sure let it be
204 -- sure to let it be
104 -- be sure let it be
21 -- is mine let it be
183 -- be sure let it to
59 -- be sure let it to
89 -- be sure let it be
162 -- be sure let it be
217 -- sure to let it to
98 -- be sure let it be
192 -- be sure let it be
```

## APPENDIX B. STATISTICS

Running time of included Python examples: 88.2 seconds.

## REFERENCES

Gertrude Stein. *An Acquaintance with Description*. Seizin Press, 1929. URL <http://books.google.com/books?id=YpFuQgAACAAJ>.

## INDEX

- n*-grams, 30
- < function, 36
- = operator, *see* assignment
- == function, 36
  
- addition, 14
- and function, 36
- arguments, 20
- arithmetic operators, 14
- assignment, 15
  
- Boolean, 36
  
- clickable links, 4
- command
  - def, 31
  - for, 35
  - if, 35
  - print, 35
  - return, 31
- commands, 14
- conditional, 35
- count function, 21
  
- data type, 16
- def command, 31
- delimiters, 17
- dict function, 33
- dictionary, 33
  
- escape sequence, 25
- exercises, 4
- expressions, 14
- extensional, 27
  
- False value, 36
- first principles, 2
- for loop, 35
- four fours, 14
- function
  - <, 36
  - ==, 36
  - and, 36
  - count, 21
  - dict, 33
  - in, 36
  - items, 34
  - join, 26
  - keys, 34
  - len, 17, 20
  - lower, 26
  - pprint, 23
  - range, 21
  - set, 29
  - sorted, 21
  - split, 26
  - strip, 27
  - values, 34
- function call, 31
- functions, 20
  
- hashable, 29
  
- if command, 35
- imitation, 2
- in function, 36
- indentation, 35
- indexing, 16
- intensional, 27
- interpreter, 14
- items function, 34
  
- join function, 26
  
- key concepts, 4
- keys, 33
- keys function, 34
  
- len function, 17, 20
- list, 16
- list comprehensions, 27
- lower function, 26
  
- mathematical notation, 20
- multiplication, 14
  
- newline, 17
  
- object notation, 20
- optional arguments, 9
  
- position, 16
- pprint function, 23

print command, 35  
problems, 4

range function, 21  
result, 20  
return command, 31

set, 29  
set function, 29  
slicing, 18  
sorted function, 21  
split function, 26  
Stein, Gertrude, 22  
string, 17  
strip function, 27

tokens, 22  
True value, 36  
tuple, 19  
types, 22

union, 29

values, 33  
values function, 34  
variables, 15

whitespace, 7

