

Programming For Biology 2017

Instructors

Simon Prochnik

Sofia Robb

Table of Contents

- Big Picture
 - Why?
 - Helpful Tips
- Unix
 - Unix 1
 - * Unix Overview
 - What is the Command-Line?
 - * The Basics
 - Logging into Your Workstation
 - Bringing up the Command-Line
 - OK. I've Logged in. What Now?
 - Command-Line Prompt
 - Issuing Commands
 - Command-Line Editing
 - Wildcards
 - Home Sweet Home
 - Getting Around
 - Essential Unix Commands
 - Getting Information About Commands
 - Finding Out What Commands are on Your Computer
 - Arguments and Command Switches
 - Spaces and Funny Characters
 - Useful Commands
 - Manipulating Directories
 - Networking
 - Standard I/O and Command Redirection
 - A Simple Example
 - Redirection Meta-Characters
 - Filters, Filenames and Standard Input
 - Standard I/O and Pipes
 - More Pipe Idioms
 - Uniquify Lines in a File
 - Concatenate Several Lists and Remove Duplicates
 - Count Unique Lines in a File
 - Page Through a Really Long Directory Listing
 - Monitor a Rapidly Growing File for a Pattern
 - * Advanced Unix
 - * Link to Unix 1 Problem Set
 - Unix 2
 - * Text Editors
 - Git for Beginners
 - * The Big Picture.
 - Collaboration
 - Storing Versions
 - Restoring Previous Versions

- Backup
 - The Details
 - * The Basics
 - Creating a new repository
 - Cloning a Repository
 - Links to Slightly less basic topics
 - * Link To Unix 2 Problem Set
- Python
 - Python 1
 - * Python Overview
 - * Running Python
 - Interactive Interpreter
 - Running Python Scripts
 - Python Script
 - * Syntax
 - Python Identifiers:
 - Naming conventions for Python Identifiers:
 - Reserved Words:
 - Lines and Indentation:
 - Comments:
 - Blank Lines:
 - Python Options:
 - * Data Types and Variables
 - Numbers and Strings
 - Lists
 - Command line parameters: A Special Built-in List
 - Tuple
 - Dictionary
 - Type Conversion
 - * Link to Python 1 Problem Set
 - Python 2
 - * Operators
 - Arithmetic Operators
 - Assignment Operators
 - Comparison Operators
 - Logical Operators
 - Membership Operators
 - Operator Precedence
 - * Truth
 - Use bool() to test for truth
 - * Logic: Control Statements
 - If Statement
 - If/Else Statement
 - if/elif
 - * Numbers
 - integer
 - floating point number
 - complex number
 - Conversion functions
 - Numeric Functions
 - * Comparing two numbers
 - * Link to Python 2 Problem Set
 - Python 3
 - * Sequences

- * What functions go with my object?
- * Strings
 - Quotation Marks
 - Strings and the print() function
 - Errors and Printing
 - Special/Escape Characters
 - Concatenation
 - Determine the length of a string:
 - Changing String Case
 - Find and Count
 - Find and Replace
 - Extracting a Substring, or Slicing
 - Locate and Report
 - Other String Methods
- * String Formatting
 - The format() mini-language
 - Summary of special formatting symbols so far
 - What's the point?
- * Lists and Tuples
 - Lists
 - Accessing Values in Lists
 - Changing Values in a List
 - Extracting a Subset of a List, or Slicing
 - List Operators
 - List Functions
 - List Methods
 - Building a List one Value at a Time
- * [Link to Python 3 Problem Set](#)
- Python 4
 - * Loops
 - while loop
 - While Loop Syntax
 - While/Else
 - For Loops
 - For Loop Syntax
 - For/Else
 - Loop Control
 - Loop Control: Break
 - Loop Control: Continue
 - Iterators
 - List Comprehension
 - * Dictionaries
 - Creating a Dictionary
 - Accessing Values in Dictionaries
 - Changing Values in a Dictionary
 - Building a Dictionary one Key/Value at a Time
 - Checking That Dictionary Keys Exist
 - Sorting Dictionary Keys
 - Dictionary Functions
 - Dictionary Methods
 - * Sets
 - Set Operators
 - Set Functions
 - Set Methods

- * [Link to Python 4 Problem Set](#)
- Python 5
 - * Iterables
 - * I/O and Files
 - Writing to the Screen
 - Reading input from the keyboard
 - Reading from a File
 - Open a File
 - Reading the contents of a file
 - Writing to a File
 - Building a Dictionary from a File
 - * [Link to Python 5 Problem Set](#)
- Python 6
 - * Regular Expressions
 - Individual Characters
 - Character Classes
 - Anchors
 - Quantifiers
 - Variables and Patterns
 - Either Or
 - Subpatterns
 - Using Subpatterns Inside the Regular Expression Match
 - Subpatterns and Greediness
 - Using Subpatterns Outside the Regular Expression Match
 - Practical Example: Codons
 - Truth and Regular Expression Matches
 - Using Regular expressions in substitutions
 - Using subpatterns in the replacement
 - * [Link to Python 6 Problem Set](#)
- Python 7
 - * Functions
 - Creating/Defining a Function to Find AT Content:
 - Using/Running/Calling Your function:
 - The details
 - Naming Arguments
 - Keyword Arguments
 - Default Values for Arguments
 - lambda
 - * Scope
 - Local Variables
 - Global
 - * Namespaces
 - * Modules
 - os.path
 - os.system
 - subprocess
 - Capturing output from a shell pipeline
 - Capturing output the long way (for a single command)
 - sys
 - re
 - collections
 - copy
 - math
 - random

- statistics
 - glob
 - argparse
 - Many more modules in many areas
 - * Link to Python 7 Problem Set
 - Python 8
 - * Exception Handling
 - try/except/else/finally
 - Getting more information about an exception
 - Raising an Exception
 - * Datastructures
 - Two-dimensional lists
 - Lists of dictionaries
 - Dictionaries of lists
 - Dictionaries of dictionaries
 - List Comprehensions
 - More complex data structures
 - * Link to Python 8 Problem Set
 - Python 9
 - * BioPython
 - BioPython Overview
 - BioPython Subtopic 1
 - BioPython Subtopic 2
 - Bioinformatic Analysis and Tools
 - Bioinformatic Analysis and Tools Overview
 - Sequence Search and Alignments
 - Assembly
 - * DNA
 - * RNA
 - NGS
 - Ontology
-

Big Picture

Why?

*Why is it important for **Biologists** to learn to program?*

You probably already know the answer to this question since you are here.

We firmly believe that knowing how to program is just as essential as knowing how to run a gel or set up a PCR reaction. The data we can now get from a single experiment can be overwhelming. This data often needs to be reformatted, filtered, and analyzed in unique ways. Programming allows you to perform these tasks in a **reproducible**, **timely**, **efficient**, and **thoughtful** way.

Helpful Tips

What are our tips for being successful in your efforts to learn to program?

1. Practice, practice, practice. Please, spend as much time possible actually coding.

2. Write only a line or two of code, then test it. If you write too many lines, it becomes more difficult to debug if there is an error.
 3. Errors are not failures. Every error you get is a learning opportunity. Every single error you debug is a major success. Fixing errors is how you will cement what you have learned.
 4. Don't spend too much time trying to figure out a problem. While it's a great learning experience to try to solve an issue on your own, it's not fun getting frustrated or spending a lot of time stuck. We are here to help you, so please ask us whenever you need help.
 5. Lectures are important, but the practice is more important.
 6. Review sessions are important, but practice is more important.
 7. It is essential that we help you to learn how to find solutions on your own.
-

Unix

Unix 1

Unix Overview

What is the Command-Line?

Underlying the pretty Mac OSX GUI is a powerful command-line operating system. The command-line gives you access to the internals of the OS, and is also a convenient way to write custom software and scripts.

Many bioinformatics tools are written to run on the command-line and have no graphical interface. In many cases, a command-line tool is more versatile than a graphical tool, because you can easily combine command-line tools into automated scripts that accomplish tasks without human intervention.

In this course, we will be writing Python scripts that are completely command-line based.

The Basics

Logging into Your Workstation

Your workstation is an iMac. To log into it, provide the following information:

Your username: the initial of your first name, followed by your full last name.

For example, if your username is **srobb** for sofia **robb**

Your password: **changeme**

Bringing up the Command-Line

To bring up the command-line, use the Finder to navigate to *Applications->Utilities* and double-click on the *Terminal* application. This will bring up a window like the following:

You can open several Terminal windows at once. This is often helpful.

You will be using this application a lot, so I suggest that you drag the Terminal icon into the shortcuts bar at the bottom of your screen.

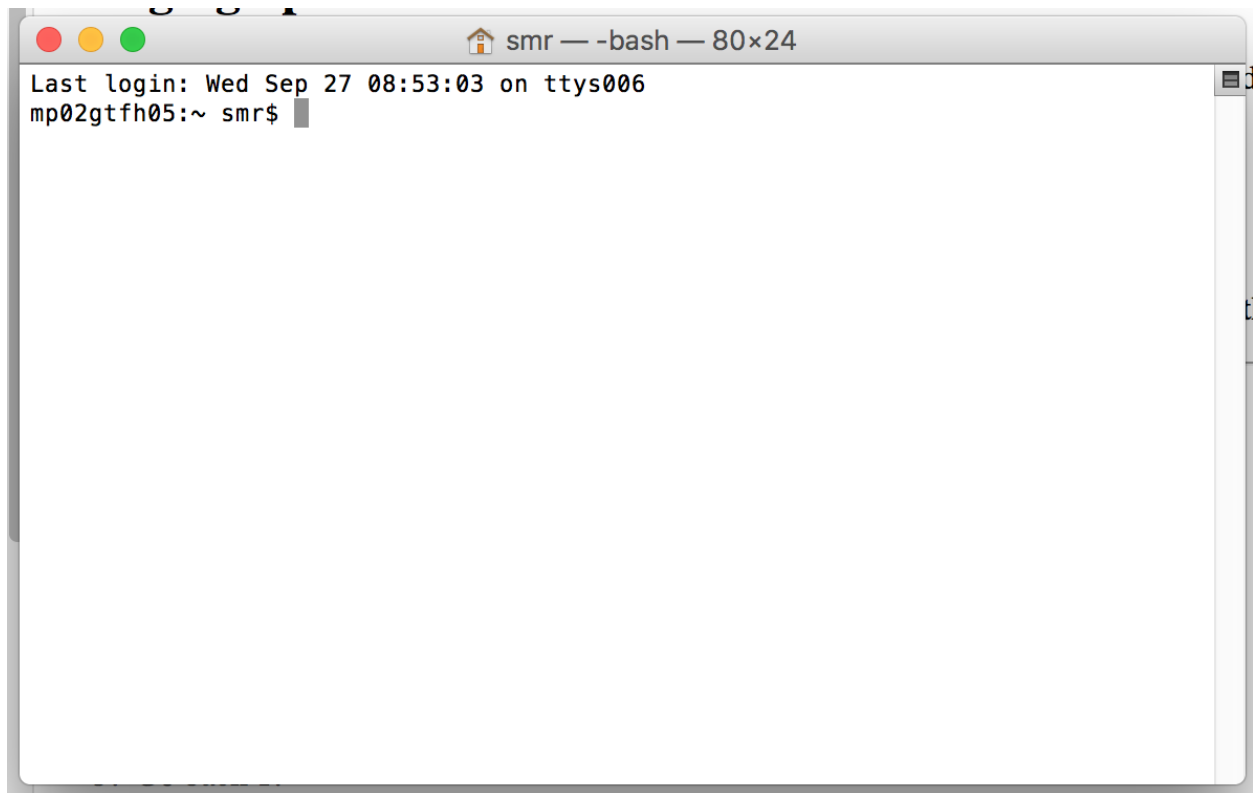


Figure 1: OSX Terminal

OK. I've Logged in. What Now?

The terminal window is running **shell** called “bash.” The shell is a loop that: 1. Prints a prompt 2. Reads a line of input from the keyboard 3. Parses the line into one or more commands 4. Executes the commands (which usually print some output to the terminal) 5. Go back 1.

There are many different shells with bizarre names like **bash**, **sh**, **csch**, **tcsh**, **ksh**, and **zsh**. The “sh” part means shell. Each shell was designed for the purpose of confusing you and tripping you up. We have set up your accounts to use **bash**. Stay with **bash** and you'll get used to it, eventually.

Command-Line Prompt

Most of bioinformatics is done with command-line software, so you should take some time to learn to use the shell effectively.

This is a command-line prompt:

```
1 bush202>
```

This is another:

```
1 (~) 51%
```

This is another:

```
1 srobb@bush202 1:12PM>
```

What you get depends on how the system administrator has customized your login. You can customize it yourself when you know how.

The prompt tells you the shell is ready to accept a command. When a long-running command is going, the prompt will not reappear until the system is ready to deal with your next request.

Issuing Commands

Type in a command and press the <Enter> key. If the command has output, it will appear on the screen. Example:

```
1 (~) 53% ls -F
2 GNUstep/          cool_elegans.movies.txt  man/
3 INBOX             docs/                   mtv/
4 INBOX~            etc/                   nsmail/
5 Mail@             games/                  pcod/
6 News/             get_this_book.txt      projects/
7 axhome/           jcod/                  public_html/
8 bin/              lib/                   src/
9 build/            linux/                  tmp/
10 ccod/
11 (~) 54%
```

The command here is `ls -F` which produces a listing of files and directories in the current directory (more on which later). After its output, the command prompt appears again.

Some programs will take a long time to run. After you issue their command names, you won't recover the shell prompt until they're done. You can either launch a new shell (from Terminal's File menu), or run the command in the background using the ampersand:

```
1 (~) 54% long_running_application &
2 (~) 55%
```

The command will now run in the background until it is finished. If it has any output, the output will be printed to the terminal window. You may wish to redirect the output as described later.

Command-Line Editing

Most shells offer command-line editing. Up until the comment you press `,` you can go back over the command-line and edit it using the keyboard. Here are the most useful keystrokes:

- *Backspace*: Delete the previous character and back up one.
- *Left arrow*, right arrow: Move the text insertion point (cursor) one character to the left or right.

- *control-a* (^a): Move the cursor to the beginning of the line. (Mnemonic: A is first letter of alphabet)
- *control-e* (^e): Move the cursor to the end of the line. Mnemonic: E for the End (^Z was already taken for something else).
- *control-d* (^d): Delete the character currently under the cursor. D=Delete.
- *control-k* (^k): Delete the entire line from the cursor to the end. k=kill. The line isn't actually deleted, but put into a temporary holding place called the "kill buffer".
- *control-y* (^y): Paste the contents of the kill buffer onto the command-line starting at the cursor. y=yank.
- *Up arrow, down arrow*: Move up and down in the command history. This lets you reissue previous commands, possibly after modifying them.

There are also some useful shell commands you can issue:

- *history*: Show all the commands that you have issued recently, nicely numbered.
- *!**<number>***: Reissue an old command, based on its number (which you can get from *history*).
- *!!*: Reissue the immediate previous command.
- *!**<partial command string>***: Reissue the previous command that began with the indicated letters. For example, *!!* would reissue the *ls -F* command from the earlier example.

bash offers automatic command completion and spelling correction. If you type part of a command and then the tab key, it will prompt you with all the possible completions of the command. For example:

```
1 (~) 51% <b>fd<tab></b>
2 (~) 51% fd
3 fd2ps      fdesign  fdformat fdlist      fdmount  fdmountd fdrawcmd
      fdumount
4 (~) 51%
```

If you hit tab after typing a command, but before pressing <Enter>, **bash** will prompt you with a list of file names. This is because many commands operate on files.

Wildcards

You can use wildcards when referring to files. "*" refers to zero or more characters. "?" refers to any single character. For example, to list all files with the extension ".txt", run **ls** with the pattern "*.txt":

```
1 (~) 56% ls -F *.txt
2 final_exam_questions.txt  genomics_problem.txt
3 genebridge.txt             mapping_run.txt
```

There are several more advanced types of wildcard patterns that you can read about in the **tcsh** manual page. For example, you can refer to files beginning with the characters "f" or "g" and ending with ".txt" this way:

```
1 (~) 57% ls -F [f-g]*.txt
2 final_exam_questions.txt  genebridge.txt
      genomics_problem.txt
```

Home Sweet Home

When you first log in, you'll be placed in a part of the system that is your personal domain, called the *home directory*. You are free to do with this area what you will: in particular you can create and delete files and other directories. In general, you cannot create files elsewhere in the system.

Your home directory lives somewhere way down deep in the bowels of the system. On our iMacs, it is a directory with the same name as your login name, located in `/Users`. The full directory path is therefore `/Users/username`. Since this is a pain to write, the shell allows you to abbreviate it as `~username` (where “username” is your user name), or simply as `~`. The weird character (technically called the “tilde” or “twiddle”) is usually hidden at the upper left corner of your keyboard.

To see what is in your home directory, issue the command `ls -F`:

```
1 (~) % ls -F
2 INBOX          Mail/          News/          nsmail/        public_html/
```

This shows one file “INBOX” and four directories (“Mail”, “News”) and so on. (The “-F” in the command turns on fancy mode, which appends special characters to directory listings to tell you more about what you’re seeing. “/” means directory.)

In addition to the files and directories shown with `ls -F`, there may be one or more hidden files. These are files and directories whose names start with a “.” (technically called the “dot” character). To see these hidden files, add an “a” to the options sent to the `ls` command:

```
1 (~) % ls -aF
2 ./          .cshrc       .login       Mail/
3 ../         .fetchhost   .netscape/  News/
4 .Xauthority .fvwmrc      .xinitrc*    nsmail/
5 .Xdefaults  .history     .xsession@   public_html/
6 .bash_profile .less        .xsession-errors
7 .bashrc     .lessrc      INBOX
```

Whoa! There’s a lot of hidden stuff there. But don’t go deleting dot files willy-nilly. Many of them are essential configuration files for commands and other programs. For example, the *.profile* file contains configuration information for the **bash** shell. You can peek into it and see all of **bash**’s many options. You can edit it (when you know what you’re doing) in order to change things like the command prompt and command search path.

Getting Around

You can move around from directory to directory using the `cd` command. Give the name of the directory you want to move to, or give no name to move back to your home directory. Use the `pwd` command to see where you are (or rely on the prompt, if configured):

```

1 (~/docs/grad_course/i) 56% cd
2 (~) 57% cd /
3 (/) 58% ls -F
4 bin/          dosc/          gmon.out       mnt/          sbin/
5 boot/         etc/          home@         net/          tmp/
6 cdrom/        fastboot      lib/          proc/         usr/
7 dev/          floppy/       lost+found/    root/         var/
8 (/) 59% <b>cd ~/docs/</b>
9 (~/docs) 60% <b>pwd</b>
10 /usr/home/lstein/docs
11 (~/docs) 62% cd ../projects/
12 (~/projects) 63% ls
13 Ace-browser/      bass.patch
14 Ace-perl/         cgi/
15 Foo/             cgi3/
16 Interface/        computertalk/
17 Net-Interface-0.02/ crypt-cbc.patch
18 Net-Interface-0.02.tar.gz fixer/
19 Pts/             fixer.tcsh
20 Pts.bak/         introspect.pl*
21 PubMed/         introspection.pm
22 SNPdb/          rhmap/
23 Tie-DBI/         sbbox/
24 ace/            sbbox-1.00/
25 atir/           sbbox-1.00.tgz
26 bass-1.30a/      zhmapper.tar.gz
27 bass-1.30a.tar.gz
28 (~/projects) 64%

```

Each directory contains two special hidden directories named “.” and “..”. “.” refers always to the directory in which it is located. “..” refers always to the parent of the directory. This lets you move upward in the directory hierarchy like this:

```
1 (~/docs) 64% cd ..
```

and to do arbitrarily weird things like this:

```
1 (~/docs) 65% cd ../../lstein/docs
```

The latter command moves upward two levels, and then into a directory named “docs”.

If you get lost, the *pwd* command prints out the full path to the current directory:

```
1 (~) 56% pwd
2 /Users/lstein
```

Essential Unix Commands

With the exception of a few commands that are built directly into the shell, all Unix commands are standalone executable programs. When you type the name of a command, the shell will search through all the directories listed in the PATH environment variable for an executable of the same name. If found, the shell will execute the command. Otherwise, it will give a “command not found” error.

Most commands live in /bin, /usr/bin, or /usr/local/bin.

Getting Information About Commands

The **man** command will give a brief synopsis of the command:

```
1 (~) 76% man wc
2 Formatting page, please wait...
3 WC(1) WC(1)
4
5 NAME
6     wc - print the number of bytes, words, and lines in files
7
8 SYNOPSIS
9     wc [-clw] [--bytes] [--chars] [--lines] [--words] [--help]
10     [--version] [file...]
11
12 DESCRIPTION
13     This manual page documents the GNU version of wc.  wc
14     counts the number of bytes, whitespace-separated words,
15     ...
```

Finding Out What Commands are on Your Computer

The **apropos** command will search for commands matching a keyword or phrase:

```
1 (~) 100% apropos column
2 showtable (1) - Show data in nicely formatted columns
3 colrm (1) - remove columns from a file
4 column (1) - columnate lists
5 fix132x43 (1) - fix problems with certain (132 column)
   graphics
6 modes
```

Arguments and Command Switches

Many commands take arguments. Arguments are often (but not inevitably) the names of one or more files to operate on. Most commands also take command-line “switches” or “options”, which fine-tune what the command does. Some commands recognize “short switches” that consist of a single character, while others recognize “long switches” consisting of whole words.

The **wc** (word count) program is an example of a command that recognizes both long and short options. You can pass it the **-c**, **-w** and/or **-l** options to count the characters, words and lines in a text file, respectively. Or you can use the longer but more readable, **-chars**, **-words** or **-lines** options. Both these examples count the number of characters and lines in the text file `/var/log/messages`:

```
1 (~) 102% wc -c -l /var/log/messages
2      23      941 /var/log/messages
3 (~) 103% wc --chars --lines /var/log/messages
4      23      941 /var/log/messages
```

You can cluster short switches by concatenating them together, as shown in this example:

```
1 (~) 104% wc -cl /var/log/messages
2      23      941 /var/log/messages
```

Many commands will give a brief usage summary when you call them with the **-h** or **-help** switch.

Spaces and Funny Characters

The shell uses whitespace (spaces, tabs and other nonprinting characters) to separate arguments. If you want to embed whitespace in an argument, put single quotes around it. For example:

```
1 mail -s 'An important message' 'Bob Ghost <bob@ghost.org>'
```

This will send an e-mail to the fictitious person Bob Ghost. The **-s** switch takes an argument, which is the subject line for the e-mail. Because the desired subject contains spaces, it has to have quotes around it. Likewise, my e-mail address, which contains embedded spaces, must also be quoted in this way.

Certain special non-printing characters have *escape codes* associated with them:

Escape Code	Description
<code>\n</code>	new line character
<code>\t</code>	tab character
<code>\r</code>	carriage return character
<code>\a</code>	bell character (ding! ding!)
<code>\nnn</code>	the character whose ASCII code is nnn

Useful Commands

Here are some commands that are used extremely frequently. Use **man** to learn more about them. Some of these commands may be useful for solving the problem set ;-)

Manipulating Directories

Command	Description
ls	Directory listing. Most frequently used as ls -F (decorated listing) and ls -l (long listing).
mv	Rename or move a file or directory.
cp	Copy a file.
rm	Remove (delete) a file.
mkdir	Make a directory
rmdir	Remove a directory
ln	Create a symbolic or hard link.
chmod	Change the permissions of a file or directory.

Command	Description
cat	Concatenate program. Can be used to concatenate multiple files together into a single file, or, much more frequently, to send the contents of a file to the terminal for viewing.
more	Scroll through a file page by page. Very useful when viewing large files. Works even with files that are too big to be opened by a text editor.
less	A version of more with more features.
head	View the head (top) of a file. You can control how many lines to view.
tail	View the tail (bottom) of a file. You can control how many lines to view. You can also use tail to view a growing file.
wc	Count words, lines and/or characters in one or more files.
tr	Substitute one character for another. Also useful for deleting characters.
sort	Sort the lines in a file alphabetically or numerically.
uniq	Remove duplicated lines in a file.
cut	Remove sections from each line of a file or files.
fold	Wrap each input line to fit in a specified width.
grep	Filter a file for lines matching a specified pattern. Can also be reversed to print out lines that don't match the specified pattern.
gzip (gunzi)	Compress (uncompress) a file.
tar	Archive or unarchive an entire directory into a single file.
emacs	Run the Emacs text editor (good for experts).
vi	Run the vi text editor (better for experts).

Networking

Command	Description
ssh	A secure (encrypted) way to log into machines.
scp	A secure way to copy (cp) files to and from remote machines.

Command	Description
ping	See if a remote host is up.
ftp/sftp (secure)	transfer files using the File Transfer Protocol.
who	See who else is logged in.
lp	Send a file or set of files to a printer.

Standard I/O and Command Redirection

Unix commands communicate via the command-line interface. They can print information out to the terminal for you to see, and accept input from the keyboard (that is, from *you!*)

Every Unix program starts out with three connections to the outside world. These connections are called “streams”, because they act like a stream of information (metaphorically speaking):

Stream Type	Description
standard input	This is a communications stream initially attached to the keyboard. When the program reads from standard input, it reads whatever text you type in.
standard output	This stream is initially attached to the command window. Anything the program prints to this channel appears in your terminal window.
standard error	This stream is also initially attached to the command window. It is a separate channel intended for printing error messages.

The word “initially” might lead you to think that standard input, output and error can somehow be detached from their starting places and reattached somewhere else. And you’d be right. You can attach one or more of these three streams to a file, a device, or even to another program. This sounds esoteric, but it is actually very useful.

A Simple Example

The **wc** program counts lines, characters and words in data sent to its standard input. You can use it interactively like this:

```

1 (~) 62% wc
2 Mary had a little lamb,
3 little lamb,
4 little lamb.
5
6 Mary had a little lamb,
7 whose fleece was white as snow.
8 ^D
9      6      20      107

```

In this example, I ran the **wc** program. It waited for me to type in a little poem. When I was done, I typed the END-OF-FILE character, control-D (^D for short). **wc** then printed out three numbers indicating the number of lines, words and characters in the input.

More often, you'll want to count the number of lines in a big file; say a file filled with DNA sequences. You can do this by *redirecting* **wc**'s standard input from a file. This uses the **<** metacharacter:

```
1 (~) 63% wc < big_file.fasta
2      2943      2998      419272
```

If you wanted to record these counts for posterity, you could redirect standard output as well using the **>** metacharacter:

```
1 (~) 64% wc < big_file.fasta > count.txt
```

Now if you **cat** the file *count.txt*, you'll see that the data has been recorded. **cat** works by taking its standard input and copying it to standard output. We redirect standard input from the *count.txt* file, and leave standard output at its default, attached to the terminal:

```
1 (~) 65% cat < count.txt
2      2943      2998      419272
```

Redirection Meta-Characters

Here's the complete list of redirection commands for **bash**:

Redirect command	Description
< filename	Redirect standard input to file
> filename	Redirect standard output to file
1> filename	Redirect just standard output to file (same as above)
2> filename	Redirect just standard error to file
> filename 2>&1	Redirect both stdout and stderr to file

These can be combined. For example, this command redirects standard input from the file named */etc/passwd*, writes its results into the file *search.out*, and writes its error messages (if any) into a file named *search.err*. What does it do? It searches the password file for a user named "root" and returns all lines that refer to that user.

```
1 (~) 66% grep root < /etc/passwd > search.out 2> search.err
```

Filters, Filenames and Standard Input

Many Unix commands act as filters, taking data from a file or standard input, transforming the data, and writing the results to standard output. Most filters are designed so that if they are called with one or more filenames on the command-line, they will use those files as input. Otherwise they will act on standard input. For example, these two commands are equivalent:


```
1 (~) 66% grep 'gatttgc' < big_file.fasta
2 (~) 67% grep 'gatttgc' big_file.fasta
```

Both commands use the **grep** command to search for the string “gatttgc” in the file *big_file.fasta*. The first one searches standard input, which happens to be redirected from the file. The second command is explicitly given the name of the file on the command line.

Sometimes you want a filter to act on a series of files, one of which happens to be standard input. Many filters let you use “-” on the command-line as an alias for standard input. Example:

```
1 (~) 68% grep 'gatttgc' big_file.fasta bigger_file.fasta -
```

This example searches for “gatttgc” in three places. First it looks in *big_file.fasta*, then in *bigger_file.fasta*, and lastly in standard input (which, since it isn’t redirected, will come from the keyboard).

Standard I/O and Pipes

The coolest thing about the Unix shell is its ability to chain commands together into pipelines. Here’s an example:

```
1 (~) 65% grep gatttgc big_file.fasta | wc -l
2 22
```

There are two commands here. **grep** searches a file or standard input for lines containing a particular string. Lines which contain the string are printed to standard output. **wc -l** is the familiar word count program, which counts words, lines and characters in a file or standard input. The **-l** command-line option instructs **wc** to print out just the line count. The **|** character, which is known as the “pipe” character, connects the two commands together so that the standard output of **grep** becomes the standard input of **wc**.

What does this pipe do? It prints out the number of lines in which the string “gatttgc” appears in the file *big_file.fasta*.

More Pipe Idioms

Pipes are very powerful. Here are some common command-line idioms.

Count the Number of Times a Pattern does NOT Appear in a File

The example at the top of this section showed you how to count the number of lines in which a particular string pattern appears in a file. What if you want to count the number of lines in which a pattern does **not** appear?

Simple. Reverse the test with the **grep -v** switch:

```
1 (~) 65% grep -v gatttgc big_file.fasta | wc -l
2 2921
```

Uniquify Lines in a File

If you have a long list of names in a text file, and you are concerned that there might be some duplicates, this will weed out the duplicates:

```
1 (~) 66% sort long_file.txt | uniq > unique.out
```

This works by sorting all the lines alphabetically and piping the result to the **uniq** program, which removes duplicate lines that occur together. The output is placed in a file named *unique.out*.

Concatenate Several Lists and Remove Duplicates

If you have several lists that might contain repeated entries among them, you can combine them into a single unique list by **cating** them together, then uniquifying them as before:

```
1 (~) 67% cat file1 file2 file3 file4 | sort | uniq
```

Count Unique Lines in a File

If you just want to know how many unique lines there are in the file, add a **wc** to the end of the pipe:

```
1 (~) 68% sort long_file.txt | uniq | wc -l
```

Page Through a Really Long Directory Listing

Pipe the output of **ls** to the **more** program, which shows a page at a time. If you have it, the **less** program is even better:

```
1 (~) 69% ls -l | more
```

Monitor a Rapidly Growing File for a Pattern

Pipe the output of **tail -f** (which monitors a growing file and prints out the new lines) to **grep**. For example, this will monitor the */var/log/syslog* file for the appearance of e-mails addressed to *mzhang*:

```
1 (~) 70% tail -f /var/log/syslog | grep mzhang
```

Advanced Unix

@sep: leave out details for now maybe later in the course

Here are a few more advanced Unix topics that are super useful and when you have time you should investigate further.

- awk
- sed
- perl
- for

Link to Unix 1 Problem Set

Unix 2

Text Editors

It is often necessary to create and write to a file while using the terminal. This makes it essential to use a terminal text editor. There are many text editors out there. Some of our favorite are Emacs and vim. We are going to start you out with a simple text editor call **nano**

The way you use nano to create a file is simply by typing the command *nano* followed by the name of the file you wish to create.

```
1 (~) 71% nano firstFile.txt
```

This is what you will see:

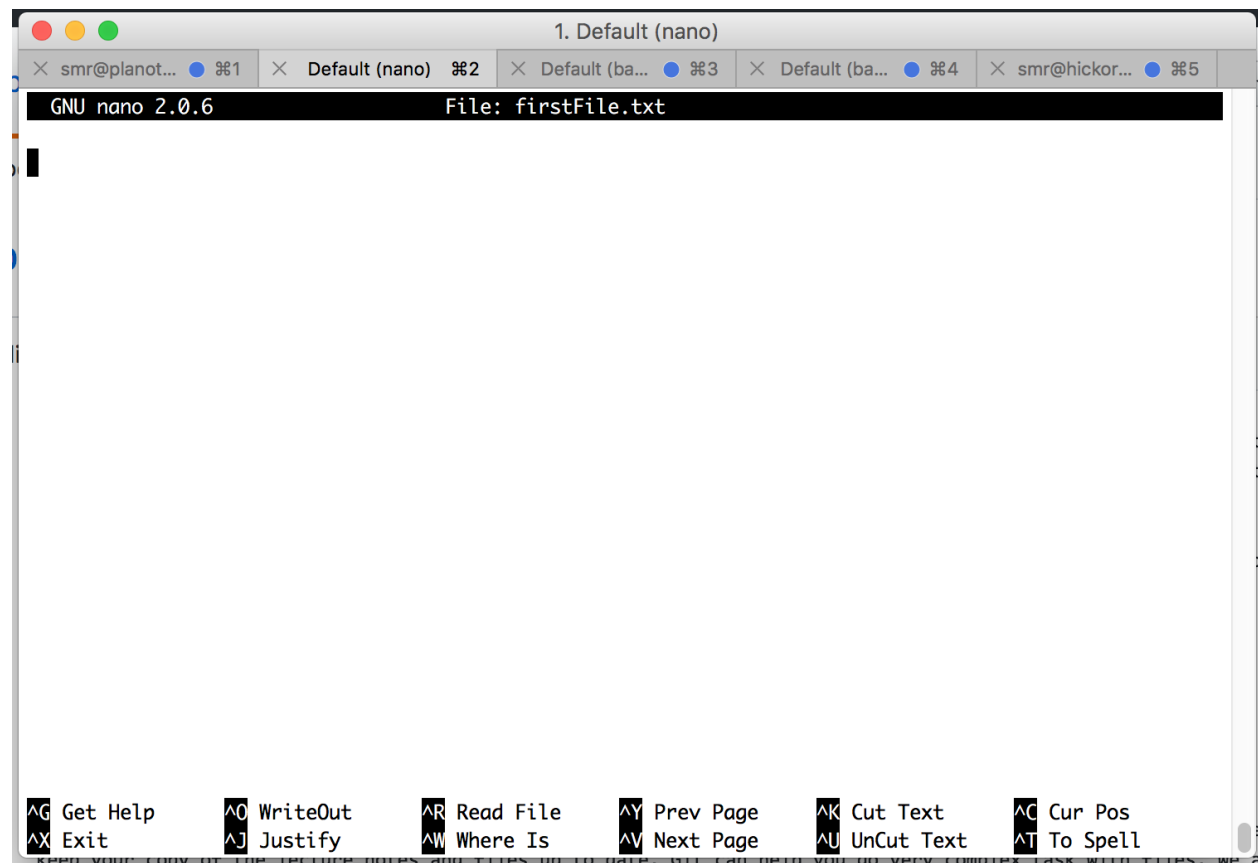


Figure 2: Create a new file with nano.

Things to notice: - At the top - the name of the program (nano) and it's version number - the name of the file you're editing - and whether the file has been modified since it was last saved. - In the middle - you will see either a blank area or text you have typed - At the bottom - A listing of keyboard commands such as Save (control + o) and Exit (control + x)

Keyboard commands are the only way to interact with the editor. You cannot use your mouse or trackpad.

```

GNU nano 2.0.6           File: firstFile.txt           Modified
line1
line2

```

Figure 3: Modified and not saced. In the top right corner it says “Modified”

Find more commands by using control g:

```

GNU nano 2.0.6           File: firstFile.txt

^G      (F1)             Display this help text
^X      (F2)             Close the current file buffer / Exit from nano
^O      (F3)             Write the current file to disk
^J      (F4)             Justify the current paragraph

^R      (F5)             Insert another file into the current one
^W      (F6)             Search for a string or a regular expression
^Y      (F7)             Move to the previous screen
^V      (F8)             Move to the next screen

^K      (F9)             Cut the current line and store it in the cutbuffer
^U      (F10)            Uncut from the cutbuffer into the current line
^C      (F11)            Display the position of the cursor
^T      (F12)            Invoke the spell checker, if available

^_      (F13) (M-G)      Go to line and column number
^\\      (F14) (M-R)      Replace a string or a regular expression
^^      (F15) (M-A)      Mark text at the cursor position
          (F16) (M-W)      Repeat last search

M-^      (M-6)          Copy the current line and store it in the cutbuffer
M-}      Indent the current line

^L      Refresh
^X      Exit

^Y      Prev Page
^V      Next Page

^P      Prev Line
^N      Next Line

M-^      First Line
M-^      Last Line

```

Figure 4: The help menu displays a listing of useful commands.

The Meta key is <esc>. To use the Meta+key, hit <esc>, release, then hit the following key

Helpful commands: - Jump to a specific line: - control + _ then line number - Copy a block of highlighted text - control + ^ then move your cursor to start to highlight a block for copying - Meta ^ to end your highlight block - Paste - control + u

Nano is a beginners text extior. vi and Emacs are better choices once you become a bit more comfortable using the terminal. These editors do cool stuff like syntax highlighting.

Git for Beginners

Git is a tool for managing files and versions of files. It is a *Version Control System*. It allows you to keep track of changes. You are going to be using Git to manage your course work and keep your copy of the lecture notes and files up to date. Git can help you do very complex task with files. We are going to keep it simple.

The Big Picture.

A Version Control System is good for Collaborations, Storing Versions, Restoring Previous Versions, and Managing Backups.

Collaboration

Using a Version Control System makes it possible to edit a document with others without the fear of overwriting someone's changes, even if more than one person is working on the same part of the document. All the changes can be merged into one document. These documents are all stored one place.

Storing Versions

A Version Control System allows you to save versions of your files and to attach notes to each version. Each save will contain information about the lines that were added or altered.

Restoring Previous Versions

Since you are keeping track of versions, it is possible to revert all the files in a project or just one file to a previous version.

Backup

A Version Control System makes it so that you work locally and sync your work remotely. This means you will have a copy of your project on your computer and the Version Control System Server you are using.

The Details

git is the Version Control System we will be using for tracking changes in our files.

GitHub is the Version Control System Server we will be using. They provide free account for all public projects.

The Basics

Creating a new repository

A repository is a project that contains all of the project files, and stores each file's revision history. Repositories can have multiple collaborators. Repositories usually have two components, one remote and one local.

Let's Do It!

Follow Steps 1 and 2 to create the remote repository. Follow Step 3 to create your local repository and link it to the remote.

1. Navigate to GitHub -> Create Account / Log In -> Go To Repositories -> Click 'New'
2. Add a name (i.e., PFB2017_problemsets) and a description (i.e., Solutions for PFB2017 Problem Sets) and click "Create Repository"
3. Create a directory on your computer and follow the instructions provided.
 - Open your terminal and navigate to the location you want to put a directory for your problem sets
 - Create a new directory directory (i.e., PFB2017_problemsets)
 - Follow the instructions provided when you created your repository. These are my instructions, yours will be different.

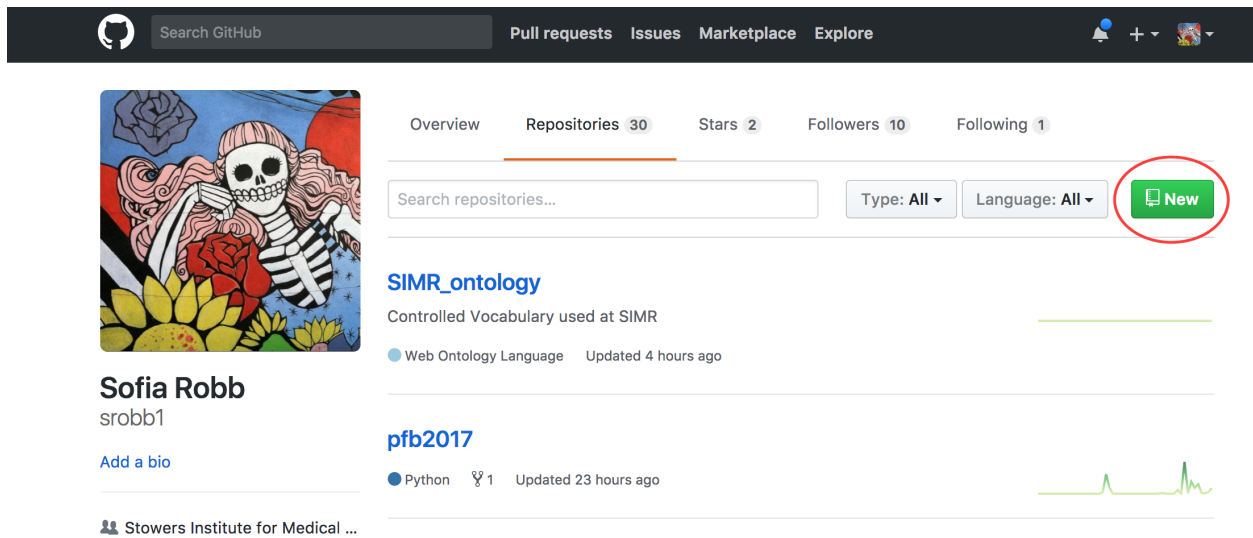


Figure 5: To create a new repository click the 'New' Button in the top right corner.

The screenshot shows the 'Create a new repository' form. The 'Owner' is 'srobb1' and the 'Repository name' is 'PFB2017_problemsets'. The 'Description (optional)' is 'Solutions for PFB2017 Problem Sets'. The 'Visibility' is set to 'Public'. The 'Initialize this repository with a README' checkbox is checked. The 'Create repository' button is at the bottom.

Figure 6: Fill in the form and click the 'Create Repository Button'

srobb1 / PFB2017_problemsets

Unwatch 1

Star 0

Fork 0

<> Code

Issues 0

Pull requests 0

Projects 0

Wiki

Settings

Insights

Quick setup — if you've done this kind of thing before

Set up in Desktop

 or

HTTPS

SSH

https://github.com/srobb1/PFB2017_problemsets.git

We recommend every repository include a [README](#), [LICENSE](#), and [.gitignore](#).

...or create a new repository on the command line

```
echo "# PFB2017_problemsets" >> README.md
git init
git add README.md
git commit -m "first commit"
git remote add origin https://github.com/srobb1/PFB2017_problemsets.git
git push -u origin master
```

...or push an existing repository from the command line

```
git remote add origin https://github.com/srobb1/PFB2017_problemsets.git
git push -u origin master
```

...or import code from another repository

You can initialize this repository with code from a Subversion, Mercurial, or TFS project.

Import code

Figure 7: Create a directory on your computer and follow these instructions.


```

1 echo "# PFB2017_problemsets" >> README.md
2 git init
3 git add README.md
4 git commit -m "first commit"
5 git remote add origin
   https://github.com/srobb1/PFB2017_problemsets.git
6 git push -u origin master

```

You now have a repository!

Let's back up a bit and talk more about git and about these commands. For basic git use, these are almost all the command you will need to know.

Every git repository has three main elements called *trees*: 1. *The Working Directory* contains your files 2. *The Index* is the staging area 3. *The HEAD* points to the last commit you made. > There are a few new words here, we will explain them as we go

command	description
<code>git init</code>	Creates your new local repository with the three trees on (local machine)
<code>git remote add remote-name URL</code>	Links your local repository to a remote repository that is often named <i>origin</i> and is found at the given URL.
<code>git add filename</code>	Propose changes and add file(s) with changes to the index or staging area (local machine)
<code>git commit -m 'message'</code>	Confirm or commit that you really want to add your changes to the HEAD (local machine)
<code>git push -u remote-name remote-branch</code>	Upload your committed changes in the HEAD to the specified remote repository to the specified branch

Let's Do it!

1. Make sure you are in your local repository
2. Create a new file with nano: `nano git_exercises.txt`
3. Add a line of text to the new file.
4. Save (control + o) and Exit (control + x)
5. (Add) Stage your changes. `git add git_exercises.txt`
6. (Commit) Become sure you want your changes your changes. `git commit -m 'added a line of text'`
7. (Push) Sync/Upload your changes to the remote repository. `git push origin master`

That is all there is to it! There are more complicated things you can do but we won't get into those. You will know when you are ready to learn more about git when you figure out there is something you want to do but don't know how. There are thousands of online tutorials for you to search and follow.

Cloning a Repository

Sometimes you want to download and use someone else's repository.

Let's clone the course material.

Let's do it!

1. Go to our PFB2017 GitHub Repository
2. Click the 'Clone or Download' Button
3. Copy the URL ~Clone PFB2017
4. *Clone* the repository to your local machine `git clone https://github.com/srobb1/pfb2017.git`

Now you have a copy of the course material on your computer!

If we make changes to any of these files and you want to update your copy you can *pull* the changes. `git pull`

Keeping track of differences between local and remote repositories

If you are ever wondering what do you need to add to your remote repository use the `git status` command. This will provide you a list of file that have been modified, deleted, and those that are untracked. Untracked files are those that have never been added to the staging area with `git add`

Links to *Slightly* less basic topics

You will KNOW if you need to use these features of git.

1. View Commit History
2. Resolving Merge Conflicts
3. Undoing Previous Commits

Link To Unix 2 Problem Set

Python

Python 1

Python Overview

- We are teaching the simplest form of the language, no ‘clever’ stuff.
- We can expand on this in subsequent years

Functions and statements we are NOT covering **assert**

Language basics: 3.x, interpreted, object-based, slow, python notebook eg jupyter

Python has - data types - functions - objects - classes - methods

Data types are just different types of data which are discussed in more detail later. Examples of data types are integer numbers and strings of letters and numbers (text). These can be stored in variables.

Functions do something with data, such as a calculation. Some functions are already built into Python. You can create your own functions as well.

****Objects*** are a way of grouping a set of data and functions (methods) that act on that data

Classes are a way to encapsulate (organize) variables and functions. Objects get their variables and methods from the class they belong to. **Methods** are just functions that belong to a Class. Objects that belong to the a Class can use Methods from that Class.

Running Python

Interactive Interpreter

Python can be run one line at a time in an interactive interpreter.

To launch the interpreter type the following into your terminal window:

```
$ python
```

Note: ‘\$’ indicates the command line prompt

First Python Commands:

```
1 $ python
```

```
1 >>> print("Hello , PFB2017!")
2 Hello , PFB2017!
```

Note: **print** is a function. Function names are followed by (), so formally, the function is **print()**

@sep writing python scripts: nano or Pycharm

Running Python Scripts

Typing the python command followed by the name of a script makes python execute the script. Recall that we just saw you can run an interactive interpreter by just typing `python` on the command line

Python Script

- The same code from above is typed into a text file.
- Python scripts are always saved in files whose names have the extension ‘.py’ (i.e. the filename ends with ‘.py’).

File Contents:

```
1 #!/usr/bin/python3
2 print ("Hello, PFB2017!")
```

Execute the Python script:

```
1 $ python test.py
```

This produces the following result:

```
1 Hello , PFB2017!
```

Syntax

@sep example code, statements, functions(), tab spacing, interactive interpreter, debugger

Python Identifiers:

A Python identifier is a name used to identify a variable, function, class, module or other object. An identifier starts with a letter A to Z or a to z or an underscore (`_`) followed by zero or more letters, underscores and digits (0 to 9).

Python does not allow punctuation characters such as `@`, `$`, and `%` within identifiers. Python is a case sensitive programming language. Thus, `seq_id` and `seq_ID` are two different identifiers in Python.

Naming conventions for Python Identifiers:

- The first character is lowercase, unless it is a name of a Class. Classes should begin with an uppercase characters. (ex. `Seq`)
- Private identifiers begin with an underscore. (ex. `_private`)
- Strong private identifiers begin with two underscores. (ex. `__private`)
- Language-defined special names begin and end with two underscores. (ex. `__special__`)

Reserved Words:

The following is a list of Python keywords. These are special words that already have a purpose in python and therefore cannot be used in identifier names.

```
1 and          exec          not
2 as           finally       or
3 assert       for           pass
4 break        from          print
5 class        global         raise
6 continue     if            return
7 def          import         try
8 del          in            while
9 elif         is            with
10 else        lambda         yield
11 except
```

Lines and Indentation:

Python denotes blocks of code by line indentation. Incorrect line spacing and/or indentation will cause an error to be reported or could make your code run in a way you don't expect. You can get help with indentation from good text editors or IDEs.

The number of spaces in the indentation need to be consistent but a specific number is not required. All lines of code, or statements, within a single block must be indented in the same way. For example:

Comments:

Comments are an essential programming practice. Making a note of what a line or block of code is doing will help the writer and readers of the code. This includes you!

#:

Comments start after a pound or hash symbol. All characters after the #, up to the end of the line are part of the comment and are ignored by python.

The first line of a script starting with **#!** is a special example of a comment that also has the special function in unix of telling the unix shell how to run the script.

```
1 #!/usr/bin/python3
2
3 # this is my first script
4 print ("Hello, PFB2017!") # this line prints
```

Blank Lines:

Blank lines are also important for increasing the readability of the code. Blank lines are ignored by the python interpreter

Python Options:

```
1 $ python -h
2 usage: python [option] ... [-c cmd | -m mod | file | -] [arg] ...
3 Options and arguments (and corresponding environment variables):
4 -c cmd : program passed in as string (terminates option list)
5 -d      : debug output from parser (also PYTHONDEBUG=x)
6 -E      : ignore environment variables (such as PYTHONPATH)
7 -h      : print this help message and exit
```

Data Types and Variables

sep note: overview of all types: start with constants: numbers, strings ‘ ’ “ ” “ ” “ ” r” = raw strings, + for concat, * for repeat iterable, string, int, float, tuple, list, dictionary, set

This is our first look at variables and data types. Each data type will be discussed in more detail in subsequent sections.

The first concept to consider is that python data types are either immutable (unchangeable) or not. Literal numbers, strings and tuples cannot be changed. Lists, dictionaries and sets can be changed. So can individual (scalar) variables. You can store data in memory by putting it in different kinds variables. You use the = sign to assign a value to a variable.

Numbers and Strings

Numbers and strings are two data types. Literal numbers and strings like this 5 or 'my name is' are immutable. However, their values can be stored in variables and then changed.

For Example:

```
1 first_variable = 5
```

Different types of data can be assigned to variables, i.e., integers (1,2,3), floats (floating point numbers, 3.1415), and strings (text).

For Example:

```
1 count      = 10      # this is an integer
2 average    = 2.5     # this is a float
3 message    = "Welcome to python" # this is a string
```

10, 2.5, and “Welcome to python” are singular pieces of data being stored in an individual variables.

Collections of data can also be stored in special data types, i.e., tuples, lists, sets, and dictionaries.

Lists

- Lists are used to store an ordered, *indexed* collection of data.
- Lists are mutable: the number of elements in the list and what's stored in each element can change
- Lists are enclosed in square brackets and items are separated by commas

```
1 [ 'atg' , 'aaa' , 'agg' ]
```

Index	Value
0	atg
1	aaa
2	agg

Command line parameters: A Special Built-in List

Command line parameters follow the name of a script or program and have spaces between them. They allow a user to pass information to a script on the command line when that script is being run. Python stores all the pieces of the command line in a special list called `sys.argv`.

You need to import the `sys` module at the beginning of your script like this

```
1 import sys
```

If you write this on the command line:

```
1 $ calculate_sum.py 5 7
```

This happens inside the script: > the script name, and the numbers 5 and 7 are contained in a list called `sys.argv`.

These are the command line parameters, or arguments you want to pass to your script.

`sys.argv[0]` is the script name.

You can access values of the other parameters by their indices, starting with 1, so `sys.argv[1]` is 5 and `sys.argv[2]` is 7.

If you wanted to calculate the sum in your script, you would add these two variables and print the result. Maybe your code would look something like this

```
1 #!/usr/bin/python3
2 import sys
3 a = sys.argv[1]
4 b = sys.argv[2]
5 print(a+b) # + is a sum operator on integers
```

Tuple

- Tuples are similar to lists and contain ordered, *indexed* collection of data.

- **Tuples are immutable: you can't change the values or the number of values**
- A tuple is enclosed in parentheses and items are separated by commas.

```
1 ( 'Jan' , 'Feb' , 'Mar' , 'Apr' , 'May' , 'Jun' , 'Jul' , 'Aug'
   , 'Sep' , 'Oct' , 'Nov' , 'Dec' )
```

Index	Value
0	Jan
1	Feb
2	Mar
3	Apr
4	May
5	Jun
6	Jul
7	Aug
8	Sep
9	Oct
10	Nov
11	Dec

Dictionary

- Dictionaries are good for storing data that can be represented as a two-column table.
- They store unordered collections of key/value pairs.
- A dictionary is enclosed in curly braces. and sets of Key/Value pairs are separated by commas
- A colon is written between each key and value. Commas separate key:value pairs.

```
1 { 'TP53' :
    'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC '
    , 'BRCA1' :
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
    }
```

Key	Value
TP53	GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC
BRCA1	GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA

Type Conversion

Sometimes you may need to convert data into a specific type. Here are some examples of functions that will help you to do this.

Function	Description
int(x)	Converts x to an integer.

Function	Description
<code>float(x)</code>	Converts x to a floating-point number.
<code>str(x)</code>	Converts x to a string.
<code>chr(x)</code>	Converts an integer to a character.
<code>tuple(s)</code>	Converts s to a tuple.
<code>list(s)</code>	Converts s to a list.
<code>set(s)</code>	Converts s to a set.
<code>dict(d)</code>	Creates a dictionary. d must be a sequence of (key,value) tuples.
<code>repr()</code>	Makes a string representation of an object. Useful for finding out what kind (Class) of object you are dealing with.

[Link to Python 1 Problem Set](#)

Python 2

Operators

@sep These depend on the type of variable, so revisit the most relevant ones in subsequent sections include `//` floor division and `%` modulus

Arithmetic Operators

Operator	Description	Example	Result
+	Addition	3+2	5
-	Subtraction	3-2	1
*	Multiplication	3*2	6
/	Division	3/2	1.5
%	Modulus (divides left operand by right operand and returns the remainder)	3%2	1
**	Exponent	3**2	9
//	Floor Division (result is the quotient with digits after the decimal point removed. If one of the operands is negative, the result is floored, i.e., rounded away from zero)	3//2 ; -11//3	1 ; -4

Assignment Operators

Operator	Equivalent to	Example	result evaluates to
=	a = 3	result = 3	3
+=	result = result + 2	result = 3 ; result += 2	5
-=	result = result - 2	result = 3 ; result -= 2	1
*	result = result * 2	result = 3 ; result *= 2	6
/=	result = result / 2	result = 3 ; result /= 2	1.5
%=	result = result % 2	result = 3 ; result %= 2	1
**=	result = result ** 2	result = 3 ; result **= 2	9
//=	result = result // 2	result = 3 ; result //= 3	1

Comparison Operators

These operators compare two values and returns true or false.

Operator	Description	Example	Result
==	equal to	3 == 2	False
!=	not equal	3 != 2	True
>	greater than	3 > 2	True
<	less than	3 < 2	False
>=	greater than or equal	3 >= 2	True
<=	less than or equal	3 <= 2	False

Logical Operators

Operator	Description	Example	Result
and	True if left operand is True and right operand is True	bool(3>=2 and 2<3)	True
or	True if left operand is True or right operand is True	bool(3==2 or 2<3)	True
not	Reverses the logical status	bool(not False)	True

Membership Operators

Operator	Description
in	True if a value is included in a list, tuple, or string
not in	True if a value is absent in a list, tuple, or string

For Example:

```
1 >>> dna =
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
2 >>> 'TCT' in dna
3 True
4 >>>
5 >>> 'ATG' in dna
6 False
7 >>> 'ATG' not in dna
8 True
9 >>> codons = [ 'atg' , 'aaa' , 'agg' ]
10 >>> 'atg' in codons
11 True
12 >>> 'ttt' in codons
13 False
```

Operator Precedence

Operators are listed in order of precedence. Highest listed first. Not all the operators listed here are mentioned above.

Operator	Description
**	Exponentiation (raise to the power)
~ + -	Complement, unary plus and minus (method names for the last two are +@ and -@)
* / % //	Multiply, divide, modulo and floor division
+ -	Addition and subtraction
>> <<	Right and left bitwise shift
&	Bitwise 'AND'
^	Bitwise exclusive 'OR' and regular 'OR'
<= < > >=	Comparison operators
<> == !=	Equality operators

Operator	Description
<code>= %= /= //=- += *= **=</code>	Assignment operators
<code>is</code>	Identity operator
<code>is not</code>	Non-identity operator
<code>in</code>	Membership operator
<code>not in</code>	Negative membership operator
<code>not, or, and</code>	logical operators

Truth

Lets take a step back, What is truth?

Everything is true, except for:

expression	TRUE/FALSE
<code>0</code>	FALSE
<code>None</code>	FALSE
<code>False</code>	FALSE
<code>''</code> (empty string)	FALSE
<code>[]</code> (empty list)	FALSE
<code>()</code> (empty tuple)	FALSE
<code>{}</code> (empty dictionary)	FALSE

Which means that these are True:

expression	TRUE/FALSE
<code>'0'</code>	TRUE
<code>'None'</code>	TRUE
<code>'False'</code>	TRUE
<code>'True'</code>	TRUE
<code>' '</code> (string of one blank space)	TRUE

Use `bool()` to test for truth

```

1 >>> bool(True)
2 True
3 >>> bool('True')
4 True
5 >>>
6 >>>
7 >>> bool(False)
8 False
9 >>> bool('False')
10 True
11 >>>
12 >>>
13 >>> bool(0)
14 False
15 >>> bool('0')
16 True
17 >>>
18 >>>
19 >>> bool('')
20 False
21 >>> bool(' ')
22 True
23 >>>
24 >>>
25 >>> bool(())
26 False
27 >>> bool([])
28 False
29 >>> bool({})
30 False

```

Logic: Control Statements

@sep while, if, elif, else, break, continue, @sep ugh: else in while loop

Control Statements are used to direct the flow of your code and create the opportunity for decision making. Control statements foundation is build on truth.

If Statement

- Use the If Statement to test for truth and to execute lines of code if true.
- When the expression evaluates to true each of the statements indented below the if statment, also known as the nested statement block, will be executed.

If Statement Syntax:

```

1 if expression :
2     statement
3     statement

```

For Example:

```

1 dna =
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
2 if 'AGC' in dna:
3     print('found AGC in your dna sequence')

```

Returns:

```

1 found AGC in your dna sequence

```

If/Else Statement

- The If portion of the if/else statement behave as before.
- The first indented block is executed if the condition is true.
- If the condition is false, the second indented else block is executed.

```

1 dna =
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
2 if 'ATG' in dna:
3     print('found ATG in your dna sequence')
4 else:
5     print('did not find ATG in your dna sequence')

```

Returns:

```

1 did not find ATG in your dna sequence

```

if/elif

- The if condition is tested as before and the indented block is executed if the condition is true.
- If it's false, the indented block following the elif is executed if the first elif condition is true.
- Any remaining elif conditions will be tested in order until one is found to be true. If none is true, the else indented block is executed.

```
1 count = 60
2 if count < 0:
3     message = "is less than 0"
4     print(count, message)
5 elif count < 50:
6     message = "is less than 50"
7     print (count, message)
8 elif count > 50:
9     message = "is greater than 50"
10    print (count, message)
11 else:
12    message = "must be 50"
13    print(count, message)
```

Returns:

```
1 60 is greater than 50
```

Let's change count to 20, which statement block gets executed?

```
1 count = 20
2 if count < 0:
3     message = "is less than 0"
4     print(count, message)
5 elif count < 50:
6     message = "is less than 50"
7     print (count, message)
8 elif count > 50:
9     message = "is greater than 50"
10    print (count, message)
11 else:
12    message = "must be 50"
13    print(count, message)
```

Returns:

```
1 20 is less than 100
```

What happens when count is 50?

```

1 count = 50
2 if count < 0:
3     message = "is less than 0"
4     print(count, message)
5 elif count < 50:
6     message = "is less than 50"
7     print(count, message)
8 elif count > 50:
9     message = "is greater than 50"
10    print(count, message)
11 else:
12    message = "must be 50"
13    print(count, message)

```

Returns:

```

1 50 must be 50

```

Numbers

Python recognizes 3 types of numbers: integers, float point numbers, and complex numbers.

integer

- known as an int
- an int can be positive or negative
- and **does not** contain a decimal point.

floating point number

- known as a float
- a floating point number can be positive or negative
- and **does** contain a decimal point

complex number

- known as complex
- is in the form of $a+bi$ where i is the imaginary part.

Conversion functions

function	Description
<code>int(x)</code>	to convert x to a plain integer
<code>float(x)</code>	to convert x to a floating-point number
<code>complex(x)</code>	to convert x to a complex number with real part x and imaginary part zero

function	Description
<code>complex(x, y)</code>	to convert x and y to a complex number with real part x and imaginary part y

```

1 >>> int(2.3)
2 2
3 >>> float(2)
4 2.0
5 >>> complex(2.3)
6 (2.3+0j)
7 >>> complex(2.3,2)
8 (2.3+2j)

```

Numeric Functions

function	Description
<code>abs(x)</code>	The absolute value of x: the (positive) distance between x and zero.
<code>round(x [,n])</code>	x rounded to n digits from the decimal point. <code>round(0.5)</code> is 1.0 and <code>round(-0.5)</code> is -1.0.
<code>max(x1, x2,...)</code>	The largest positive argument is returned
<code>min(x1, x2,...)</code>	The smallest argument is returned

```

1 >>> abs(2.3)
2 2.3
3 >>> abs(-2.9)
4 2.9
5 >>> round(2.3)
6 2
7 >>> round(2.5)
8 2
9 >>> round(2.9)
10 3
11 >>> round(-2.9)
12 -3
13 >>> round(-2.3)
14 -2
15 >>> round(-2.009,2)
16 -2.01
17 >>> max(4,-5,5,1,11)
18 11
19 >>> min(4,-5,5,1,11)
20 -5

```

Many numeric functions are not built into the Python core and need to be included in our script if we want

to use them. To include them at the tip of the script type: `import math`

These next functions are found in the math module and need to be imported. To use these function, prepend the function with the module name, i.e, `math.ceil(15.5)`

math.function	Description
<code>ceil(x)</code>	The smallest integer not greater than x is returned
<code>floor(x)</code>	the largest integer not greater than x is returned.
<code>exp(x)</code>	The exponential of x: e^x is returned
<code>log(x)</code>	the natural logarithm of x, for $x > 0$ is returned
<code>log10(x)</code>	The base-10 logarithm of x for $x > 0$ is returned
<code>modf(x)</code>	The fractional and integer parts of x are returned in a two-item tuple.
<code>pow(x, y)</code>	The value of $x^{**}y$ is returned
<code>sqrt(x)</code>	The square root of x for $x > 0$ is returned

```

1 >>> import math
2 >>>
3 >>> math.ceil(2.3)
4 3
5 >>> math.ceil(2.9)
6 3
7 >>> math.ceil(-2.9)
8 -2
9 >>> math.floor(2.3)
10 2
11 >>> math.floor(2.9)
12 2
13 >>> math.floor(-2.9)
14 -3
15 >>> math.exp(2.3)
16 9.974182454814718
17 >>> math.exp(2.9)
18 18.17414536944306
19 >>> math.exp(-2.9)
20 0.05502322005640723
21 >>>
22 >>> math.log(2.3)
23 0.8329091229351039
24 >>> math.log(2.9)
25 1.0647107369924282
26 >>> math.log(-2.9)
27 Traceback (most recent call last):
28   File "<stdin>", line 1, in <module>
29 ValueError: math domain error
30 >>>
31 >>> math.log10(2.3)
32 0.36172783601759284
33 >>> math.log10(2.9)
34 0.4623979978989561
35 >>>
36 >>> math.modf(2.3)
37 (0.29999999999999998, 2.0)
38 >>>
39 >>> math.pow(2.3,1)
40 2.3
41 >>> math.pow(2.3,2)
42 5.2899999999999999
43 >>> math.pow(-2.3,2)
44 5.2899999999999999
45 >>> math.pow(2.3,-2)
46 0.18903591682419663
47 >>>
48 >>> math.sqrt(25)
49 5.0
50 >>> math.sqrt(2.3)
51 1.51657508881031
52 >>> math.sqrt(2.9)
53 1.70293863659264

```

Comparing two numbers

The simple function `cmp(x,y)` is no longer available in python3.

Use this idiom instead:

```
1 cmp = (x>y)-(x<y)
```

It returns three different values depending on x and y

- if $x < y$ -1 is returned
- if $x > y$ 1 is returned
- $x == y$ 0 is returned

[Link to Python 2 Problem Set](#)

Python 3

Sequences

In the next section, we will learn about strings, tuples and lists. These are all examples of python sequences. A string is a sequence of characters 'ACGTGA', a tuple (0.23, 9.74, -8.17, 3.24, 0.16) and a list ['dog', 'cat', 'bird'] are sequences of any kind of data. We'll see much more detail in a bit.

In python a type of object gets operations that belong to that type. Sequences have sequence operations so strings can also use sequence operations. Strings also have their own specific operations.

You can ask what the length of any sequence is

```
1 >>>len('ACGTGA') # length of a string
2 6
3 >>>len ( (0.23, 9.74, -8.17, 3.24, 0.16) ) # length of a tuple,
    needs two parentheses (( ))
4 5
5 >>>len(['dog', 'cat', 'bird']) # length of a list
6 3
```

You can also use string-specific functions on strings, but not on lists and vice versa. We'll learn a lot more about this later on. `rstrip()` is a string method or function. You get an error if you try to use it on a list.

```
1 >>> 'ACGTGA'.rstrip('A')
2 'ACGTG'
3 >>> ['dog', 'cat', 'bird'].rstrip()
4 Traceback (most recent call last):
5   File "<stdin>", line 1, in <module>
6 AttributeError: 'list' object has no attribute 'rstrip'
```

What functions go with my object?

How do you find out what functions work with an object? There's a handy function `dir()`. As an example what functions can you call on our string 'ACGTGA'?

```

1 >>> dir('ACGTGA')
2 ['__add__', '__class__', '__contains__', '__delattr__', '__dir__',
  '__doc__', '__eq__', '__format__', '__ge__', '__getattribute__',
  '__getitem__', '__getnewargs__', '__gt__', '__hash__',
  '__init__', '__init_subclass__', '__iter__', '__le__', '__len__',
  '__lt__', '__mod__', '__mul__', '__ne__', '__new__',
  '__reduce__', '__reduce_ex__', '__repr__', '__rmod__',
  '__rmul__', '__setattr__', '__sizeof__', '__str__',
  '__subclasshook__', 'capitalize', 'casefold', 'center', 'count',
  'encode', 'endswith', 'expandtabs', 'find', 'format',
  'format_map', 'index', 'isalnum', 'isalpha', 'isdecimal',
  'isdigit', 'isidentifier', 'islower', 'isnumeric', 'isprintable',
  'isspace', 'istitle', 'isupper', 'join', 'ljust', 'lower',
  'lstrip', 'maketrans', 'partition', 'replace', 'rfind', 'rindex',
  'rjust', 'rpartition', 'rsplit', 'rstrip', 'split', 'splitlines',
  'startswith', 'strip', 'swapcase', 'title', 'translate', 'upper',
  'zfill']

```

You can call `dir()` on any object.

Strings

- A string is a series of characters starting and ending with a quotation mark.
- Strings are an example of a Python Sequence. A sequence is defined as a positionally ordered set. This means each element in the set has a position, starting with zero, i.e. 0,1,2,3 and so on until you get to the end of the string. If this is confusing, think about a string as being made up of individual characters: character 0, character 1, character 2 and so on.

Quotation Marks

- Single (')
- Double (")
- Triple (""" or """)

Notes about quotes:

- Single and double quotes are the same.
- A variable will not be replaced with its value (sometimes called 'interpolation') if placed inside of quotes.
- Triple quotes are used before and after a string that spans multiple lines.

Use of quotation examples:

```
1 word = 'word'
2 sentence = "This is a sentence."
3 paragraph = """This is a paragraph. It is
4 made up of multiple lines and sentences.
5 """
```

Strings and the `print()` function

We saw examples of `print()` earlier. Lets talk about it a bit more.

`print()` is a function that takes one or more comma-separated arguments.

Let's use the `print()` function to print a string.

```
1 >>>print("ATG")
2 ATG
```

We get ATG printed to the screen like we expect.

Let's assign a string to a variable and print the variable.

```
1 >>>dna = 'ATG'
2 ATG
3 >>> print(dna)
4 ATG
```

We get ATG printed to the screen

What happens if we put the variable in quotes?

```
1 >>>dna = 'ATG'
2 ATG
3 >>> print("dna")
4 dna
```

The literal value of 'dna' is printed to the screen. The variable called 'dna' is not interpolated when it is inside of quotes.

Let's see what happens when we give `print()` two literal strings as arguments.

```
1 >>> print("ATG","GGTCTAC")
2 ATG GGTCTAC
```

We get the two literal strings printed to the screen separated by a space

What if you do not want your strings separated by a space? Use the concatenation operator to concatenate the two strings before or within the `print()` function.

```
1 >>> print("ATG"+"GGTCTAC")
2 ATGGGTCTAC
3 >>> combined_string = "ATG"+"GGTCTAC"
4 ATGGGTCTAC
5 >>> print(combined_string)
6 ATGGGTCTAC
```

We get the two strings printed to the screen without being separated by a space.
You can also use this

```
1 >>> print('ATG','GGTCTAC',sep='')
2 ATGGGTCTAC
```

Now, lets print a variable and a literal string.

```
1 >>> dna = 'ATG'
2 ATG
3 >>> print(dna,'GGTCTAC')
4 ATG GGTCTAC
```

We get the value of the variable and the literal string printed to the screen separated by a space

How would we print the two without a space?

```
1 >>> dna = 'ATG'
2 ATG
3 >>> print(dna + 'GGTCTAC')
4 ATGGGTCTAC
```

Something to think about: Values of variable are variable. Or in other words, they are mutable, changeable.


```

1 >>> dna = 'ATG'
2 ATG
3 >>> print(dna)
4 ATG
5 >>> dna = 'TTT'
6 TTT
7 >>> print(dna)
8 TTT

```

The new value of the variable 'dna' is printed to the screen when 'dna' is an argument for the `print()` function.

Errors and Printing

Let's look at the typical errors you will encounter when you use the `print()` function.

What will happen if you forget to close your quotes?

```

1 >>> print("GGTCTAC)
2 File "<stdin>", line 1
3     print("GGTCTAC)
4         ^
5 SyntaxError: EOL while scanning string literal

```

We get a 'SyntaxError' if the closing quote is not used

What will happen if you forget to enclose your literal string in quotes in a print statement?

```

1 print(GGTCTAC)
2 Traceback (most recent call last):
3   File "<stdin>", line 1, in <module>
4 NameError: name 'GGTCTAC' is not defined

```

We get a 'NameError' when the literal string is not enclosed in quotes.

```

1 >>> print "boo"
2 File "<stdin>", line 1
3     print "boo"
4         ^
5 SyntaxError: Missing parentheses in call to 'print'

```

In python2, the command was `print`, but this changed in python3 to `print()`, so don't forget the parentheses!

Special/Escape Characters

How would you include a new line, carriage return, or tab in your string?

Escape Character	Description
<code>\n</code>	New line
<code>\r</code>	Carriage Return
<code>\t</code>	Tab

Let's include some escape characters in our strings and `print()` functions.

```
1 >>> string_with_newline = 'this sting has a new line\nthis is the
    second line'
2 >>> print(string_with_newline)
3 this sting has a new line
4 this is the second line
```

We printed a new line to the screen

Generally, you don't have to worry about adding whitespace between arguments because `print()` adds space between arguments and a new line at the end for you. You can change these with `sep=` and `end=` `print('one line', 'second line', 'third line', sep='\n', end = '')`

A neater way to do this is to express a multi-line string enclosed in triple quotes (`"""`).

```
1 >>> print("""this sting has a new line
2 ... this is the second line""")
3 this sting has a new line
4 this is the second line
```

Let's print a tab character (`\t`).

```
1 >>> line = "value1\tvalue2\tvalue3"
2 >>> print(line)
3 value1  value2  value3
```

We get the three words separated by tab characters. A common format for data is tab separated.

You can add a backslash before any character to force it to be printed as a literal. This is called 'escaping'. This is only really useful for printing literal quotes `'` and `"`

```

1 >>> print("this is a \'word\'")
2 this is a 'word'
3 >>> print("this is a 'word'")
4 this is a 'word'

```

In both cases the single quote is printed to the screen as a quote.

If you want every character in your string to remain exactly as it is, declare your string a raw string literal with 'r' before the first quote. This looks ugly, but it works.

```

1 >>> line = r"value1\tvalue2\tvalue3"
2 >>> print(line)
3 value1\tvalue2\tvalue3

```

Our escape characters '\t' remain as we typed them, they are not converted to actual tab characters.

Concatenation

To concatenate strings use the concatenation operator '+'

```

1 >>> promoter= 'TATAAA'
2 >>> upstream = 'TAGCTA'
3 >>> downstream = 'ATCATAAT'
4 >>> dna = upstream + promoter + downstream
5 >>> print(dna)
6 TAGCTATATAAAATCATAAT

```

The concatenation operator can be used to combine strings. The newly combined string can be stored in a variable. What happens if you use + with numbers (these are integers or ints)?

```

1 >>> 4+3
2 7

```

For strings, + concatenates; for integers, + adds.

You need to convert the numbers to strings before you can concatenate them

```
1 >>> str(4) + str(3)
2 '43'
```

Determine the length of a string:

Use the `len()` function to calculate the length of a string. This function takes a sequence as an argument and returns an int

```
1 >>> print(dna)
2 TAGCTATATAAAATCATAAT
3 >>> len(dna)
4 20
```

The length of the string, including spaces, is calculated and returned.

The value `len()` returns can be stored in a variable.

```
1 >>> dna_length = len(dna)
2 >>> print(dna_length)
3 20
```

You can mix strings and ints in `print()`, but not in concatenation.

```
1 >>> print("The length of the DNA sequence:" , dna , "is" , dna_length)
2 The length of the DNA sequence: TAGCTATATAAAATCATAAT is 20
```

Changing String Case

Changing the case of a string is a bit different that you might first expect. For example, to lowercase a string we need to use a method. A method is a function that is specific to a object. When we assign a string to a variable we are creating an instance of a string object. This object has a series of methods that will work on the data that is stored in the object. `lower()` function is one of these object methods.

Let's do create a new string object.

```
1 dna = "ATGCTTG"
```

Look familiar? It should. Creating a string object is what we have been doing all along.

Now that we have a string object we can use string methods. The way you use a method is to append the method with a `.` to the variable name.

```
1 >>> dna = "ATGCTTG"
2 >>> dna.lower()
3 'atgcttg'
```

the `lower()` method returns the contents stored in the `'dna'` variable in lowercase.

The contents of the `'dna'` variable have not been changed. If you want to reuse the returned lowercased value, store it in a new variable.

```
1 >>> print(dna)
2 ATGCTTG
3 >>> dna_lowercase = dna.lower()
4 >>> print(dna)
5 ATGCTTG
6 >>> print(dna_lowercase)
7 atgcttg
```

The string method can be nested inside of other functions.

```
1 >>> dna = "ATGCTTG"
2 >>> print(dna.lower())
3 atgcttg
```

The value of `'dna'` is lowercased and returned. The `print()` function takes the returned value from the `lower()` method and prints it. If you try to use a string method on a object that is not a string you will get an error.

```
1 >>> nt_count = 6
2 >>> dna_lc = nt_count.lower()
3 Traceback (most recent call last):
4   File "<stdin>", line 1, in <module>
5 AttributeError: 'int' object has no attribute 'lower'
```

You get an `AttributeError` when you use a method on the an incorrect object type. We are told that the `int` object (an `int` is returned by `len()`) does not have a function called `lower`.

Now let's uppercase a string.

```

1 >>> dna = 'attgct'
2 >>> dna.upper()
3 'ATTGCT'
4 >>> print(dna)
5 attgct

```

The contents of the variable 'dna', 'attgct' was returned in upper case. The actual contents of 'dna' were not altered.

Find and Count

`count(str)` returns the number of exact matches of `str` it found (as an int)

```

1 >>> dna = 'ATGCTGCATT'
2 >>> dna.count('T')
3 4

```

The number of times 'T' is found is returned. The string stored in 'dna' is not altered.

Find and Replace

`replace(str1,str2)` returns a new string with all matches of `str1` in a string replaced with `str2`.

```

1 >>> dna = 'ATGCTGCATT'
2 >>> dna.replace('T','U')
3 'AUGCUGCAUU'
4 >>> print(dna)
5 ATGCTGCATT
6 >>> rna = dna.replace('T','U')
7 >>> print(rna)
8 AUGCUGCAUU

```

All occurrences of T are replaced by U. The new string is returned. The original string has not actually been altered. If you want to reuse the new string, store it in a variable.

Extracting a Substring, or Slicing

Parts of a string can be located based on position and returned. This is because a string is a sequence. Coordinates start at 0. You add the coordinate in square brackets after the string's name.

This string 'ATTAAAGGGCCC' is made up of the following sequence of characters, and positions (starting at zero).

Position/Index	Character
0	A
1	T
2	T
3	A
4	A
5	A
6	G
7	G
8	G
9	C
10	C
11	C

Let's return the 4th, 5th, and 6th nucleotides. To do this, we need to count like a computer and start our string at 0 and return the 3rd, 4th, and 5th characters. This will be everything from 3 to 6.

```
1 >>> dna = 'ATTAAAGGGCCC'
2 >>> sub_dna = dna[3:6]
3 >>> print(sub_dna)
4 AAA
```

The characters with the positional index of 3, 4, 5 are returned. Or in other words, every character starting at index 3 and up to but not including the index of 6 are returned.

Let's return the first 6 characters.

```
1 >>> dna = 'ATTAAAGGGCCC'
2 >>> sub_dna = dna[0:6]
3 >>> print(sub_dna)
4 ATTAAA
```

Every character starting at index 0 and up to but not including index 6 are returned. This is the same as `dna[:6]`

Let's return every character from index 6 to the end of the string.

```
1 >>> dna = 'ATTAAAGGGCCC'
2 >>> sub_dna = dna[6:]
3 >>> print(sub_dna)
4 GGGCCC
```

When the second argument is left blank, every character from index 6 and greater is returned.

Let's return the last 3 characters.

```
1 >>> sub_dna = dna[-3:]
2 >>> print(sub_dna)
3 CCC
```

When the second argument is left blank and the first argument is negative (-X), X characters from the end of the string are returned.

Locate and Report

The positional index of an exact string in a larger string can be found and returned with the string method `find`. A exact string is given as an argument and the index of its first occurrence is returned. -1 is returned if it is not found.

```
1 >>> dna = 'ATTAAAGGGCCC'
2 >>> dna.find('T')
3 1
4 >>> dna.find('N')
5 -1
```

The substring 'T' is found for the first time at index 1 in the string 'dna' so 1 is returned. The substring 'N' is not found, so -1 is returned.

Other String Methods

Since these are methods, be sure to use in this format `string.method()`.

function	Description
<code>s.strip()</code>	returns a string with the whitespace removed from the start and end
<code>s.isalpha()</code>	tests if all the characters of the string are alphabetic characters. Returns True or False.
<code>s.isdigit()</code>	tests if all the characters of the string are numeric characters. Returns True or False.
<code>s.startswith('other_string')</code>	tests if the string starts with the string provided as an argument. Returns True or False.
<code>s.endswith('other_string')</code>	tests if the string ends with the string provided as an argument. Returns True or False.
<code>s.split('delim')</code>	splits the string on the given exact delimiter. Returns a list of substrings. If no argument is supplied, the string will be split on whitespace.

function	Description
s.join(list)	opposite of split(). The elements of a list will be concatenated together using the string stored in 's' as a delimiter.

String Formatting

Strings can be formatted using the `format()` function. Pretty intuitive! For example, if you want to include literal strings and variables in your print statement and do not want to concatenate or use multiple arguments in the `print()` function you can use string formatting.

```

1 >>> string = "This sequence: {} is {} nucleotides long and is found
  in {}."
2 >>> string.format(dna,dna_len,gene_name)
3 'This sequence: TGAACATCTAAAAGATGAAGTTT is 23 nucleotides long and
  is found in Brca1.'
4 >>> print(string)
5 This sequence: {} is {} nucleotides long and is found in {}.
6 >>> new_string = string.format(dna,dna_len,gene_name)
7 >>> print(new_string)
8 This sequence: TGAACATCTAAAAGATGAAGTTT is 23 nucleotides long and is
  found in Brca1.
```

We put together the three variables and literal strings into a single string using the function `format()`. The original string is not altered, a new string is returned that incorporates the arguments. You can save the returned value in a new variable. Each `{}` is a placeholder for the strings that need to be inserted.

Something nice about `format()` is that you can print int and string variable types without converting first.

You can also directly call the `format` function on a string inside a `print` function. Here are two examples

```

1 >>> string = "This sequence: {} is {} nucleotides long and is found
  in {}."
2 >>> print(string.format(dna,dna_len,gene_name))
3 This sequence: TGAACATCTAAAAGATGAAGTTT is 23 nucleotides long and is
  found in Brca1.
```

Or you can create a string on the fly and use the `format()` function.

```

1 >>> print( "This sequence: {} is {} nucleotides long and is found in
  {}.".format(dna,dna_len,gene_name))
2 This sequence: TGAACATCTAAAAGATGAAGTTT is 23 nucleotides long and is
  found in Brca1.
```

There is no need to store the string in a variable.

The `format()` mini-language

So far, we have just used `{}` to show where to insert the value of a variable in a string. You can add special characters inside the `{}` to change the way the variable is formatted when it's inserted into the string.

You can number these, not necessarily in order.

```
1 >>> '{0}, {1}, {2}'.format('a', 'b', 'c')
2 'a, b, c'
3 >>> '{2}, {1}, {0}'.format('a', 'b', 'c')
4 'c, b, a'
```

To change the spacing of strings and the way numbers are formatted, you add `:` and other special characters like this `{:>5}` to right-justify a string in a five-character field.

Lets right justify some numbers.

```
1 >>> print( "{:>5}".format(2) )
2     2
3 >>> print( "{:>5}".format(20) )
4    20
5 >>> print( "{:>5}".format(200) )
6   200
```

The arguments: 2, 20, 200 have all been right justified in a field that is 5 characters wide by using `'{:>5}'`

How about padding with zeroes? This means the five-character field will be filled as needed with zeroes to the left of any numbers you want to display

```
1 >>> print( "{:>05}".format(2) )
2 00002
3 >>> print( "{:>05}".format(20) )
4 00020
```

Now all of the arguments: 2 ,20 are right justified to a width of 5 and any empty space is filled with a zero.

Use a `<` to indicate left-justification.

```

1 >>> print( "{:<5} genes".format(2) )
2 2      genes
3 >>> print( "{:<5} genes".format(20) )
4 20     genes
5 >>> print( "{:<5} genes".format(200) )
6 200    genes

```

The arguments: 2, 20, 200 have all been left justified by a width of 5 columns by using ‘{:<5}’

If you want to pad with a non-zero character, you can add this too. Between : and the symbol for the justification. Let’s try padding with _, an underscore.

```

1 >>> print( "{:_<5} next".format(2) )
2 2_ _ _ _ next
3 >>> print( "{:_<5} next".format(20) )
4 20_ _ _ next
5 >>> print( "{:_<5} next".format(200) )
6 200_ _ next

```

The arguments: 2, 20, 200 have all been left justified in a field that is 5 characters wide and any empty space is filled with a _

Center aligning is done with ^

```

1 >>> print( "{:_^10}".format(2) )
2 ##### _2_
3 >>> print( "{:_^10}".format(20) )
4 ##### _20_
5 >>> print( "{:_^10}".format(200) )
6 ##### _200_

```

Text can be center aligned by using ‘:^10’. 10 of course is your column width. The ‘^’ indicates center justification. In our example an underscore is used to illustrate the empty spaces.

Summary of special formatting symbols so far

Here are some of the ALIGNMENT options:

Option	Meaning
'<'	Forces the field to be left-aligned within the available space (this is the default for most objects).
'>'	Forces the field to be right-aligned within the available space (this is the default for numbers).
'='	Forces the padding to be placed after the sign (if any) but before the digits. This is used for printing fields in the form '+000000120'.
'^'	This alignment option is only valid for numeric types. Forces the field to be centered within the available space.

Here's an example

```
{ : x < 10 s}
```

fill with x

left justify <

10 a field of ten characters s a string

Common Types

type	description
b	convert to binary
d	decimal integer
e	exponent, default precision is 6, uses e
E	exponent, uses E

type	description
f	floating point, default precision 6 (also F)
g	general number, float for values close to 0, exponent for others; also G
s	string, default type
x	convert to hexadecimal, also X
%	converts to % by multiplying by 100

What's the point?

So much can be done with the `format()` function. Here is one last example, but not the last functionality of this function. Let truncate a long floating point number. The default is 6 decimal places. Note that the function rounds to the nearest decimal place.

```
1 '{:f}'.format(3.141592653589793)
2 '3.141593 '
3 >>> '{:.4f}'.format(3.141592653589793)
4 '3.1416 '
```

Lists and Tuples

Lists

Lists are valuable data types that can store a collection of data in a single variable.

- Lists are used to store an ordered, *indexed* collection of data.
- Values are separated by commas
- Values are enclosed in square brackets `[]`
- Lists can grow and shrink
- Values are mutable

Accessing Values in Lists

To retrieve a single value in a list use the value's index in this format `list[index]`. This will return the value at the specified index.

Here is our list:

```
1 >>> codons = [ 'atg' , 'aaa' , 'agg' ]
```

There are 3 values with the indices of 0, 1, 2

Index	Value
0	atg
1	aaa
2	agg

Let's access the 0th value.

```
1 >>> codons = [ 'atg' , 'aaa' , 'agg' ]
2 >>> codons[0]
3 'atg'
```

The the value of the 0th index is returned by using the syntax list[index]

The value cana be saved for later use by storing in a variable.

```
1 >>> codons = [ 'atg' , 'aaa' , 'agg' ]
2 >>> codon_0 = codons[0]
3 >>> print(codon_0)
4 atg
```

Each value can be saved in a new variable to use later.

The values can be retrieved an used directly.

```
1 >>> codons = [ 'atg' , 'aaa' , 'agg' ]
2 >>> print(codons[0])
3 atg
4 >>> print(codons[1])
5 aaa
6 >>> print(codons[2])
7 agg
```

The 3 values are independently accesses and immediately printed. They are not stored in a variable.

If you want to access the values in reverse, use negative indices.

```
1 >>> codons = [ 'atg' , 'aaa' , 'agg' ]
2 >>> print(codons[-1])
3 agg
4 >>> print(codons[-2])
5 aaa
```

Using a negative index will return the values from the end of the list. For example, -1 is the index of the last value 'agg'. This value also has an index of 2.

Changing Values in a List

Individual values can be changed using the value's index and the assignment operator.

```
1 >>> print(codons)
2 ['atg', 'aaa', 'agg']
3 >>> codons[2] = 'cgc'
4 >>> print(codons)
5 ['atg', 'aaa', 'cgc']
```

codon[2] used to contain 'agg'. We reassigned codon[2] to contain the new value 'cgc'

What about trying to assign a value to an index that does not exist?

```
1 >>> codons[5] = 'aac'
2 Traceback (most recent call last):
3   File "<stdin>", line 1, in <module>
4 IndexError: list assignment index out of range
```

codon[5] does not exist, and when we try to assign a value to this index we get an IndexError.

Extracting a Subset of a List, or Slicing

This works in exactly the same way with lists as it does with strings. This is because both are Sequences, or ordered collections of data with positional information.

Index	Value
0	atg
1	aaa
2	agg
3	aac
4	cgc
5	acg

```
1 >>> codons = [ 'atg' , 'aaa' , 'agg' , 'aac' , 'cgc' , 'acg' ]
2 >>> print (codons[1:3])
3 ['aaa', 'agg']
4 >>> print (codons[3:])
5 ['aac', 'cgc', 'acg']
6 >>> print (codons[:3])
7 ['atg', 'aaa', 'agg']
8 >>> print (codons[0:3])
9 ['atg', 'aaa', 'agg']
```

codons[1:3] returns every value starting with the value of codons[1] up to but not including codons[3] codons[3:] returns every value starting with the value of codons[3] and every value after. codons[:3] returns every value up to but not including codons[3] codons[0:3] is the same as codons[:3]

List Operators

Operator	Description	Example
+	Concatenation	[10, 20, 30] + [40, 50, 60] returns [10, 20, 30, 40, 50, 60]
*	Repetition	['atg'] * 4 returns ['atg','atg','atg','atg']
in	Membership	20 in [10, 20, 30] returns True

List Functions

Functions	Description	Example
len(list)	returns the length or the number of values in list	len([1,2,3]) returns 3
max(list)	returns the value with the largest ascii value	max(['a','A','z']) returns 'z'
min(list)	returns the value with the smallest ascii value	min(['a','A','z']) returns 'A'
list(seq)	converts a tuple into a list	list(('a','A','z')) returns ['a', 'A', 'z']
sorted(list, key=None, reverse=False)	returns a sorted list based on the key provided	sorted(['a','A','z']) returns ['A', 'a', 'z'] sorted(['a','A','z'],key=str.lower) returns ['a', 'A', 'z']

List Methods

Remember methods are apart of the object and are used in the following format list.method().

For these examples use: list = [1,2,3] and codons = ['atg' , 'aaa' , 'agg']

Method	Description	Example
list.append(obj)	appends an object to the end of a list	list.append(9) ; print(list) ; returns [1,2,3,9]
list.count(obj)	counts the occurance of an object in a list	list.count(2) returns 1
list.index(obj)	returns the lowest index where the given object is found	list.index(2) returns 1
list.pop()	removes and returns the last value in the list. The list is now 1 value shorter	list.pop() returns 3
list.insert(index, obj)	inserts a value at the given index	list.insert(0,100) ; print(list) returns [100, 1, 2, 3]
list.extend(new_list)	adds the provided list to the end of list	list.extend(['a', 'z']) ; print(list) returns [1, 2, 3, 'a', 'z']

Method	Description	Example
<code>list.pop(index)</code>	removes and returns the value of the index argument. The list is now 1 value shorter	<code>list.pop(0)</code> returns 1
<code>list.remove(obj)</code>	finds the lowest index of the given object and removes the value. The list is now 1 value shorter	<code>codons.remove('aaa') ; print(codons)</code> returns ['atg' , 'agg']
<code>list.reverse()</code>	reverses the order of the list	<code>list.reverse() ; print(list)</code> returns [3,2,1]
<code>list.copy()</code>	Returns a shallow copy of list. Shallow vs Deep only matters in multidimensional datastructures.	
<code>list.sort([func])</code>	sorts a list using the provided function. Does not return a list. The list has been changed. Advanced list sort will be covered once writing your own functions has been discussed.	<code>codons.sort() ; print(codons)</code> returns ['aaa', 'agg', 'atg']

Be careful how you make a copy of your list

```

1 >>> list=['a', 'one', 'two']
2 >>> l2=list
3 >>> l2.append('1')
4 >>> print(list)
5 ['a', 'one', 'two', '1']
6 >>> print(l2)
7 ['a', 'one', 'two', '1']

```

Not what you expected?! Boht lists have changed because we only copied a pointer to the original list when we wrote `l2=list`.

Let's copy the list using the `copy()` method.

```

1 >>> list=['a', 'one', 'two']
2 >>> l2=list.copy()
3 >>> l2.append('1')
4 >>> print(list)
5 ['a', 'one', 'two']

```

There we go, we get what we expect this time!

Building a List one Value at a Time

Now that you have seen the `append()` function we can go over how to build a list one value at a time.

```
1 >>> words = []
2 >>> print(words)
3 []
4 >>> words.append('one')
5 >>> words.append('two')
6 >>> print(words)
7 ['one', 'two']
```

We start with a an empty list called 'words'. We use `append()` to add the value 'one' then to add the value 'two'. We end up with a list with two values.

[Link to Python 3 Problem Set](#)

Python 4

Loops

@sep > for i in list: > (ugh: else in for loop) > iterators, range()

All of the coding that we have gone over so far has been executed line by line. Sometimes there are blocks of code that we want to execute more than once. Loops let us do this.

There are two loop types: 1. while loop 2. for loop

while loop

The while loop will continue to execute the while loop block as long as a given condition returns True.

While Loop Syntax

```
1 while expression:
2     statement(s)
```

The condition is found in the expression. The while loop block of code is the indented statements following the expression.

This example is not using the python interpreter. It is a python script that will be executed on the command line.

Code:

```
1 #!/usr/bin/python3
2
3 count = 0
4 while count < 5:
5     print("count:" , count)
6     count+=1
7 print("Done")
```

Output:

```
1 $ python while.py
2 count: 0
3 count: 1
4 count: 2
5 count: 3
6 count: 4
7 Done
```

The while condition was true 5 times and the while block of code was executed 5 times.

- count is equal to 0 when we begin
- 0 is less than 5 so we execute the while block 1 time
- count is printed
- count is incremented (count = count + 1)
- count is now equal to 1.
- 1 is less than 5 so we execute the while block for the 2nd time.
- this continues until count is 5.
- 5 is not less than 5 so the while block is not executed.
- The first line following the while statement is executed, “Done” is printed

An infinite loop is when a while condition is always true. Here is an example of an infinite loop.

```
1 #!/usr/bin/python3
2
3 count = 0
4 while count < 5:
5     print("count:" , count)
6 print("Done")
```

Output:

```
1 $ python infinit.py
2 count: 0
3 count: 0
4 count: 0
5 count: 0
6 count: 0
7 count: 0
8 count: 0
9 count: 0
10 ...
11 ...
```

What caused this code to never be true? The statement that increments the count is missing. To stop the code from forever printing use Ctrl+C.

While/Else

An Else statement can be used with a while statement. It behaves in the same way as with an If statement. When the while statement is false, the else block is executed ONE TIME.

```
1 #!/usr/bin/python3
2
3 count = 0
4 while count < 5:
5     print("count:" , count)
6     count+=1
7 else:
8     print("count:", count, "is now not less than 5")
9 print("Done")
```

Output:

```
1 $ python while_else.py
2 count: 0
3 count: 1
4 count: 2
5 count: 3
6 count: 4
7 count: 5 is now not less than 5
8 Done
```

The while loop was executed five times like before. Now when count is equal to 5 and therefore not less than 5, the else block is executed. Finally the lines of code outside the while/else are executed.

For Loops

A for loop is a loop that executes the for block of code for every iteration of a sequence. Remember a sequence is an ordered collection of data.

For Loop Syntax

```
1 for iterating_variable in sequence:
2     statement(s)
```

An example of a sequence is a list. Let's use a for loop list of words.

Code:

```

1 #!/usr/bin/python3
2
3 words = ['zero','one','two','three','four']
4 for word in words:
5     print(word)

```

Output:

```

1 python3 list_words.py
2 zero
3 one
4 two
5 three
6 four

```

This is an example of iterating over a string. Remember a string is a sequence like a list. It is data where position is important. Look back at “Extracting a Substring, or Slicing” in the Strings section to see other ways of how strings can be treated like lists.

Code:

```

1 #!/usr/bin/python3
2
3 dna =
4     'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
5 for nt in dna:
6     print(nt)

```

Output:

```

1 $ python3 for_string.py
2 G
3 T
4 A
5 C
6 C
7 T
8 T
9 ...
10 ...

```

This is a easy way to access each character in a string. It is especially nice for DNA sequences.

Another example of iterating over a list of variables, this time numbers.

Code:

```
1 #!/usr/bin/python3
2
3 numbers = [0,1,2,3,4]
4 for num in numbers:
5     print(num)
```

Output:

```
1 $ python3 list_numbers.py
2 0
3 1
4 2
5 3
6 4
```

Python has a function called `range()` that will return numbers that can be converted to a list.

```
1 >>> range(5)
2 range(0, 5)
3 >>> list(range(5))
4 [0, 1, 2, 3, 4]
```

This can be used in conjunction with a for loop to iterate over a range of numbers Code:

```
1 #!/usr/bin/python3
2
3 for num in range(5):
4     print(num)
```

Output:

```
1 $ python list_range.py
2 0
3 1
4 2
5 3
6 4
```

As you can see this is the same output as using the list `numbers = [0, 1, 2, 3, 4]` And this has the same functionality as a while loop with a condition of `count = 0 ; count < 5`.

For/Else

An else statement can be used with a for loop as well. The else block of code will be executed when the for loop exits normally.

Can we just skip this. it seems pretty stupid

Loop Control

Loops control statements allow for altering the normal flow of execution.

Control Statement	Description
break	A loop is terminated when a break statement is executed. All the lines of code after the break, but within the loop block are not executed. No more iteration of the loop are performed
continue	A single iteration of a loop is terminated when a continue statement is executed. The next iteration will proceed normally.

Loop Control: Break

Code:

```
1 #!/usr/bin/python3
2
3 count = 0
4 while count < 5:
5     print("count:" , count)
6     count+=1
7     if count == 3:
8         break
9 print("Done")
```

Output:

```
1 $ python break.py
2 count: 0
3 count: 1
4 count: 2
5 Done
```

when the count is equal to 3, all iterations of the while loop are terminated

Loop Control: Continue

Code:

```
1 #!/usr/bin/python3
2
3 count = 0
4 while count < 5:
5     print("count:" , count)
6     count+=1
7     if count == 3:
8         continue
9     print("line after our continue")
10 print("Done")
```

Output:

```
1 $ python continue.py
2 count: 0
3 line after our continue
4 count: 1
5 line after our continue
6 count: 2
7 count: 3
8 line after our continue
9 count: 4
10 line after our continue
11 Done
```

When the count is equal to 3 the continue is executed. This causes all the lines within the loop block to be skipped. “line after our continue” is not printed when count is equal to 3. The next loop is executed normally.

Iterators

An iterable is any data type that is iterable, or can be used in iteration. An iterable can be made into an iterator with the `iter()` function. This means you can use the `next()` function.

```

1 >>> codons = [ 'atg' , 'aaa' , 'agg' ]
2 >>> codons_iterator=iter(codons)
3 >>> next(codons_iterator)
4 'atg'
5 >>> next(codons_iterator)
6 'aaa'
7 >>> next(codons_iterator)
8 'agg'
9 >>> next(codons_iterator)
10 Traceback (most recent call last):
11   File "<stdin>", line 1, in <module>
12 StopIteration

```

An iterator allows you to get the next element in the iterator until there are no more elements. If you want to go through each element again, you will need to redefine the iterator.

List Comprehension

List comprehension is a way to make a list without typing out each element. There are many many ways to use list comprehension to generate lists. Some are quite complex, yet useful.

Here is an simple example:

```

1 dna_list = ['TAGC', 'ACGTATGC', 'ATG', 'ACGGCTAG']
2 lengths = [len(dna) for dna in dna_list]

```

This will produce this list [4, 8, 3, 8]

Using conditions:

```

1 dna_list = ['TAGC', 'ACGTATGC', 'ATG', 'ACGGCTAG']
2 lengths = [len(dna) for dna in dna_list if dna.startswith('A')]

```

This generates the following list: [8, 3, 8]

Using conditions:

```

1 dna_list = ['TAGC', 'ACGTATGC', 'ATG', 'ACGGCTAG']
2 lengths = [len(dna) for dna in dna_list if dna.startswith('A')]

```

This generates the following list: [8, 3, 8]

Here is an example of using mathematical operators to generate a list:

```
1 l = [2 ** x for x in range(100000)]
```

This creates a list of list of the first one hundred thousand powers of two

Dictionaries

Dictionaries are another iterable, like a string and list. Unlike strings and lists, dictionaries are not a sequence, or in other words, the position is not important.

Dictionaries are a collection of key/value pairs. In python, each key is separated from its value by a colon (:), the items are separated by commas, and the whole thing is enclosed in curly braces. An empty dictionary without any items is written with just two curly braces, like this: {}

Each key in a dictionary is unique, while values may not be. The values of a dictionary can be of any type, but the keys must be of an immutable data type such as strings, numbers, or tuples.

Data that is appropriate for dictionaries are two pieces of information that naturally go together, like gene name and sequence.

Key	Value
TP53	GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAGTC
BRCA1	GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA

Creating a Dictionary

```
1 genes = { 'TP53' :  
            'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAGTC '  
            , 'BRCA1' :  
            'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '  
            }
```

Breaking up the key/value pairs over multiple lines make them easier to read.

```

1 genes = {
2     'TP53' :
3         'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGT',
4     'BRCA1' :
5         'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAA'
6 }

```

Accessing Values in Dictionaries

To retrieve a single value in a dictionary use the value's key in this format `dict[key]`. This will return the value at the specified key.

```

1 >>> genes = { 'TP53' :
2     'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC '
3     , 'BRCA1' :
4     'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
5     }
6 >>>
7 >>>
8 >>> genes['TP53']
9 GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC

```

The sequence of the gene TP53 is stored as a value of the key 'TP53'. We can assess the sequence by using the key in this format `dict[key]`

The value can be accessed and passed directly to a function or stored in a variable.

```

1 >>> print(genes['TP53'])
2 GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC
3 >>>
4 >>> seq = genes['TP53'];
5 >>> print(seq)
6 GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC

```

Changing Values in a Dictionary

Individual values can be changed via the key and the assignment operator.

```

1 >>> genes = { 'TP53' :
    'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC '
    , 'BRCA1' :
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
    }
2 >>> print(genes)
3 {'BRCA1':
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA ',
    'TP53':
    'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC '}
4 >>> genes['TP53'] = 'atg'
5 >>> print(genes)
6 {'BRCA1':
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA ',
    'TP53': 'atg'}

```

The contents of the dictionary have changed.

Other assignment operators can also be used to change a value of a dictionary key.

```

1 >>> genes = { 'TP53' :
    'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC '
    , 'BRCA1' :
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
    }
2 >>> genes['TP53'] +=
    'TAGAGCCACCGTCCAGGGAGCAGGTAGCTGCTGGGCTCCGGGGACACTTTGCGTTCGGGCTGGGAGCGTG '
3 >>> print(genes)
4 {'BRCA1':
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA ',
    'TP53':
    'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTCTAGAGCCACCGTGGGCTCCGGGGACACTTTGCGTTCGGGCTGGGAGCGTG '}

```

Here we have used the '+' concatenation assignment operator. This is equivalent to `genes['TP53'] = genes['TP53'] + 'TAGAGCCACCGTCCAGGGAGCAGGTAGCTGCTGGGCTCCGGGGACACTTTGCGTTCGGGCTGGGAGCGTG'`.

Building a Dictionary one Key/Value at a Time

Building a dictionary one key/value at a time is akin to what we just saw when we change a key's value. Normally you won't do this. We'll talk about ways to build a dictionary from a file in a later lecture.

```

1 >>> genes = {}
2 >>> print(genes)
3 {}
4 >>> genes['Brca1'] =
      'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '
5 >>> genes['TP53'] =
      'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC '
6 >>> print(genes)
7 {'Brca1':
   'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA ',
   'TP53':
   'GATGGGATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC '}

```

We start by creating an empty dictionary. Then we add each key/value pair using the same syntax as when we change a value.
`dict[key] = new_value`

Checking That Dictionary Keys Exist

Python generates an error (NameError) if you try to access a key that does not exist.

```

1 >>> print(genes['HDAC'])
2 Traceback (most recent call last):
3   File "<stdin>", line 1, in <module>
4 NameError: name 'HDAC' is not defined

```

Operator	Description
in	key in dict returns True if the key exists in the dictionary
not in	key not in dict returns True if the key does not exist in the dictionary

Because python generates a NameError if you try to use a key that doesn't exist in the dictionary, you probably need to check whether a key exists before trying to use it. The best way to check whether a key exists is to use `in`

```

1 >>> gene = 'TP53'
2 >>> if gene in genes: print('found')
3 ...
4 found
5 >>>
6 >>> if gene in genes:
7 ...     print(genes[gene])
8 ...
9 GATGGGATTGGGGTTTTCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC
10 >>>

```

Sorting Dictionary Keys

If you want to print the contents of a dictionary, you probably want to sort the keys then iterate over the keys with a for loop. Why do you want to sort the keys?

```

1 for key in sorted(genes):
2     print(key, '=>' , genes[key])

```

This will print the same order of keys every time you run your script. This is good for finding bugs and means you don't have to check whether the key exists explicitly because `in sorted(genes)` returns a list of all the keys. It also sorts it.

Dictionary Functions

Function	Description
<code>len(dict)</code>	returns the total number of key/value pairs
<code>str(dict)</code>	returns a string representation of the dictionary
<code>type(variable)</code>	Returns the type or class of the variable passed to the function. If the variable is dictionary, then it would return a dictionary type.

Dictionary Methods

Method	Description
<code>dict.clear()</code>	Removes all elements of dictionary dict
<code>dict.copy()</code>	Returns a shallow copy of dictionary dict. Shallow vs Deep only matters in multidementional datastructures.
<code>dict.fromkeys(seq,value)</code>	Create a new dictionary with keys from seq (python sequence type) and values set to value.
<code>dict.items()</code>	Returns a list of (key, value) tuple pairs
<code>dict.keys()</code>	Returns list of keys
<code>dict.setdefault(key, default = None)</code>	Similar to <code>get()</code> , but will set <code>dict[key] = default</code> if key is not already in dict
<code>dict.update(dict2)</code>	Adds dictionary dict2's key-values pairs to dict

Method	Description
<code>dict.values()</code>	Returns list of dictionary dict's values

Sets

A set is another Python data type. It is essentially a dictionary with keys but no values.

- A set is unordered
- A set is a collection of data with no duplicate elements.
- Common uses include looking for differences and eliminating duplicates in data sets.

Curly braces or the `set()` function can be used to create sets.

Note: to create an empty set you have to use `set()`, not `{}` the latter creates an empty dictionary.

```
1 >>> basket = {'apple', 'orange', 'apple', 'pear', 'orange', 'banana'}
2 >>> print(basket)
3 {'orange', 'banana', 'pear', 'apple'}
```

Look, duplicates have been removed

Test to see if an value is in the set

```
1 >>> 'orange' in basket
2 True
3 >>> 'crabgrass' in basket
4 False
```

The `in` operator works the same with sets as it does with lists and dictionaries

Union, intersection, difference and symmetric difference can be done with sets

```
1 >>> a = set('abracadabra')
2 >>> b = set('alacazam')
3 >>> a
4 {'a', 'r', 'b', 'c', 'd'}
```

Sets contain unique elements, therefore, even if duplicate elements are provided they will be removed.

Set Operators

Difference:

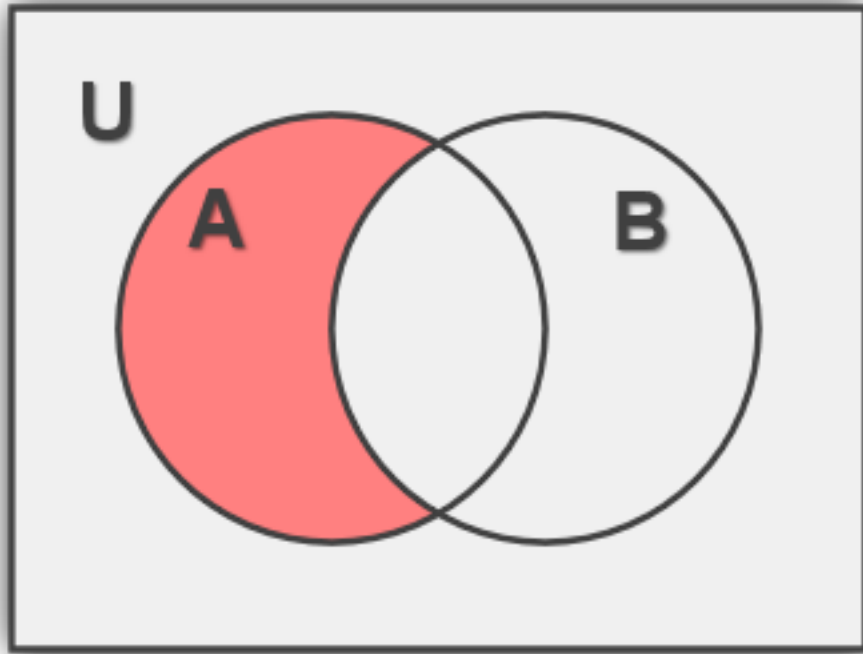


Figure 8: Set Difference

```
1 >>> a - b
2 {'r', 'd', 'b'}
```

This returns the letters that are in a but not in b

Union:

```
1 >>> a | b
2 {'a', 'c', 'r', 'd', 'b', 'm', 'z', 'l'}
```

This returns letters that are in a or b both

Intersection:

```
1 >>> a & b
2 {'a', 'c'}
```

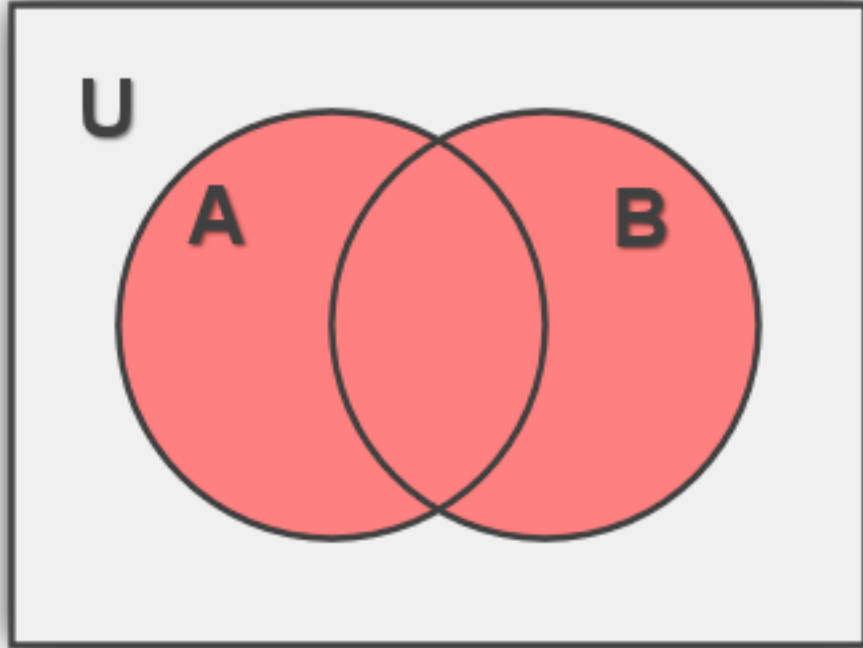


Figure 9: Set Union

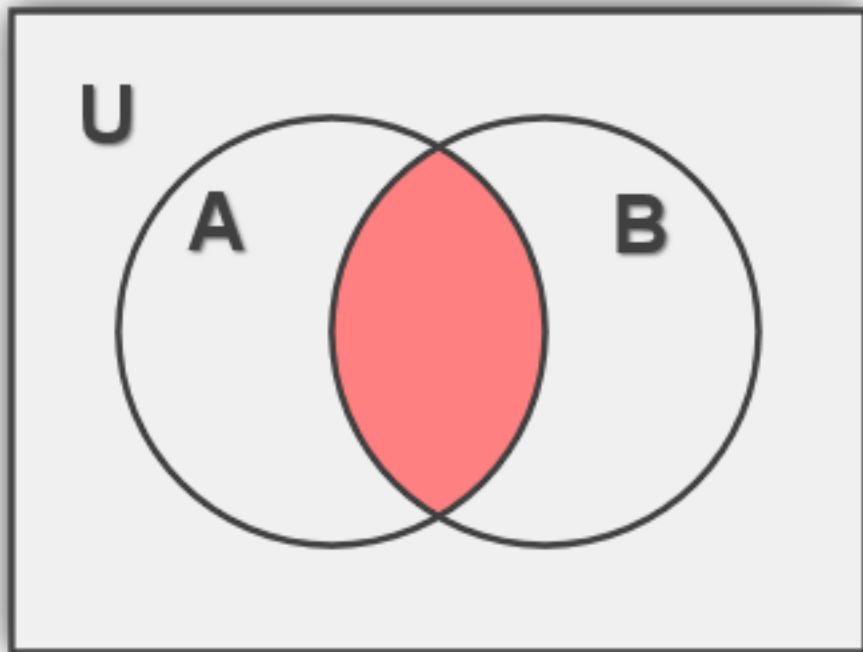


Figure 10: Set Intersection

This returns letters that are in both a and b

Symmetric Difference:

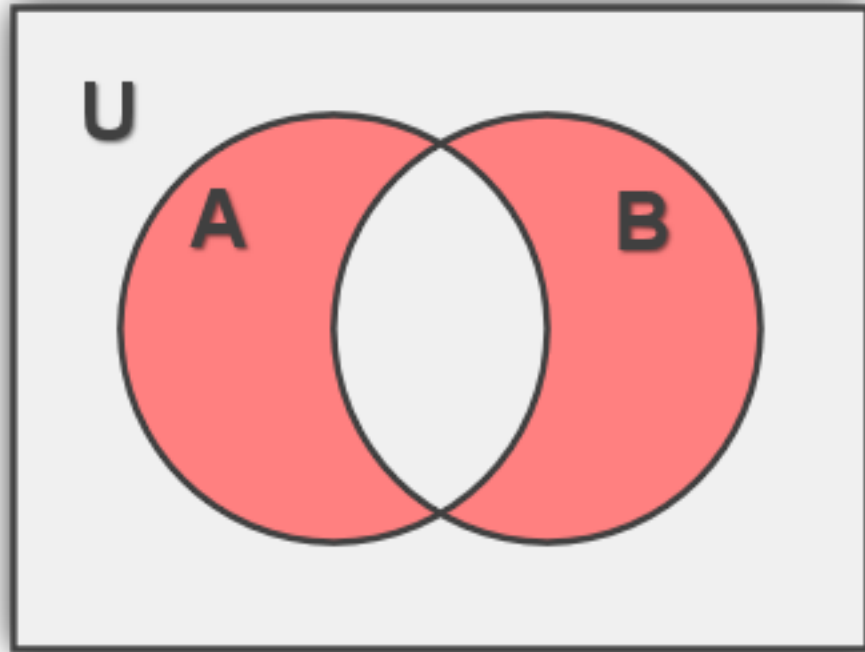


Figure 11: Set Symmetric Difference

```
1 >>> a ^ b
2 {'r', 'd', 'b', 'm', 'z', 'l'}
```

This returns the letters that are in a or b but not in both

Set Functions

Function	Description
<code>all()</code>	returns True if all elements of the set are true (or if the set is empty).
<code>any()</code>	returns True if any element of the set is true. If the set is empty, return False.
<code>enumerate()</code>	return an enumerate object. It contains the index and value of all the items of set as a pair.
<code>len()</code>	return the number of items in the set.
<code>max()</code>	return the largest item in the set.
<code>min()</code>	return the smallest item in the set.

Function	Description
<code>sorted()</code>	returns a new sorted list from elements in the set (does not alter the original set).
<code>sum()</code>	retrune the sum of all elements in the set.

Set Methods

Method	Description
<code>set.add(new)</code>	adds a new element
<code>set.clear()</code>	remove all elements
<code>set.copy()</code>	returns a shallow copy of a set
<code>set.difference(set2)</code>	returns the difference of set and set2
<code>set.difference_update(set2)</code>	removes all elements of another set from this set
<code>set.discard(element)</code>	removes an element from set if it is found in set. (Do nothing if the element is not in set)
<code>set.intersection(sets)</code>	return the intersection of set and the other provided sets
<code>set.intersection_update(sets)</code>	updates set with the intersection of set and the other provided sets
<code>set.isdisjoint(set2)</code>	returns True if set and set2 have no intersection
<code>set.issubset(set2)</code>	returns True if set2 contains set
<code>set.issuperset(set2)</code>	returns True if set contains set2
<code>set.pop()</code>	removes and returns an arbitrary element of set.
<code>set.remove(element)</code>	removes element from a set.
<code>set.symmetric_difference(set2)</code>	returns the symmetric difference of set and set2
<code>set.symmetric_difference_update(set2)</code>	updates set with the symmetric difference of set and set2
<code>set.union(sets)</code>	returns the union of set and the other provided sets
<code>set.update(set2)</code>	update set with the union of set and set2

[Link to Python 4 Problem Set](#)

Python 5

Iterables

@sep some theory and general ideas to solidify previous section

@sep comparison between iterables and sequences, dir()

@smr I mentioned Iterators at the end of loops, just before dictionaries. @smr is an iterator really any method that has a next() method

I/O and Files

I/O stands for input/output. The in and out refer to getting data into and out of your script. It might be a little surprising at first, but writing to the screen, reading from the keyboard, reading from a file, and writing to a file are all examples of I/O.

Writing to the Screen

You should be well versed in writing to the screen. We have been using the `print()` function to do this.

```
1 >>> print ("Hello, PFB2017!")
2 Hello, PFB2017!
```

Remember this example from one of our first lessons?

Reading input from the keyboard

This is something new. There is a function which prints a message to the screen and waits for input from the keyboard. This input can be stored in a variable.

```
1 >>> user_input = input("Type Something Now: ")
2 Type Something Now: Hi
3 >>> print(user_input)
4 Hi
```

All inputted text will be treated as a string. If you are to do math with the input, convert to an int or float first.

Reading from a File

Most of the data we will be dealing with will be contained in files.

The first thing to do with a file is open it. We can do this with the `open()` function. The `open()` function takes the file name, access mode as arguments and returns a file object.

The most common access modes are reading (r) and writing (w).

Open a File

```
1 >>> file_object = open("seq.nt.fa", "r")
```

'file_object' is a name of a variable. This can be anything.

Reading the contents of a file

Now that we have opened a file and created a file object we can do things with it, like read it. Lets read all the contents at once.

Let's cat the contents on the *command line* to see what's in it before we open it in python

```
1 $ cat seq.nt.fa
2 ACAAATACGTTTTGTAAATGTTGTGCTGTTAACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAG
3 ACCGGTTTCCAAAGACAGTCTTCTAATTCCTCATTAGTAATAAGTAAAATGTTTATTGTTGTAGCTCTGG
```

Now, lets print the contents to the screen with Python

```
1 >>> file_object = open("seq.nt.fa", "r")
2 >>> contents = file_object.read()
3 >>> file_name = file_object.name
4 >>> print(file_name)
5 seq.nt.fa
6 >>> print(contents)
7 ACAAATACGTTTTGTAAATGTTGTGCTGTTAACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAG
8 ACCGGTTTCCAAAGACAGTCTTCTAATTCCTCATTAGTAATAAGTAAAATGTTTATTGTTGTAGCTCTGG
9
10 >>> file_object.close()
```

The complete contents can be retrieved with the `read()` method. Another attribute of `file_object` that can be accessed is the file name. Notice the formatting including newlines are maintained when `contents` is printed to the screen. `print()` adds another new line when it is finished printing. It is good practice to close your file. Use the `close()` method.

A for loop can be used to iterate through the `file_object` one line at a time.

```

1 #!/usr/bin/python3
2
3 file_object = open("seq.nt.fa","r")
4 for line in file_object:
5     print(line)

```

Output:

```

1 $ python3 file_for.py
2 ACAAATACGTTTTGTAAATGTTGTGCTGTTAACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAG
3
4 ACCGGTTTCCAAAGACAGTCTTCTAATTCCTCATTAGTAATAAGTAAATGTTTATTGTTGTAGCTCTGG

```

Notice the blank line at after each line we print. `print()` automatically adds a newline and we have a newline at the end of each line in our file. Use `rstrip()` method to remove the newline from each line.

Let's use `rstrip()` method to remove the newline from our file input.

```

1 $ cat file_for_rstrip.py
2 #!/usr/bin/python3
3
4 file_object = open("seq.nt.fa","r")
5 for line in file_object:
6     line = line.rstrip()
7     print(line)

```

`rstrip()` without any parameters returns a string with whitespace removed from the end.

Output:

```

1 $ python3 file_for_rstrip.py
2 ACAAATACGTTTTGTAAATGTTGTGCTGTTAACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAG
3 ACCGGTTTCCAAAGACAGTCTTCTAATTCCTCATTAGTAATAAGTAAATGTTTATTGTTGTAGCTCTGG

```

Writing to a File

Writing to a file is nothing more than opening a file for writing then using the `write()` method.

The `write()` method is like the `print()` function. The biggest difference is that it writes to your file object instead of the screen. `write()` can only take one string type argument.

Let's write a few lines to a file named "writing.txt".

```

1 #!/usr/bin/python3
2
3 fo = open("writing.txt" , "w")
4 fo.write("One line.\n")
5 fo.write("2nd line.\n")
6 fo.write("3rd line" + " has extra text\n")
7 some_var = 5
8 fo.write("4th line has " + str(some_var) + " words\n")
9 fo.close()

```

Output:

```

1 $ python3 file_write.py
2 $ cat writing.txt
3 One line.
4 2nd line.
5 3rd line has extra text
6 4th line has 5 words

```

As you can see, literal strings, return values from operators, return values from functions and variables can be written to a file with the `write()` method.

Now, lets get crazy! Lets read from one file a line at a time. Do something to each line and write the results to a new file.

```

1 #!/usr/bin/python3
2
3 seq_read  = open("seq.nt.fa","r")
4 seq_write = open("nt.counts.txt","w")
5
6 total_nts = 0
7 for line in seq_read:
8     line = line.rstrip()
9     nt_count = len(line)
10    total_nts += nt_count
11    seq_write.write(str(nt_count) + "\n")
12
13 seq_write.write("Total: " + str(total_nts))
14
15 seq_read.close()
16 seq_write.close()

```

Output:


```

1 $ python file_read_write.py
2 $ cat nt.counts.txt
3 71
4 71
5 Total: 142

```

The file we are reading from is named, “seq.nt.fa”
The file we are writing to is named, “nt.counts.txt”
We read each line, calculate the length of each line and print the length
We also create a variable to keep track of the total nt count
Once we go through each line of our file we are reading, we print out the total count of nts
Finally we close each of the files

Building a Dictionary from a File

This is a very common task. It will use a loop, file I/O and a dictionary.

Assume we have a file called “sequence_data.txt” that contains tab-delimited gene names and sequences that looks something like this

```

1 TP53
   GATGGGATTGGGGTTTTCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC
2 BRCA1
   GTACCTTGATTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGTTTCCGTGGCAACGGAAAA

```

How can we read this whole file in to a dictionary?

```

1 #!/usr/bin/python3
2
3 seq_read = open("sequence_data.txt","r")
4 genes = {}
5 for line in seq_read:
6     line = line.rstrip()
7     id,seq = line.split() #split on whitespace
8     genes[id] = seq
9
10 print(genes)

```

Output:

```

1 {'TP53':
   'GATGGGATTGGGGTTTTCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTC ',
  'BRCA1':
   'GTACCTTGATTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGTTTCCGTGGCAACGGAAAA '}

```

Link to Python 5 Problem Set

Python 6

Regular Expressions

Regular Expressions is a language for pattern matching. Many different computer languages incorporate regular expressions as well as some unix commands like `grep` and `sed`. So far we have seen a few functions for finding exact matches in strings, but this is not always sufficient.

Functions that utilize regular expressions allow for non-exact pattern matching.

These specialized functions are not included in the core of Python. We need to import them by typing `import re` at the top of your script

```
1 #!/usr/bin/python3
2
3 import re
```

First we will go over a few examples then go into the mechanics in more detail.

Let's start simple and find an exact match for the EcoRI restriction site in a string.

```
1 >>> dna =
    'ACAAAATACGTTTTGTAAATGTTGTGCTGTTAACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAGGAATTCACCC'
2 >>> if re.search(r"GAATTC",dna):
3 ...     print("Found an EcoRI site!")
4 ...
5 Found an EcoRI site!
6 >>>
```

Since we can search for control characters like a tab (`\t`) it is good to get in the habit of using the raw string function `r` when defining patterns.

Here we used the `search()` function with two arguments, 1) our pattern and 2) the string we want to search.

Let's find out what is returned by the `search()` function.

```
1 >>> dna =
    'ACAAAATACGTTTTGTAAATGTTGTGCTGTTAACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAGGAATTCACCC'
2 >>> found=re.search(r"GAATTC",dna)
3 >>> print(found)
4 <_sre.SRE_Match object; span=(70, 76), match='GAATTC'>
```

Information about the first match is returned

How about a non-exact match. Let's search for a methylation site that has to match the following criteria:
- G or A - followed by C - followed by one of anything or nothing - followed by a G

This could match any of these:

GCAG
GCTG
GCGG
GCCG
GCG
ACAG
ACTG
ACGG
ACCG
ACG

We could test for each of these, or use regular expressions. This is exactly what regular expressions can do for us.

```
1 >>> dna =  
    'ACAAAATACGTTTTGTAAATGTTGTGCTGTTAACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAGGAATTCACCG' 1  
2 >>> found=re.search(r"[GA]C.?G",dna)  
3 >>> print(found)  
4 <_sre.SRE_Match object; span=(7, 10), match='ACG'>
```

Here you can see in the returned information that ACG starts at string position 7 (nt 8). The first position not in the match is at string position 10 (nt 11).

What about other potential matches in our DNA string? We can use `findall()` function to find all matches.

```
1 >>> dna =  
    'ACAAAATACGTTTTGTAAATGTTGTGCTGTTAACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAGGAATTCACCG' 1  
2 >>> found=re.findall(r"[GA]C.?G",dna)  
3 >>> print(found)  
4 ['ACG', 'GCTG', 'ACTG', 'ACCG', 'ACAG', 'ACCG', 'ACAG']
```

`findall()` returns a list of all the pieces of the string that match the regex.

A quick count of all the matching sites can be done by counting the length of the returned list.

```
1 >>> len (re.findall(r"[GA]C.?G",dna))  
2 7
```

There are 7 methylation sites. Here we have another example of nesting. We call the `findall()` function, searching for all the matches of a methylation site. This function returns a list, the list is past to the `len()` function, which in turn returns the number of elements in the list.

Let's talk a bit more about all the new characters we see in the pattern.

The pattern is made up of atoms.

Individual Characters

Atom	Description
a-z, A-Z, 0-9 and some punctuation	These are ordinary characters that match themselves
"."	The dot, or period. This matches any single character except for the newline.

Character Classes

A group of characters that are allowed to be matched one time. There are a few predefined classes, symbol that means a series of characters.

Atom	Description
[]	A bracketed list of characters, like [GA]. This indicates a single character can match any character in the bracketed list.
\d	Digits. Also can be written [0-9]
\w	Word character. Also can be written [A-Za-z_0-9]
\s	White Space character. Also can be written [\t\n\r]
\D	A Non-Digit character. Also can be written [^0-9]
\W	A Non-word character. Also can be written [^A-Za-z_0-9]
\S	A Non-White Space Character. Also can be written [^ \t\n\r]

Anchors

A pattern can be anchored to a region in the string

Atom	Description
^	Matches the beginning of the string
\$	Matches the end of the string
\b	Matches a word boundary. This is between a \w and a \W

Examples:

```
1 g..t
```

matches “gaat”, “goat”, and “gotta get a goat” (twice)

```
1 g[ga][ga]t
```

matches “gaat”, “gttt”, “gatt”, and “gotta get an agatt” (once)

```
1 \d\d\d-\d\d\d\d`
```

matches 376-8380, and 5128-8181 but not 055-98-2818.

```
1 ^\d\d\d-\d\d\d\d`
```

matches 376-8380 and 376-83801 but not 5128-8181.

```
1 ^\d\d\d-\d\d\d\d$
```

only matches telephone numbers

```
1 \bcat
```

matches “cat”, “catsup” and “more catsup please” but not “scat”.

```
1 \bcat\b
```

only text containing the word “cat”.

Quantifiers

Quantifiers quantify how many atoms are to be found. By default an atom matches only once. This behaviour can be modified following an atom with a quantifier.

Quantifier	Description
?	atom matches zero or exactly once
*	atom matches zero or more times
+	atom matches one or more times
{3}	atom matches exactly 3 times
{2,4}	atom matches between 2 and 4 times, inclusive
{4,}	atom matches at least 4 times

Examples:

```
1 goa?t
```

matches “goat” and “got”. Also any text that contains these words.

```
1 g.+t
```

matches “goat”, “goot”, and “grant”, among others.

```
1 g.*t
```

matches “gt”, “goat”, “goot”, and “grant”, among others.

```
1 ^\d{3}-\d{4}$
```

matches US telephone numbers (no extra text allowed).

Something to think about.

1) What would be a pattern to recognize an email address? 2) What would be a pattern to recognize the ID portion of a sequence record in a FASTA file?

Variables and Patterns

Variables can be used to store patterns.

```

1 >>> pattern = r"[GA]C.?G"
2 >>> len (re.findall(pattern,dna))
3 7

```

In this example we stored our methylation pattern in the variable named 'pattern' and used it as the first argument to findall.

Either Or

A pipe '|' can be used to indicate that either the pattern before or after the '|' can match. Enclose the two options in parenthesis.

```

1 big bad (wolf|sheep)

```

This pattern must match a string that contains "big" followed by a space followed by "bad" followed by a space followed by *either* "wolf" or "sheep" This would match, "big bad wolf" Or "big bad sheep"

Something to think about. 1) What would a pattern to match 'ATG' followed by a C or a T look like?

Subpatterns

Subpatterns, or parts of the pattern enclosed in parenthesis can be extracted and stored for later use.

```

1 Who's afraid of the big bad w(+)f

```

This pattern has only one subpattern (+)

You can combine parenthesis and quantifiers to quantify entire subpatterns

```

1 Who's afraid of the big (bad )?wolf\?

```

This matches "Who's afraid of the big bad wolf?". As well as "Who's afraid of the big wolf?". The 'bad' is optional, it can be present 0 or 1 times in our string. This also shows how to literally match special characters. Use a \" in to escape them.

Something to think about: How would you create a pattern to capture the ID in a sequence record of a FASTA file in a subpattern.

Example FASTA sequence record.


```
1 >ID Optional Description
2 SEQUENCE
3 SEQUENCE
4 SEQUENCE
```

Using Subpatterns Inside the Regular Expression Match

This is helpful when you want to find a subpattern and then match the contents again

Once a subpattern matches, you can refer to it within the same regular expression. The first subpattern becomes \1, the second \2, the third \3, and so on.

```
1 Who's afraid of the big bad w(.)\1f
```

This would match “Who’s afraid of the big bad woof” “Who’s afraid of the big bad weef” “Who’s afraid of the big bad waaf” But Not, “Who’s afraid of the big bad wolf” Or, “Who’s afraid of the big bad wife”

In a similar vein,

```
1 \b(\w+)s love \1 food\b
```

This pattern will match “dogs love dog food” But not “dogs love monkey food”.

```

1 >>> str="dogs love dog food"
2 >>> found=re.search(r"\b(\w+)s love \1 food\b", str)
3 >>> if found:
4 ...     print(found.group(1))
5 ... else:
6 ...     print("Not Found")
7 ...
8 dog
9 >>> str="cats love cat food"
10 >>> found=re.search(r"\b(\w+)s love \1 food\b", str)
11 >>> if found:
12 ...     print(found.group(1))
13 ... else:
14 ...     print("Not Found")
15 ...
16 cat
17 >>> str="dogs love monkey food"
18 >>> found=re.search(r"\b(\w+)s love \1 food\b", str)
19 >>> if found:
20 ...     print(found.group(1))
21 ... else:
22 ...     print("Not Found")
23 ...
24 Not Found

```

This pattern matches:

- “dogs love dog food” - “cats love cat food” - but not “dogs love monkey food”

Subpatterns and Greediness

By default, regular expressions are “greedy”. They try to match as much as they can. Use the quantifier ‘?’ to make the match not greedy. The not greedy match is called ‘lazy’

```

1 >>> str = 'The fox ate my box of doughnuts '
2 >>> found = re.search(r"(f.+x)",str)
3 >>> print(found.group(1))
4 fox ate my box

```

The pattern `f.+x` does not match what you might expect, it matches past ‘fox’ all the way out to ‘fox ate my box’.

The ‘.’ is greedy. As many characters as possible are found that are between the ‘f’ and the ‘x’.

Let’s make this match lazy by using ‘?’

```

1 >>> found = re.search(r"(f.+?x)",str)
2 >>> print(found.group(1))
3 fox

```

The match is now lazy and will only match 'fox'

Using Subpatterns Outside the Regular Expression Match

Using the captured subpattern in code that follows the regular expression.

Outside the regular expression match statement, the matched subpatterns can be access with the `group()` method.

Example:

```

1 >>> dna =
    'ACAAAATACGTTTTGTAAATGTTGTGCTGTTAACACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAGGAATTCACCC
2 >>> found=re.search( r"(.{50})TATTAT(.{25})" , dna )
3 >>> upstream = found.group(1)
4 >>> print(upstream)
5 TCTAATTCCTCATTAGTAATAAGTAAAATGTTTATTGTTGTAGCTCTGGA
6 >>> downstream = found.group(2)
7 >> print(downstream)
8 CCGGTTTCCAAAGACAGTCTTCTAA

```

- 1) This pattern will recognize a consensus transcription start site (TATTAT)
- 2) And store the 50 base pairs upstream of the site
- 3) And the 25 base pairs downstream of the site

If you want to find the upstream and downstream sequence of ALL 'TATTAT' sites, use the `findall()` function.

```

1 >>>
    dna="ACAAAATACGTTTTGTAAATGTTGTGCTGTTAACACTGCAAATAAACTTGGTAGCAAACACTTCCAAAAGGAATTC
2 >>> found = re.findall( r"(.{50})TATTAT(.{25})" , dna )
3 >>> print(found)
4 [('TCTAATTCCTCATTAGTAATAAGTAAAATGTTTATTGTTGTAGCTCTGGA ',
    'CCGGTTTCCAAAGACAGTCTTCTAA '),
    ('TCTAATTCCTCATTAGTAATAAGTAAAATGTTTATTGTTGTAGCTCTGGA ',
    'CCGGTTTCCAAAGACAGTCTTCTAA ')]

```

The subpatterns are stored in tuples within a list. More about this type of datastructure later.

Another option for retrieving the upstream and downstream subpatterns is to put the `findall()` in a for loop

```
1 >>>
   dna="ACAAAATACGTTTTGTAAATGTTGTGCTGTTAACACTGCAAATAAACTTGGTAGCAAACACTTCCTCAAAAGGAATT"
2 >>> for (upstream, downstream) in re.findall(
   r"(.{50})TATTAT(.{25})" , dna ):
3 ...     print("upstream:" , upstream)
4 ...     print("downstream:" , downstream)
5 ...
6 upstream: TCTAATTCCTCATTAGTAATAAGTAAAAATGTTTATTGTTGTAGCTCTGGA
7 downstream: CCGGTTTCCAAAGACAGTCTTCTAA
8 upstream: TCTAATTCCTCATTAGTAATAAGTAAAAATGTTTATTGTTGTAGCTCTGGA
9 downstream: CCGGTTTCCAAAGACAGTCTTCTAA
```

- 1) This code executes the `findall()` function once
- 2) The subpatterns are returned
- 3) The subpatterns are stored in the variables `upstream` and `downstream`
- 4) The for block of code is executed
- 5) The `findall()` searches again
- 6) A match is found
- 7) New subpatterns are returned
- 8) The for block of code gets executed again
- 9) The `findall()` searches again, but no match is found
- 10) The for loop ends

One other way to get this done is with the `finditer()` function in a for loop

```
1 >>>
   dna="ACAAAATACGTTTTGTAAATGTTGTGCTGTTAACACTGCAAATAAACTTGGTAGCAAACACTTCCTCAAAAGGAATT"
2 >>> for match in re.finditer(r"(.{50})TATTAT(.{25})" , dna):
3 ...     print("upstream:" , match.group(1))
4 ...     print("downstream:" , match.group(2))
5 ...
6 upstream: TCTAATTCCTCATTAGTAATAAGTAAAAATGTTTATTGTTGTAGCTCTGGA
7 downstream: CCGGTTTCCAAAGACAGTCTTCTAA
8 upstream: TCTAATTCCTCATTAGTAATAAGTAAAAATGTTTATTGTTGTAGCTCTGGA
9 downstream: CCGGTTTCCAAAGACAGTCTTCTAA
```

- 1) This code executes `finditer()` function once.
- 2) The match object is returned. A match object will have all the information about the match
- 3) In the for block we call the `group()` method on the first match object returned
- 4) We print out the first and second subpattern using the `group()` method
- 5) The `finditer()` function is executed a second time and a match is found
- 6) The second match object is returned
- 7) The second subpatterns are retrieved from the match object using the `group()` method
- 8) The `finditer()` function is executed again, but no matches found, so the loop ends

FYI: `match()` function is another regular expression function that looks for patterns. It is similar to `search()` but it only looks at the beginning of the string for the pattern while `search()` looks in the entire string. Usually `search()` and `findall()` will be more useful.

Practical Example: Codons

Extracting codons from a string of DNA can be accomplished by using a subpattern in a `findall()` function. Remember the `findall()` function will return a list of the matches.

```
1 >>> dna = 'GTTGCCTGAAATGGCGGAACCTTGAA'
2 >>> codons = re.findall(r"(.{3})",dna)
3 >>> print(codons)
4 ['GTT', 'GCC', 'TGA', 'AAT', 'GGC', 'GGA', 'ACC', 'TTG']
```

Or you can use a for loop to do something to each match.

```
1 >>> for codon in re.findall(r"(.{3})",dna):
2 ...     print(codon)
3 ...
4 GTT
5 GCC
6 TGA
7 AAT
8 GGC
9 GGA
10 ACC
11 TTG
12 >>>
```

`finditer()` would also work in this for loop. Each codon can be accessed by using the `group()` method.

Truth and Regular Expression Matches

The `search()`, `match()`, `findall()`, and `finditer()` can be used in conditional tests. If a match is not found an empty list or 'None' is returned. These both are False.

```
1 >>> found=re.search( r"(.{50})TATTATZ(.{25})" , dna )
2 >>> if found:
3 ...     print("found it")
4 ... else:
5 ...     print("not found")
6 ...
7 not found
8 >>> print(found)
9 None
```

None is False so the else block is executed and “not found” is printed

Using Regular expressions in substitutions

Earlier we went over how to find an exact pattern and replace it using the `replace()` method. To find a pattern and make a replacement the regular expression `sub()` function is used. This function takes the pattern, the replacement, the string to be searched, the number of times to do the replacement, and flags.

```
1 >>> str = "Who's afraid of the big bad wolf?"
2 >>> re.sub(r'w.+f' , 'goat', str)
3 "Who's afraid of the big bad goat?"
4 >>> print(str)
5 Who's afraid of the big bad wolf?
```

The `sub()` function returns “Who’s afraid of the big bad goat?”
The value of variable `str` has not been altered
The new string can be stored in a new variable for later use.

Let’s save the returned new string in a variable

```

1 >>> str = "He had a wife."
2 >>> new_str = re.sub(r'w.+f' , 'goat', str)
3 >>> print(new_str)
4 He had a goate.
5 >>> print(str)
6 He had a wife.

```

The characters between 'w' and 'f' have been replaced with 'goat'.
The new string is saved in new_str

Using subpatterns in the replacement

Sometimes you want to find a pattern and use it in the replacement.

```

1 >>> str = "Who's afraid of the big bad wolf?"
2 >>> new_str = re.sub(r"(\w+) (\w+) wolf" , r"\2 \1 wolf" , str)
3 >>> print(new_str)
4 Who's afraid of the bad big wolf?

```

We found two words before 'wolf' and swapped the order. \2 refers to the second subpattern
\1 refers to the first subpattern

Something to think about.

How would you use regular expressions to find all occurrences of 'ATG' and replace with '-M-' in this sequence
'GCAGAGGTGATGGACTCCGTAATGGCCAAATGACACGT'?

[Link to Python 6 Problem Set](#)

Python 7

Functions

@sep: good place for list comprehensions > **def**, **return**, **positional arguments**, **keyword arguments**, ***arguments** ****keywords** > **lambda** e.g. > **squares** = **list**(**map**(**lambda** x: x**2, **range**(10)))

Writing your own functions is a way to group lines of code together. Why would you want to group lines of code together?

1. Let's say that you have a few lines of code that are required to produce some functionality that you reuse in your code. You can group these all together and have that functionality become one line. You can then just use one line of code where ever you need this functionality.
2. You can simplify the main body of your code if you have one human readable and informative line to represent many lines of code.

Let's see some examples:

Code to find AT content:

```
1 dna =  
    'GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA '  
2 a_count = dna.count('A')  
3 t_count = dna.count('T')  
4 dna_len = len(dna)  
5 at_content = (a_count + t_count) / dna_len  
6 print(at_content)
```

Here are 5 lines of code to do one thing, calculate AT content in a DNA string.

Creating/Defining a Function to Find AT Content:

```
1 def calculate_at_content(dna):  
2     a_count = dna.count('A')  
3     t_count = dna.count('T')  
4     dna_len = len(dna)  
5     at_content = (a_count + t_count) / dna_len  
6     return at_content
```

Here is a custom function that you can use like a built in python function

Using/Running/Calling Your function:


```

1 dna_string =
  "GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA "
2 print(calculate_at_content(dna_string))

```

This code will print “0.4857142857142857” to the screen. You can save this returned value in a variable to use later in your code. As you can see we now have one line, that is very explicit in the functionality, replace 5 lines of code.

The details

1. All functions must be defined.
 - We do this with the word **def**.
2. The function must have a name.
 - “calculate_at_content” is our name
3. Functions define the variables required in the function body. These are placeholders for the arguments that get passed into our function. Functions don’t need to have arguments but they always need ‘()’
 - We do this there “(dna)”.
 - Our function needs one variable, and we can refer to it in the function body as ‘dna’.
 - The variable ‘dna’ and its contents are only visible within the body of the function
4. The first line of the function must end with a ‘:’
5. The function body needs to be indented.
6. Do your work.
7. You can return a value, but this is not required
 - We return the value of ‘at_content’

Naming Arguments

You can name your argument variables anything you want, you just have to use the same term within the function body. Let’s use ‘string’ instead of ‘dna’

```

1 def calculate_at_content(string):
2     a_count = string.count('A')
3     t_count = string.count('T')
4     dna_len = len(string)
5     at_content = (a_count + t_count) / dna_len
6     return at_content

```

Now everywhere we had previously used ‘dna’ we need to use ‘string’

Keyword Arguments

Arguments can be named and these names can be used when the function is called. This name is called a ‘keyword’

```

1 dna_string =
  "GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA "
2 print(calculate_at_content(dna_string))
3 print(calculate_at_content("GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA "))
4 print(calculate_at_content(dna="GTACCTTGATTTCGTATTCTGAGAGGCTGCTGCTTAGCGGTAGCCCCTTGGTTTCCGTGGCAACGGAAAA "))
5 print(calculate_at_content(dna=dna_string))

```

This code will print “0.4857142857142857” (4x) to the screen. As you can see, we supplied the function arguments in a few different ways: 1) a variable, 2) a literal, 3) a literal with a keyword, and 4) a variable with a keyword. The keyword must be the same as the defined function argument. If a function has multiple arguments, using the keyword allows for calling the function with the arguments in any order.

Default Values for Arguments

When you define your function you can also define default values for arguments

```

1 def calculate_at_content(dna='ATGC'):
2     a_count = dna.count('A')
3     t_count = dna.count('T')
4     dna_len = len(dna)
5     at_content = (a_count + t_count) / dna_len
6     return at_content

```

If you call the function with no arguments, the default will be used. In this case a default is pretty useless, and the function will return ‘0.5’ if called without providing a DNA sequence.

lambda

Lambda expressions can be used when your function only contains one line. It is a special kind of custom function that you don’t have to give a name. Lambda expressions are useful because some functions can take a lambda as an argument.

Here is a one line custom function, like the functions we have already talked about:

```

1 def get_first_codon(dna):
2     return dna[0:3]
3
4 print(get_first_codon('ATGTTT'))

```

This will return ‘ATG’

Here is the same function written as a lambda

```
1 get_first_codon = lambda dna : dna[0:3]
2 print(get_first_codon('ATGTTT'))
```

This also returns 'ATG'. lambdas need to have only one line. They do not require the usage of 'return' to have the result returned.

Lambda expressions can be used as arguments in functions:

```
1 print((lambda dna : dna[0:3])('ATGTTT'))
```

This prints 'ATG'. In this example, the code works, but might not be the best way to get this job done since it obfuscates the code, or in other words, makes the code harder to understand. There are cases where using a lambda is the only way to get a task done.

Scope

Most all variables can be seen, retrieved, and used anywhere in your code. This is called global. One of the only non-global or local variables are those used as function arguments.

```
1 #!/usr/bin/python3
2 x = 100
3 y = 20;
4 if x > y:
5     z = 10
6     x = 30
7     print("x (inside if block):", x)
8     print("y (inside if block):", y)
9     print("z (inside if block):", z)
10
11 print("x (outside if block):", x)
12 print("y (outside if block):", y)
13 print("z (outside if block):", z)
```

Let's Run it:

```

1 $ python scripts/scope.py
2 x (inside if block): 30
3 y (inside if block): 20
4 z (inside if block): 10
5 x (outside if block): 30
6 y (outside if block): 20
7 z (outside if block): 10

```

Local Variables

Function argument variables are local and therefore can only be accessed from within the function block.

```

1 #!/usr/bin/python3
2
3 def scope_function(x):
4     print("x (inside function block):", x)
5     x = 5
6     print("x (inside function block):", x)
7     print("y (inside function block):", y)
8     print("z (inside function block):", z)
9
10 x = 100
11 y = 20;
12 if x > y:
13     z = 10
14     x = 30
15     print("x (inside if block):", x)
16     print("y (inside if block):", y)
17     print("z (inside if block):", z)
18
19 print("x (outside if block):", x)
20 print("y (outside if block):", y)
21 print("z (outside if block):", z)
22
23 scope_function(500)
24
25 print("x (outside if block after function call):", x)

```

Here we have added a function with an argument named 'x'. This variable exists only within the function. It does not matter that there is a variable of the same name outside the function block.

Let's run it:

```

1 $ python3 scripts/scope_w_function.py
2 x (inside if block): 30
3 y (inside if block): 20
4 z (inside if block): 10
5 x (outside if block): 30
6 y (outside if block): 20
7 z (outside if block): 10
8 x (inside function block): 1
9 x (inside function block): 5
10 y (inside function block): 20
11 z (inside function block): 10
12 x (outside if block after function call): 30

```

As you can see, x is 30 inside and outside the if block, but inside the function x is 1, then x is 5. Once we have completed the function call x is 30 again. Variables within a function block are local to that block.

Global

You can make a local variable global with the function `global()`. Now a variable defined in a function can be visible to the rest of the code. It is best not to define any variables as global until you know you need it. Most of the time the default scope is what you want.

Here is an example use of `global()`. This can be a bit mind bending, just knowing this exists is probably good enough for right now.

```

1 #!/usr/bin/python3
2
3 def scope_function():
4     global var_1
5     var_1 = "I say hello"
6     var_2 = "You say good-bye"
7     print("var_1 (inside function block):", var_1)
8     print("var_2 (inside function block):", var_2)
9
10 var_1 = "Hello, hello"
11 var_2 = "Good-Bye"
12 print("var_1 (outside function block before function call):", var_1)
13 print("var_2 (outside function block before function call):", var_2)
14 scope_function()
15 print("var_1 (outside function block after function call):", var_1)
16 print("var_2 (outside function block after function call):", var_2)

```

By default all variables are global, except for those within functions. Here we mess with the default behavior by defining 'var_1' as global. We the value of 'var_1' within the function will the value outside the function block.

```

1 $ python3 scripts/scope_global.py
2 var_1 (outside function block before function call): Hello, hello
3 var_2 (outside function block before function call): Good-Bye
4 var_1 (inside function block): I say hello
5 var_2 (inside function block): You say good-bye
6 var_1 (outside function block after function call): I say hello
7 var_2 (outside function block after function call): Good-Bye

```

If we did not make 'var_1' global the value of 'var_' outside the function block would have been "Hello, hello"

Namespaces

@smr I don't think we need this

Modules

@sep > dir() > import > from bla import boo

The standard library: os, sys, glob, shutil, math, random, statistics (scipy), zlib, sqlite3 (DB access)? urllib.request

os.path

`os.path` has common utilities for working file paths (filenames and directories). A path is either a relative or absolute list of directories (often ending with a filename) that tells you where to find a file or directory.

function	description
<code>os.path.basename(path)</code>	what's the last element of the path? Note <code>/home/tmp/</code> returns <code>' '</code> , rather than <code>tmp</code>
<code>os.path.dirname(path)</code>	what's the directory the file is in?
<code>os.path.exists(path)</code>	does the path exist?
<code>os.path.getsize(path)</code>	returns path (file) size in bytes or error
<code>os.path.isfile(path)</code>	does the path point to a file?
<code>os.path.isdir(path)</code>	does the path point to a directory?

os.system

Run a system command from python. Replaced by subprocess

```

1 import os
2 os.system("ls -l")

```

subprocess

```
1 import subprocess
2 run(["ls", "-l"]) # same as running ls -l on the command line
```

more complex than `os.system()`. You need to specify where input and output go

```
“rtn =
“
```

Capturing output from a shell pipeline

```
ls -l | grep amanda
```

becomes this ‘shortcut’ which will capture the output of the two unix commands in the variable `output`

```
1 import subprocess
2 output = subprocess.check_output('ls -l | grep amanda', shell = True)
```

This is better than alternatives with `subprocess.run()`.

`output` contains a bytes object (more or less a string of ASCII character encodings)

```
1 b'-rw-r--r-- 1 amanda staff      161952 Oct  2 18:03
   test.subreads.fa\n-rw-r--r-- 1 amanda staff      126 Oct  2
   13:23 test.txt\n'
```

You can covert by decoding the bytes object into a string

```
1 >>>output.decode('utf-8')
2 '-rw-r--r-- 1 amanda staff      161952 Oct  2 18:03
   test.subreads.fa\n-rw-r--r-- 1 amanda staff      126 Oct  2
   13:23 test.txt\n'
```

Capturing output the long way (for a single command)

Let’s assume that `ls -l` generates some output something like this

```
1 total 112
2 -rw-r--r-- 1 amanda staff      69 Jun 14 17:41 data.cfg
3 -rw-r--r-- 1 amanda staff    161952 Oct  2 18:03
   test.subreads.fa
4 -rw-r--r-- 1 amanda staff      126 Oct  2 13:23 test.txt
```

How do we run `ls -l` in python and capture the output (stdout)?

```
1 import subprocess
2 rtn = subprocess.run(['ls','-l'], stdout=subprocess.PIPE )
3 bytes = rtn.stdout
4 stdout = bytes.decode('utf-8')
5 # something like
6 lines = stdout.splitlines()
```

`lines` now contains elements from every line of the `ls -l` output, including the header line, which is not a file

```
1 >>> lines[0]
2 'total 112'
3 >>> lines[1]
4 '-rw-r--r--  1 amanda  staff      69 Jun 14 17:41 data.cfg'
```

sys

A couple of useful variables for beginners. Many more advanced system parameters and settings that we are not covering here.

function	description
<code>sys.argv</code>	list of command line parameters
<code>sys.path</code>	where python should look for modules

re

see notes on regular expressions

collections

`from collections import deque`

copy

`copy.copy()`

and

`copy.deepcopy()`

math

function	description
<code>math.exp()</code>	e^{**x}
<code>math.log2()</code>	
<code>math.log10()</code>	

function	description
math.sqrt()	
math.sin()	
math.pi(), math.e()	constants
etc	

see also numpy

random

statistics

glob

argparse

Great (and overly-complicated) tool for parsing command line arguments and automatically generating help messages for scripts (very handy!). Here's a simple script that explains a little of what it does.

```

1 #!/usr/bin/python3
2 import argparse
3 parser = argparse.ArgumentParser(description="A test program that
   reads in some number of lines from an input file. The output can
   be screen or an output file")
4 # we want the first argument to be the filename
5 parser.add_argument("file", help="path to input fasta filename")
6 # second argument will be line number
7 # default type is string, need to specify if expecting an int
8 parser.add_argument("lines", type=int, help="how many lines to
   print")
9 # optional outfile argument specified with -o or --out
10 parser.add_argument("-o", "--outfile", help="optional: supply
   output filename, otherwise write to screen", dest='out')
11 args = parser.parse_args()
12 # arguments appear in args
13 filename = args.file
14 lines = args.line
15 if args.out:
16     print("writing output to", args.out)

```

Many more modules in many areas

time, HTML, XML, email, CGI, sockets, audio, GUIs with Tk, debugging, testing, unix utils

Also, non-core: BioPython for bioinformatics, Numpy for mathematics, statistics

[Link to Python 7 Problem Set](#)

Python 8

Exception Handling

@sep > try/except/finally (see with/as)

There are a few different types of errors when coding. Syntax errors, logic errors, and exceptions. You have probably encountered all three. Syntax and logic errors are issues you need to deal with while coding. An exception is a special type of error that can be informative and used to write code to respond to this type of error. This is especially relevant when dealing with user input. What if they don't give you any, or it is the wrong kind of input. We want our code to be able to detect these types of errors and respond accordingly.

```
1 #!/usr/bin/python3
2
3 import sys
4 file = sys.argv[1]
5
6 print("User provided file:" , file)
```

This code takes user provided input and prints it

Run it.

```
1 $ python scripts/exceptions.py test.txt
2 User provided file: test.txt
```

What happens if the user does not provide any input and we try to print it?

```
1 $ python scripts/exceptions.py
2 Traceback (most recent call last):
3   File "scripts/exceptions.py", line 4, in <module>
4     file = sys.argv[1]
5 IndexError: list index out of range
```

We get an IndexError exception

We have already seen quite a few Exceptions, here are some:

- ValueError: math domain error - AttributeError: 'list' object has no attribute 'rstrip' - SyntaxError: EOL while scanning string literal - NameError: name 'GGTCTAC' is not defined - SyntaxError: Missing parentheses in call to 'print' - AttributeError: 'int' object has no attribute 'lower' - IndexError: list assignment index out of range - NameError: name 'HDAC' is not defined

We can use the exception to our advantage to help out our users. We can use a try/except condition to look for exceptions and to do something if we do not have an exception and do something different if we do have

an exception.

```
1 #!/usr/bin/python3
2 import sys
3
4 file = ''
5 try:
6     file = sys.argv[1]
7     print("User provided file:" , file)
8 except:
9     print("Please provide a file name")
```

We need to “try” to get a user provided argument. If we are successful then we can print it out. If we try and fail, we execute the code in the except portion of our try/except and print that we need a file name.

Let’s run it WITH user input

```
1 $ python3 scripts/exceptions_try.py test.txt
2 User provided file: test.txt
```

It runs as expected

Let’s run it WITHOUT user input

```
1 $ python scripts/exceptions_try.py
2 Please provide a file name
```

Yeah, the user is informed that they need to provide a file name to the script

What if the user provides input but it is not a valid file or the path is incorrect? Or if you want to check to see if the user provided input as well as if it can open the input.

We can add multiple exception tests as well as specify what kind of exception for which the exception block should be executed.

```

1 import sys
2
3 file = ''
4 try:
5     file = sys.argv[1]
6     print("User provided file name:" , file)
7     FASTA = open(file, "r")
8     for line in FASTA:
9         line = line.rstrip()
10        print(line)
11 except IndexError:
12     print("Please provide a file name")
13 except IOError:
14     print("Can't find file:" , file)

```

Here we test for an `IndexError` and a `IOError`. The `IndexError` occurs when we try to access a list element that does not exist. The `IOError` happens when we try to access a file that does not exist.

Let's run it with no input.

```

1 $ python scripts/exceptions_try_files.py test.txt
2 User provided file name: test.txt
3 Can't find file: test.txt

```

This informs the user that they did provide input but that the file listed can not be found.

Let's run it with no input

```

1 $ python scripts/exceptions_try_files.py
2 Please provide a file name

```

This informs the user that they need to provide a file.

try/except/else/finally

Let's summarize what we have covered and add on `else` and `finally`.

```

1 try:
2     # try block is executed until an exception is raised
3 except _ExceptionType_:
4     # if there is an exception of "ExceptionType" this block will be
       executed
5     # there can be more than one except block, just like an elif
6 except:
7     # if there are any exceptions that are not of _ExceptionType_ this
       except block will be executed
8 else:
9     # the else block is executed after the try block has been completed
10 finally:
11     # the finally block is executed no matter if there are exceptions
       or not.

```

Getting more information about an exception

Some exceptions can be thrown for multiple reasons, for example, `ErrorIO` will occur if the file does not exist as well you don't have permissions to read it. We can get more information by viewing the contents of our Exception Object. Yes, an exception is an object too! To access the object use `as` and supply a variable name, like 'ex'

```

1 file = ''
2 try:
3     file = sys.argv[1]
4     print("User provided file name:" , file)
5     FASTA = open(file, "r")
6     for line in FASTA:
7         line = line.rstrip()
8         print(line)
9 except IndexError:
10    print("Please provide a file name")
11 except IOError as ex:
12    print("Can't find file:" , file , ': ' , ex.strerror )

```

Here we added `except IOError as ex` and now we can get the 'strerror' message from `ex`.

Run it.

```

1 $ python scripts/exceptions_try_files_as.py test.txt
2 User provided file name: test.txt
3 Can't find file: test.txt : No such file or directory

```

Now we know that this file name or path is not valid

Raising an Exception

We can call or raise exceptions too!! This is accomplished by using a `raise` statement.

1. First, create a new Exception Object, i.e., `ValueError()`
2. Use the Exception Object in a Raise statment `raise ValueError('your message')`

Let's raise an exception if the file name does not end in 'fa'

```
1 import sys
2
3 file = ''
4 try:
5     file = sys.argv[1]
6     print("User provided file name:" , file)
7     if not file.endswith('.fa'):
8         raise ValueError("Not a FASTA file")
9     FASTA = open(file, "r")
10    for line in FASTA:
11        print(line)
12 except IndexError:
13     print("Please provide a file name")
14 except IOError as ex:
15     print("Can't find file:" , file , ': ' , ex.strerror )
```

Here we raise a known exception, 'ValueError', if the file does not end with (uses `endswith()` method).

Let's run it.

```
1 $ python scripts/exceptions_try_files_raise.py test.txt
2 User provided file name: test.txt
3 Traceback (most recent call last):
4   File "scripts/exceptions_try_files_raise.py", line 10, in <module>
5     raise ValueError("Not a FASTA file")
6 ValueError: Not a FASTA file
```

Our exception get's raised, now lets do something with it.

```

1 import sys
2
3 file = ''
4 try:
5     file = sys.argv[1]
6     print("User provided file name:" , file)
7     if not file.endswith('.fa'):
8         raise ValueError("Not a FASTA file")
9     FASTA = open(file, "r")
10    for line in FASTA:
11        print(line)
12 except IndexError:
13    print("Please provide a file name")
14 except ValueError:
15    print("File needs to be a FASTA file and end with .fa")
16 except IOError as ex:
17    print("Can't find file:" , file , ': ' , ex.strerror )

```

Here we created an except to catch any ValueError

Let's Run it.

```

1 $ python scripts/exceptions_try_files_raise_value.py test.txt
2 User provided file name: test.txt
3 File needs to be a FASTA file and end with .fa

```

We get a great error message now.

But what if there is another ValueError, how can we tell if it is do to the FASTA file extension or not. We can create our own custom exceptoin. We will need to create a new class of exception. Below is the syntax to do this.


```

1 import sys
2
3 class NotFASTAError(Exception):
4     pass
5
6
7 file = ''
8 try:
9     file = sys.argv[1]
10    print("User provided file name:" , file)
11    if not file.endswith('.fa'):
12        raise NotFASTAError("Not a FASTA file")
13    FASTA = open(file, "r")
14    for line in FASTA:
15        print(line)
16 except IndexError:
17    print("Please provide a file name")
18 except NotFASTAError:
19    print("File needs to be a FASTA file and end with .fa")
20 except IOError as ex:
21    print("Can't find file:" , file , ': ' , ex.strerror )

```

Here we created a new class of exception called 'NotFASTAError'. Then we raised this new exception.

Let's Run it.

```

1 $ python scripts/exceptions_try_files_raise_try.py test.txt
2 User provided file name: test.txt
3 File needs to be a FASTA file and end with .fa

```

Our new class of exception, NotFASTAError, works just like the built in exceptions.

Datastructures

@sep: how you initialize from a file > import json

Sometimes a *simple* list or dictionary just doesn't do what you want. Sometimes you need to organize data in a more *complex* way. You can nest any data type inside any other type. This lets you build multidimensional data tables easily.

Two-dimensional lists

Often times a list of lists, often called a matrix are important for organizing and accessing data

Here's a way to make a 3 x 3 table of values.

```
1 >>> M = [[1,2,3], [4,5,6], [7,8,9]]
2 >>> M[1] # second row (starts with index 0)
3 [4,5,6]
4 >>> M[1][2] # second row, third element
5 6
```

Here's a way to store sequence alignment data:

Four sequences aligned:

```
1 AT-TG
2 AATAG
3 T-TTG
4 AA-TA
```

The alignment in a list of lists.

```
1 aln = [['A', 'T', '-', 'T', 'G'],
2 ['A', 'A', 'T', 'A', 'G'],
3 ['T', '-', 'T', 'T', 'G'],
4 ['A', 'A', '-', 'T', 'A']]
```

Get an the full length of one sequence:

```
1 >>> seq = aln[2]
2 >>> seq
3 ['T', '-', 'T', 'T', 'G']
```

Use the outer most index to access each sequence

Retrieve the nucleotide at a particular position in a sequence.

```
1 >>> nt = aln[2][3]
2 >>> nt
3 'T'
```

Use the outer most index to access the sequence of interest and the inner most index to access the position

Get every nucleotide in a single column: an entire column by retrieving a particular position for each inner

```

1 >>> col = [seq[3] for seq in aln]
2 >>> col
3 ['T', 'A', 'T', 'T']

```

Here we can access a set position (3) for each list (here called 'seq') in our 'aln' list.

Lists of dictionaries

You can nest dictionaries in lists as well:

```

1 >>> records = [
2 ... {'name' : 'actgctagt', 'accession' : 'ABC123', 'genetic_code' :
3 ...   1},
4 ... {'name' : 'ttaggttta', 'accession' : 'XYZ456', 'genetic_code' :
5 ...   1},
6 ... {'name' : 'cgcgatcgt', 'accession' : 'HIJ789', 'genetic_code' :
7 ...   5}
8 ... ]
9 >>> records[0]['name']
10 'actgctagt'
11 >>> records[0]['accession']
12 'ABC123'
13 >>> records[0]['genetic_code']
14 5

```

Here you can retrieve the accession of one record at a time by using a combination of the outer index and the key 'accession'

Dictionaries of lists

And, if you haven't guessed, you can nest lists in dictionaries

Here is a dictionary of kmers. The key is the kmer and its values is a list of positions

```

1 >>> kmers = {'ggaa': [4, 10], 'aatt': [0, 6, 12], 'gaat': [5, 11],
2 ...         'tgga': [3, 9], 'attg': [1, 7, 13], 'ttgg': [2, 8]}
3 >>> kmers
4 {'tgga': [3, 9], 'ttgg': [2, 8], 'aatt': [0, 6, 12], 'attg': [1, 7,
5 ...         13], 'ggaa': [4, 10], 'gaat': [5, 11]}
6 >>>
7 >>> kmers['ggaa']
8 [4, 10]
9 >>> len(kmers['ggaa'])
10 2

```

Here we can get a list of the positions of a kmer by using the kmer as the key. We can also do things to the returned list, like determining its length. The length will be the total count of this kmers.

You can also use the dictionary `get()` method to retrieve records.

```

1 >>> kmers['ggaa']
2 [4, 10]
3 >>> kmers.get('ggaa')
4 [4, 10]

```

These two statements returns the same results

Dictionaries of dictionaries

Dictionaries of dictionaries is my favorite!! You can do so many useful things with this datastructure.

```

1 >>> genes = {
2 ... 'gene1' : {
3 ...     'seq' : "TATGCC",
4 ...     'desc' : 'something',
5 ...     'len' : 6,
6 ...     'nt_comp' : {
7 ...         'A' : 1,
8 ...         'T' : 2,
9 ...         'G' : 1,
10 ...        'C' : 2,
11 ...     }
12 ... },
13 ...
14 ... 'gene2' : {
15 ...     'seq' : "CAAATG",
16 ...     'desc' : 'something',
17 ...     'len' : 6,
18 ...     'nt_comp' : {
19 ...         'A' : 3,
20 ...         'T' : 1,
21 ...         'G' : 1,
22 ...         'C' : 1,
23 ...     }
24 ... }
25 ... }
26 >>> genes
27 {'gene1': {'nt_comp': {'C': 2, 'G': 1, 'A': 1, 'T': 2}, 'desc':
28     'something', 'len': 6, 'seq': 'TATGCC'}, 'gene2': {'nt_comp':
29     {'C': 1, 'G': 1, 'A': 3, 'T': 1}, 'desc': 'something', 'len': 6,
30     'seq': 'CAAATG'}}
31 >>> genes['gene2']['nt_comp']
32 {'C': 1, 'G': 1, 'A': 3, 'T': 1}

```

Here we store a gene name as the outermost key, with a second level of keys for qualities of the gene, like sequence, length, nucleotide composition. We can retrieve a quality by using the gene name and quality in conjunction.

There are also specific data table and frame handling libraries like Pandas.

[Link to Python 8 Problem Set](#)

Python 9

BioPython

BioPython Overview

BioPython Subtopic 1

BioPython Subtopic 2

Bioinformatic Analysis and Tools

Bioinformatic Analysis and Tools Overview

- What you want to do:
 - tools to do it

DNA

RNA