

```
1. """
2. Toward higher-order functions
3.
4. Dan Armendariz
5. danallan@cs.harvard.edu
6. """
7.
8. # create an arbitrary list of numbers
9. nums = [1, 5, 10, 8, -4, 432]
10.
11. print "Initial:", nums
12.
13. # now let's add 1 to each number
14. result = []
15. for num in nums:
16.     result.append(num + 1)
17.
18. print "Add 1: ", result
19.
20. # perhaps we now want to multiply each by some value
21. result = []
22. for num in nums:
23.     result.append(num * 2)
24.
25. print "Double: ", result
```

```
1. """
2. Toward higher-order functions
3.
4. Dan Armendariz
5. danallan@cs.harvard.edu
6. """
7.
8. # create an arbitrary list of numbers
9. nums = [1, 5, 10, 8, -4, 432]
10.
11. print "Initial:", nums
12.
13. def incrementor(num):
14.     """ Increment a number by 1 """
15.     return num + 1
16.
17. # iterate over every item in the list, applying our function to it
18. result = []
19. for num in nums:
20.     result.append(incrementor(num))
21.
22. print "Add 1: ", result
23.
24.
25. def doubler(num):
26.     """ Double a number """
27.     return num * 2
28.
29. # iterate over every item in the list, applying our function to it
30. result = []
31. for num in nums:
32.     result.append(doubler(num))
33.
34. print "Double: ", result
```

```
1. """
2. Toward higher-order functions
3.
4. Dan Armendariz
5. danallan@cs.harvard.edu
6. """
7.
8. # create an arbitrary list of numbers
9. nums = [1, 5, 10, 8, -4, 432]
10.
11. print "Initial:", nums
12.
13. def apply_to_all(f, items):
14.     """ Apply a function f to every item in a list """
15.     result = []
16.     for item in items:
17.         new_item = f(item)
18.         result.append(new_item)
19.     return result
20.
21. def incrementor(num):
22.     """ Increment a number by 1 """
23.     return num + 1
24.
25. def doubler(num):
26.     """ Double a number """
27.     return num * 2
28.
29. print "Add 1: ", apply_to_all(incrementor, nums)
30. print "Double: ", apply_to_all(doubler, nums)
```

```
1. """
2. Using Python's built-in Higher-Order Functions (HOFs): map
3.
4. Dan Armendariz
5. danallan@cs.harvard.edu
6. """
7.
8. # create an arbitrary list of numbers
9. nums = [1, 5, 10, 8, -4, 432]
10.
11. print "Initial:", nums
12.
13. def incrementor(num):
14.     """ Increment a number by 1 """
15.     return num + 1
16.
17. def doubler(num):
18.     """ Double a number """
19.     return num * 2
20.
21. # the `map` function is a HOF that applies a function to an entire list
22. print "Add 1: ", map(incrementor, nums)
23. print "Double: ", map(doubler, nums)
```

```
1. """
2. Using Python's built-in Higher-Order Functions (HOFs): reduce
3.
4. Dan Armendariz
5. danallan@cs.harvard.edu
6. """
7.
8. # create an arbitrary list of numbers
9. nums = [1, 5, 10, 8, -4, 432]
10.
11. print "Initial:", nums
12.
13. def add(x, y):
14.     """ Sum two numbers """
15.     return x + y
16.
17. def multiply(x, y):
18.     return x * y
19.
20. # the `reduce` function is a HOF that collapses a list to a single value
21. print "Sum of all values: ", reduce(add, nums, 0)
22. print "Multiply all values: ", reduce(multiply, nums, 1)
23.
24. """
25. Reduce operates in this way:
26. We want to sum all of the objects, so we provide our function `sum` to reduce()
27. It then applies sum on the first element and the initial value, and the
28. output of that function is then used as input to the next iteration.
29.
30. Let's say we have a list [1, 5, 8]
31. The first operation would be:
32. sum(0, 1) -- "0" is our initial value. The output is 1
33.
34. 1 is used as one of the next inputs; the other is the next item in the list
35. sum(1, 5) -- outputs "6"
36.
37. sum(6, 8) -- outputs our final value, 14!
38.
39. Note: the order is not guaranteed!
40. """
```

```
1. """
2. Use MapReduce (Hadoop) style programming to count the most common letter
3. unigrams (single characters) and bigrams (character pairs) in a text.
4. See: http://en.wikipedia.org/wiki/Bigram
5.
6. Hadoop assumes that the input data is a text file. This is convenient,
7. because text files can be easily broken into chunks at individual lines
8. and distributed to many computers. It then roughly operates in this fashion:
9.
10. MAP STEP: runs a developer-provided function on one line at a time of the input
11. files. The Map step is therefore generally used to convert the data represented
12. in the line into one or more key/value pairs. Normally, we could think of
13. "key/value" pairs as a dictionary, but Hadoop lets us have duplicate keys.
14. So, in essence, if we're writing code for Hadoop (which is normally done in the
15. Java programming language, but we're using Python for familiarity), we want
16. to write a map function that will accept a line of text and output a dictionary.
17.
18. SHUFFLE STEP: As key/value pairs come in from the map step, hadoop sorts them
19. and redirects all values that have the same key to the same computer to run
20. the Reduce step. The map step must finish completely before moving on to
21. the Reduce step.
22.
23. REDUCE STEP: Another developer-provided function, the reducer is given a
24. key and a list of all of the values for that key by Hadoop. It then
25. performs some computation and outputs a *single* value for that key.
26.
27. SORT STEP: Hadoop will generally sort the final output data.
28.
29. It is not unusual to have to run several Map/Reduce steps on a big data set
30. before receiving the output you want. In our example below, we'll only need a
31. single Map and single Reduce step to complete our goal.
32.
33. See more information on MapReduce applications, of which Hadoop is a part, here:
34. http://en.wikipedia.org/wiki/MapReduce
35.
36. For a good exercise, consider how you could modify the code below to count
37. unigram WORDS and bigram WORDS. In other words, instead of counting individual
38. characters or character pairs, counting whole words and counting word pairs.
39.
40. Dan Armendariz, Dan Bradley
41. danallan@cs.harvard.edu, dbradley@cs50.harvard.edu
42. """
43.
44. # itemgetter will assist in sorting
45. from operator import itemgetter
46.
47. def unigram_mapper(line):
48.     """Given a line of text, remove extraneous characters and output a key/value
```

```
49.     pair for each character. The character will be the key and the value will
50.     be 1.
51.
52.     For example, given the string "Hello, world" we will output:
53.     [("h",1), ("e",1), ("l",1), ("l",1), ("o",1), ("w",1), ("o",1), ("r",1),
54.      ("l",1), ("d",1)]
55.     """
56.     output = []
57.
58.     # perform the step for every character in the string
59.     for char in line:
60.         # skip the character if it is not a letter
61.         if not char.isalpha():
62.             continue
63.
64.         # create our keyvalue pair, converting the letter to lower case
65.         keyvalue_pair = (char.lower(), 1)
66.         output.append(keyvalue_pair)
67.
68.     return output
69.
70.
71. def bigram_mapper(line):
72.     """Given a line of text, remove extraneous characters and output a key/value
73.     pair for each bigram. The bigram will be the key and the value will be 1.
74.
75.     For example, given the string "Hello, world" we will output:
76.     [("he",1), ("el",1), ("ll",1), ("lo",1), ("wo",1), ("or",1), ("rl",1),
77.      ("ld", 1)]
78.     """
79.     output = []
80.
81.     # iterate over each character in the line
82.     for i in range(len(line)-1):
83.         # fetch a bigram by splicing the line at i to grab 2 neighboring chars
84.         bigram = line[i:i+2]
85.
86.         # skip the bigram if either character isn't a letter
87.         if not bigram.isalpha():
88.             continue
89.
90.         # create the key/value pair, and ensure the bigram is lower case
91.         keyvalue_pair = (bigram.lower(), 1)
92.         output.append(keyvalue_pair)
93.
94.     return output
95.
96.
```

```
97. def count_reducer(key, values):
98.     """Given a list of values for a key, we'll return the sum of all values."""
99.     total = 0
100.    for value in values:
101.        total += value
102.    return (key, total)
103.
104.
105. def hadoop(file, mapper, reducer):
106.     """A higher-order function that accepts a mapper and a reducer and
107.     performs Hadoop/MapReduce-style processing on a given file.
108.
109.     Essentially, this code (in a very simple way) performs on a single computer
110.     what Hadoop and other MapReduce implementations do across a wide number
111.     of computers.
112.     """
113.     # a structure to store the final output
114.     data = []
115.
116.     # MAP STEP
117.     pairs = []
118.     with open(file) as f:
119.         # iterate over every line in the opened file
120.         for line in f:
121.             # ask our map step to process the line
122.             output = mapper(line)
123.
124.             # store the newly computed key/value pairs in our `pairs` structure
125.             pairs.extend(output)
126.
127.     # SHUFFLE STEP: aggregate all values for a given key
128.     collapsed = {}
129.     for pair in pairs:
130.         # split the pair into its key and value
131.         key, value = pair
132.
133.         # if the key doesn't exist yet in our dictionary, add it
134.         if collapsed.get(key, None) is None:
135.             collapsed[key] = []
136.
137.         # append the value to the list of values for a given key
138.         collapsed[key].append(value)
139.
140.
141.     # REDUCE STEP: provide all of the values for a key to the reducer
142.     for key, values in collapsed.iteritems():
143.         # the key is different; run the reduce step on our accumulated vals
144.         output = reducer(key, values)
```



```
145.         data.append(output)
146.
147.         # SORT STEP: sort the output by values in descending order
148.         data.sort(key=itemgetter(1), reverse=True)
149.
150.         # all done!
151.         return data
152.
153. # Using "Hadoop" and our reducer and unigram mapper, count the number of
154. # instances each character appears in Carroll's "Alice and Wonderland"
155. print "Unigram character counts:"
156. result = hadoop("aliceandwonderland.txt", unigram_mapper, count_reducer)
157. for count in result:
158.     print count
159.
160. print
161.
162. # On to bigrams. Our reducer is so simple we can use the same one for both!
163. print "Bigram character counts:"
164. result = hadoop("aliceandwonderland.txt", bigram_mapper, count_reducer)
165.
166. # there are many bigrams, so we'll only print the first 20
167. for count in result[0:10]:
168.     print count
```