

# Data Structures and Algorithms with Python

## Workshop 2

In this workshop you'll implement search and sort algorithms in Python. You'll also practice reading Python code written by others, and work on the concept of Big-Oh complexity. A skeleton code for these exercises is provided in the Python file available on the Hub. Open this file in Spyder and you're ready to go. This is the core part of the workshop.

The main workshop is followed by two extra problems — this time, you can explore a dataset of Fortune 500 companies, as well as the merge sort algorithm, which greatly improves sorting speed compared to selection sort. This is intended for those of you who wish to challenge themselves and go further than the required course material. The exercises may be open to different approaches and sometimes involve new concepts — Google may be your friend here.

### 1 Searching

In the lecture we started to delve into the mechanics of searching for an item in a list. We first looked at how linear search goes through each element  $x$  in list  $A$ , as defined below.

```
1 def linearSearch(A,x):
2     for elem in A:
3         if elem == x:
4             return True
5     return False
6
7 y = 6
8 linearSearch(A,y)
```

Try it out and make sure you understand how the function works. Notice that the elements of the list do not need to be integers (or indeed numbers).

We saw in the lecture that the worst-case running time of this algorithm is  $O(n)$ , where  $n$  is the length of  $A$ : that is, if the element we're looking for happen not to be in the list, we have to go

through all  $n$  elements.

If your list is sorted, a better solution is to use **binary search**. We defined the algorithm in the lecture using the following “recipe” to search for item  $x$  in list  $L$ .

- Pick an index  $i$  roughly dividing  $L$  in half
- If there is nothing left to search, return False
- If  $L[i] == x$ , return True
- If not:
  - If  $L[i] > x$ , recursively search left half of  $L$
  - Otherwise recursively search right half

Complete the skeleton code given in the Python file. Your function should recursively call itself, narrowing down the section of the list that it is searching. It may be helpful to use `print`-statements in your code to see what is happening in the recursion.

```
1 def binarySearch(A,x,low,high):
2     """Assumes A is a list with elements in ascending order.
3     x is the element we're searching for.
4     low and high are indices of the list in between which we search.
5     Returns True if x is in A and False otherwise"""
6     if low >= high: # if range is empty: nothing left to search
7         if A[low] == x: # if just found x, return True
8             return True
9         else: # otherwise did not find x
10             return False
11     midpoint = (low + high) // 2 # divide list roughly in half using integer division
12     if A[midpoint] == x:
13         return True # found x!
14     elif A[midpoint] > x:
15         return binarySearch(A,x,low,midpoint-1) # search left side of midpoint recursively
```

```

16     else:
17         return binarySearch(A,x,?,?) # code: search right side of midpoint recursively

```

Verify that your search works and compare the speed of your implementation with linear search. You can use the `%timeit` command to time the performance of your functions. Note that this command works in the “enhanced” IPython console that Spyder uses, but not in a standard Python console.

```

1 %timeit linearSearch(A,y)

```

## 2 Sorting

We saw earlier how searching becomes much faster when we can exploit the fact that a list is sorted. Today we’ve looked at how sorting algorithms work under the hood. You’ll implement the two methods you saw in the lecture, selection sort and merge sort.

The selection sort algorithm uses searching as a subroutine. In each iteration of the main loop, it searches for the smallest element in the unsearched part of the list and moves it to the end of the searched part. The algorithm is as follows:

**Selection sort** of list  $L$  of length  $n$ :

- Initialization step: divide the list into a “prefix”  $P$  and a “suffix”  $S$  with initially  $P$  empty and  $S = L$
- Main loop:
  - Search for the smallest element of  $S$  and move it to the end of  $P$
  - Repeat until  $S$  is empty

Complete the skeleton code in the Python file. The idea in this implementation is to keep track of the prefix and a suffix using a single index, `prefixSize`. Each go of the main loop goes adds another element to the prefix, increasing `prefixSize` by one. Your task is to make sure that the algorithm finds the minimum element of the suffix and moves it to this index.

```

1 def selectionSort(A):
2     """Assumes A is a list (whose elements can be compared).
3     Returns sorted list"""
4     # see code in the Python file provided

```

Although the selection sort algorithm is intuitive, it is not as fast as we would hope for: its running time is  $O(n^2)$  for  $n$  items. We can do better by other algorithms, for example using merge sort, which is a recursive “divide-and-conquer” algorithm. It significantly improves the complexity of sorting to  $O(n \log n)$ . You can explore this algorithm in a challenge problem below.

### 3 Big-Oh practice

**Exercise.** How would you order the following functions in increasing order of complexity?

$$T_1(n) = 20n \log n$$

$$T_2(n) = (n - 1)^2$$

$$T_3(n) = 2^n$$

$$T_4(n) = \log n^4$$

The following exercises are intended for you to both practice reading Python code, and developing a sense for the big-Oh complexity of different loops. It may be helpful to run the code in Spyder for different inputs.

What do the following code snippets do? What is their big-Oh complexity in terms of the size of the input (eg a list length)?

```

1 def lsum(L):
2     """ assumes L is list of length n """
3     listSum = 0
4     for item in L:

```

```
5         listSum = listSum + item
6     print(listSum)
```

```
1 def addOne(L):
2     """ assumes L is list of length n """
3     for index in range(len(L)):
4         L[index] = L[index] + 1
5     print(L)
```

```
1 def mult(L1,L2):
2     """ assumes L1,L2 are lists of length n,m """
3     s = 0
4     for i in L1:
5         for j in L2:
6             s = s + i*j
7     print(s)
```

```
1 def testSubset(x,y):
2     """ assumes x and y are lists of length n,m"""
3     for a in x:
4         found = False
5         for b in y:
6             if a == b:
7                 found = True
8                 break
9     if not found:
10        return False
```

```
11     return True
```

The following exercises are more challenging...

```
1 def intToStr(i):
2     """ assumes i is integer """
3     digits = '0123456789'
4     if i==0:
5         return '0'
6     result = ''
7     while i>0:
8         result = digits[i%10] + result
9         i = i//10
10    return result
```

```
1 def allIntsToStr(A):
2     """ assumes A is a list of integers """
3     r = []
4     for i in A:
5         r.append(intToStr(i))
6     return r
```

```
1 def F(n):
2     """ assume n is a nonnegative integer """
3     y = 0
4     i = n
5     while i > 0:
6         j = 0
7         while j < i:
```

```

8         y += 1
9         j += 1
10        i -= 1
11    return y

```

## 4 Built-in Python functions for searching and sorting

Searching and sorting are common building blocks in more complex programs and their performance is therefore crucial. Whilst designing and writing these algorithms is both fun and useful for developing your algorithmic and coding skills, in practice we usually rely on Python's built-in methods.

Python gives you an easy method of checking whether an item can be found in a list through the keyword `in`.

```

1  A = [3,8,11,9,4]
2  4 in A

```

You can sort a list using `list.sort()` and create a sorted copy using `sorted`. If you want to change the order, you can simply reverse the list.

```

1  A.sort()
2  B = sorted(A)
3  B = sorted(A, reverse=True)
4  B.reverse()
5  help(list)

```

But how does Python do its sorting? It uses an algorithm called **timsort**, which has worst-case running time  $O(n \log n)$  just like your merge sort, but is faster on average. In fact, it turns out that  $O(n \log n)$  is a theoretical bound for how fast you can expect to sort in the worst case. Timsort's is faster on average because it exploits the fact that many lists have some order to begin with, and finding such ordered sequences saves time compared to straight-up sorting. However,

the cost is that timsort is conceptually quite a bit more complex than merge sort (just like merge sort is conceptually more complex than selection sort). If you want to understand a state-of-the-art sorting algorithm, look up timsort on Wikipedia. While you're at it, there are many other well-known sorting algorithms that we don't go through here: for example insertion sort, bubble sort, quicksort, radix sort, just to name a few.

## 5 Challenge problem: Fortune 500

The file `fortune500.csv` contains a list of Fortune 500 companies every year from 1955 to 2005 and Fortune 1000 companies from 2006 to 2009.<sup>1</sup> The Fortune 500 list comprises the largest American companies by yearly revenue.

The file is comma-delimited, with the first line containing column titles (Year, Rank, Company, Revenue (in millions), Profit (in millions)). The subsequent lines contain this information for the companies in each year's list.

Use Python to read through the file and explore the questions below. The Python file provides two ways by which you could proceed: by methods you already know, or a different but convenient Python library. That is, you can either read the file into Python lists and work from there (using eg loops), or alternatively use the Python library `pandas`. Code for reading the file into Python is provided for both options.

`pandas` is a convenient library for dealing with data in Python. We will cover it in more detail in later tutorials. For now, the code in the Python file explains some useful commands in dealing with data — you can find `pandas` help at <http://pandas.pydata.org/pandas-docs/version/0.18.1/tutorials.html>.

When reading data in Spyder, make sure your data is in your current working directory (visible in the Spyder toolbar).

Questions to explore:

1. What were the largest profits and revenues of a company, both overall and in each decade covered by the data? What were the largest losses?
2. Which companies have held the number one slot? Which companies have been on the list for the longest periods? Are there companies that have stayed on the list during the entire time?

---

<sup>1</sup>Cleaned data courtesy of Robert Sedgewick, Princeton University.



3. How has the average profit of the companies on the list developed over the years?
4. Visualize the position of different companies on the list in a plot. For example, how has the position of Apple and Kodak changed over the years? We will cover a Python plotting library, `matplotlib`, in a later tutorial. For now, you can find a help on plotting here: [http://matplotlib.org/users/pyplot\\_tutorial.html](http://matplotlib.org/users/pyplot_tutorial.html).

## 6 Challenge problem: Merge sort

Although the selection sort algorithm is intuitive, it is not as fast as we would hope for: its running time is  $O(n^2)$  for  $n$  items. We can do better by other algorithms, for example using merge sort, which is a recursive “divide-and-conquer” algorithm. It significantly improves the complexity of sorting to  $O(n \log n)$ .

The merge sort algorithm is given below. The idea bears some resemblance to binary search. You start with a list you want to sort. You identify a *base case* for which sorting is easy: if the length of a list is no greater than one, the list is obviously sorted. So you recursively divide the list roughly in half as many times as it requires to reach this base case - approximately  $\log n$  times. You then combine these results together at each level of recursion and return the result. This combining is called *merging* lists, and the exact procedure is given in the algorithm below. As it turns out, each level of recursion requires  $O(n)$  work. So merge sort in total has  $O(n \log n)$  complexity.

### Merge sort:

- Base case: if list length  $n \leq 1$ , the list is sorted: return the list itself
- Divide: if list length  $n > 1$ , split into two lists and recursively sort each
- Combine (merge) the two (sorted) lists:
  - Initialize an empty **result** list
  - Look at first element of each list, add smaller to end of **result**
  - When one of the lists becomes empty, copy the rest of the other list to **result**
- Return the result of merging

The Python file contains a full implementation of the `merge` subroutine, as well as a skeleton code for the `mergesort` routine which you'll have to complete. For an example of how merge sort works, see the extra slides in the lecture folder, or this wonderful video: [https://www.youtube.com/watch?v=XaqR3G\\_NVoo](https://www.youtube.com/watch?v=XaqR3G_NVoo).