

## Complexity of Python Operations

In this lecture we will learn the complexity classes of various operations on Python data types. Then we will learn how to combine these complexity classes to compute the complexity class of all the code in a function, and therefore the complexity class of the function. This is called "static" analysis, because we do not need to run any code to perform it (contrasted with Dynamic or Empirical Analysis, when we do run code and take measurements of its execution).

## Python Complexity Classes

In ICS-46 we will write low-level implementations of all of Python's data types and see/understand WHY these complexity classes apply. For now we just need to try to absorb (not memorize) this information, with some -but minimal- justification.

Binding a value to any name (copying a reference) is  $O(1)$ . Simple operators on integers (whose values are small: e.g., under 12 digits) like  $+$  or  $==$  are also  $O(1)$ . You should assume small integers in problems unless explicitly told otherwise.

In all these examples,  $N = \text{len}(\text{data-type})$ . The operations are organized by increasing complexity class

## Lists:

Operation	Example	Complexity Class	Notes
Index	<code>l[i]</code>	$O(1)$	
Store	<code>l[i] = 0</code>	$O(1)$	
Length	<code>len(l)</code>	$O(1)$	
Append	<code>l.append(5)</code>	$O(1)$	mostly: ICS-46 covers details
Pop	<code>l.pop()</code>	$O(1)$	same as <code>l.pop(-1)</code> , popping at end
Clear	<code>l.clear()</code>	$O(1)$	similar to <code>l = []</code>
Slice	<code>l[a:b]</code>	$O(b-a)$	<code>l[1:5]:O(1)/l[:]:O(len(l)-0)=O(N)</code>
Extend	<code>l.extend(...)</code>	$O(\text{len}(...))$	depends only on len of extension
Construction	<code>list(...)</code>	$O(\text{len}(...))$	depends on length of ... iterable
check <code>==, !=</code>	<code>l1 == l2</code>	$O(N)$	
Insert	<code>l[a:b] = ...</code>	$O(N)$	
Delete	<code>del l[i]</code>	$O(N)$	depends on $i$ ; $O(N)$ in worst case
Containment	<code>x in/not in l</code>	$O(N)$	linearly searches list
Copy	<code>l.copy()</code>	$O(N)$	Same as <code>l[:]</code> which is $O(N)$
Remove	<code>l.remove(...)</code>	$O(N)$	
Pop	<code>l.pop(i)</code>	$O(N)$	$O(N-i)$ : <code>l.pop(0):O(N)</code> (see above)
Extreme value	<code>min(l)/max(l)</code>	$O(N)$	linearly searches list for value
Reverse	<code>l.reverse()</code>	$O(N)$	
Iteration	<code>for v in l:</code>	$O(N)$	Worst: no return/break in loop
Sort	<code>l.sort()</code>	$O(N \log N)$	key/reverse mostly doesn't change
Multiply	<code>k*l</code>	$O(kN)$	<code>5*l</code> is $O(N)$ : <code>len(l)*l</code> is $O(N^2)$

Tuples support all operations that do not mutate the data structure (and they have the same complexity classes).

## Sets:

Operation	Example	Complexity Class	Notes
-----------	---------	---------------------	-------

Length	<code>len(s)</code>	$O(1)$	
Add	<code>s.add(5)</code>	$O(1)$	
Containment	<code>x in/not in s</code>	$O(1)$	compare to list/tuple - $O(N)$
Remove	<code>s.remove(..)</code>	$O(1)$	compare to list/tuple - $O(N)$
Discard	<code>s.discard(..)</code>	$O(1)$	
Pop	<code>s.pop()</code>	$O(1)$	popped value "randomly" selected
Clear	<code>s.clear()</code>	$O(1)$	similar to <code>s = set()</code>
Construction	<code>set(...)</code>	$O(\text{len}(...))$	depends on length of ... iterable
check <code>==, !=</code>	<code>s != t</code>	$O(\text{len}(s))$	same as <code>len(t)</code> ; False in $O(1)$ if the lengths are different
<code>&lt;= / &lt;</code>	<code>s &lt;= t</code>	$O(\text{len}(s))$	issubset
<code>&gt;= / &gt;</code>	<code>s &gt;= t</code>	$O(\text{len}(t))$	issuperset <code>s &lt;= t == t &gt;= s</code>
Union	<code>s   t</code>	$O(\text{len}(s) + \text{len}(t))$	
Intersection	<code>s &amp; t</code>	$O(\text{len}(s) + \text{len}(t))$	
Difference	<code>s - t</code>	$O(\text{len}(s) + \text{len}(t))$	
Symmetric Diff	<code>s ^ t</code>	$O(\text{len}(s) + \text{len}(t))$	
Iteration	<code>for v in s:</code>	$O(N)$	Worst: no return/break in loop
Copy	<code>s.copy()</code>	$O(N)$	

Sets have many more operations that are  $O(1)$  compared with lists and tuples. Not needing to keep values in a specific order in a set (which lists/tuples require an order) allows for faster implementations of set operations.

Frozen sets support all operations that do not mutate the data structure (and they have the same complexity classes).

Dictionaries: dict and defaultdict

		Complexity	
Operation	Example	Class	Notes
Index	<code>d[k]</code>	$O(1)$	
Store	<code>d[k] = v</code>	$O(1)$	
Length	<code>len(d)</code>	$O(1)$	
Delete	<code>del d[k]</code>	$O(1)$	
get/setdefault	<code>d.get(k)</code>	$O(1)$	
Pop	<code>d.pop(k)</code>	$O(1)$	popped key "randomly" selected
Pop item	<code>d.popitem()</code>	$O(1)$	popped item "randomly" selected
Clear	<code>d.clear()</code>	$O(1)$	similar to <code>s = {}</code> or <code>= dict()</code>
View	<code>d.keys()</code>	$O(1)$	same for <code>d.values()</code>
Construction	<code>dict(...)</code>	$O(\text{len}(...))$	depends # (key,value) 2-tuples
Iteration	<code>for k in d:</code>	$O(N)$	all forms: keys, values, items Worst: no return/break in loop

So, most dict operations are  $O(1)$ .

defaultdicts support all operations that dicts support, with the same complexity classes (because it inherits all those operations); this assumes that calling the constructor when a values isn't found in the defaultdict is  $O(1)$  - which is true for `int()`, `list()`, `set()`, ... (the things we commonly use)

Note that `for i in range(...)` is  $O(\text{len}(...))$ ; so `for i in range(1,10)` is  $O(1)$ . If `len(alist)` is  $N$ , then

```
for i in range(len(alist)):
```

is  $O(N)$  because it loops  $N$  times. Of course even

```
for i in range (len(alist)//2):
```

is  $O(N)$  because it loops  $N/2$  times, and dropping the constant  $1/2$  makes it  $O(N)$ : the work doubles when the list length doubles. By this reasoning,

```
for i in range (len(alist)//1000000):
```

is  $O(N)$  because it loops  $N/1000000$  times, and dropping the constant 1000000 makes it  $O(N)$ : the work doubles when the list length doubles.

Finally, when comparing two lists for equality, the complexity class above shows as  $O(N)$ , but in reality we would need to multiply this complexity class by  $O==(\dots)$  where  $O==(\dots)$  is the complexity class for checking whether two values in the list are  $==$ . If they are ints,  $O==(\dots)$  would be  $O(1)$ ; if they are strings,  $O==(\dots)$  in the worst case it would be  $O(\text{len}(\text{string}))$ . This issue applies any time an  $==$  check is done. We mostly will assume  $==$  checking on values in lists is  $O(1)$ : e.g., checking ints and small/fixed-length strings.

## Composing Complexity Classes: Sequential and Nested Statements

In this section we will learn how to combine complexity class information about simple operations into complexity class information about complex operations (composed from simple operations). The goal is to be able to analyze all the statements in a function/method to determine the complexity class of executing the function/method. As with computing complexity classes themselves, these rules are simple and easy to apply once you understand how to use them.

### Law of Addition for big-O notation

$O(f(n)) + O(g(n))$  is  $O(f(n) + g(n))$

That is, when adding complexity classes we bring the two complexity classes inside the  $O(\dots)$ . Ultimately,  $O(f(n) + g(n))$  results in the bigger of the two complexity classes (because we always drop the lower-complexity added term). So,

$O(N) + O(\log N) = O(N + \log N) = O(N)$

because  $N$  is the faster growing term:  $\lim_{N \rightarrow \infty} \log N / N = 0$ .

This rule helps us understand how to compute the complexity class of doing any SEQUENCE of operations: executing a statement that is  $O(f(n))$  followed by executing a statement that is  $O(g(n))$ . Executing both statements SEQUENTIALLY is  $O(f(n)) + O(g(n))$  which is  $O(f(n) + g(n))$  by the rule above.

For example, if some function call  $f(\dots)$  is  $O(N)$  and another function call  $g(\dots)$  is  $O(N \log N)$ , then doing the sequence

```
f(...)
g(...)
```

is  $O(N) + O(N \log N) = O(N + N \log N) = O(N \log N)$ . Of course, executing the sequence (calling  $f$  twice)

```
f(...)
f(...)
```

is  $O(N) + O(N)$  which is  $O(N + N)$  which is  $O(2N)$  which is  $O(N)$ .

Note that an if statement sequentially evaluates test and then one of the blocks.

```
if test:      assume complexity class of computing test is O(T)
    block 1   assume complexity class of executing block 1 is O(B1)
else:
    block 2   assume complexity class of executing block 2 is O(B2)
```

The complexity class for the if is  $O(T) + \max(O(B1), O(B2))$ . The test is always evaluated, and one of the blocks is always executed afterward (so, a sequence of evaluating a test followed by executing a block). In the worst case, the if will execute the block with the largest complexity class. So, given

```
if test:      complexity class is O(N)
    block 1   complexity class is O(N**2)
else:
    block 2   complexity class is O(N)
```

The complexity class for the if is  $O(N) + \max(O(N**2), O(N)) = O(N) + O(N**2) = O(N + N**2) = O(N**2)$ .

If the test had complexity class  $O(N**3)$ , then the complexity class for the if is  $O(N**3) + \max(O(N**2), O(N)) = O(N**3) + O(N**2) = O(N**3 + N**2) = O(N**3)$ .

-----

Law of Multiplication for big-O notation

$O(f(n)) * O(g(n))$  is  $O(f(n) * g(n))$

If we repeat an  $O(f(N))$  process  $O(N)$  times, the resulting complexity class is  $O(N) * O(f(N)) = O(N * f(N))$ . An example of this is, if some function call  $f(...)$  is  $O(N**2)$ , then executing that call  $N$  times (in the following loop)

```
for i in range(N):
    f(...)
```

is  $O(N) * O(N**2) = O(N * N**2) = O(N**3)$

This rule helps us understand how to compute the complexity class of doing some statement INSIDE A BLOCK controlled by a statement that is REPEATING it. We multiply the complexity class of the number of repetitions by the complexity class of the statement (sequence; using the summing rule) being repeated.

Compound statements can be analyzed by composing the complexity classes of their constituent statements. For sequential statements (including if tests and their block bodies) the complexity classes are added; for statements repeated in a loop the complexity classes are multiplied.

Let's use the data and tools discussed above to analyze (determine their complexity classes) three different functions that each compute the same result: whether or not a list contains only unique values (no duplicates). We will assume in all three examples that  $\text{len}(\text{alist})$  is  $N$  and that we can compare the list elements in  $O(1)$ : e.g., they are small ints or strs.

1) Algorithm 1: A list is unique if each value in the list does not occur in any later indexes:  $\text{alist}[i+1:]$  is a list slice containing all values after the one at index  $i$ .

```
def is_unique1 (alist : [int]) -> bool:
    for i in range(len(alist)):      O(N) - for every index
        if alist[i] in alist[i+1:]:  O(N) - index+add+slice+in: O(1)+O(1)+O(N)+O(N) = O(N)
            return False             O(1) - never executed in worst case; ignore
    return True                      O(1)
```

The complexity class for executing the entire function is  $O(N) * O(N) + O(1) = O(N**2)$ . So we know from the previous lecture that if we double the length of  $\text{alist}$ , this function takes 4 times as long to execute.

-----

Many students want to write this as  $O(N) * (O(N) + O(1)) + O(1)$  because the if statement's complexity is  $O(N) + O(1)$ : complexity of test + complexity of block when test is True. But in the worst case, the return is NEVER EXECUTED

(the loop keeps executing) so it should not appear in the formula. Although, even if it appears in this formula, the formula still computes the same complexity class (because  $O(N) + O(1)$  is still  $O(N)$ ):  $O(N^2)$ .

So, in the worst case, we never return False and keep executing the loop, so this  $O(1)$  does not appear in the formula. Also, in the worst case the list slice is `alist[1:]` which is  $O(N-1) = O(N)$ , although when `i` is `len(alist)` the slice contains 0 values: is empty. The average list slice taken in the if has  $N/2$  values, which is still  $O(N)$ .

-----

2) Algorithm 2: A list is unique if when we sort its values, no ADJACENT values are equal. If there were duplicate values, sorting the list would put these duplicate values right next to each other (adjacent). Here we copy the list so as to not mutate (change the order of) the parameter's list by sorting it (functions generally shouldn't mutate their arguments unless that is the purpose of the function): it turns out that copying the list does not increase the complexity class of the method, because the  $O(N)$  used for copying is not the largest added term.

```
def is_unique2 (alist : [int]) -> bool:
    copy = list(alist)           O(N)
    copy.sort()                  O(N Log N) - for fast Python sorting
    for i in range(len(alist)-1): O(N) - really N-1, but that is O(N); len and - are
both O(1)
        if copy[i] == copy[i+1]: O(1): +, 2 [i], and == on ints: all O(1)
            return False         O(1) - never executed in worst case
    return True                  O(1)
```

The complexity class for executing the entire function is given by the sum  $O(N) + O(N \log N) + O(N) \cdot O(1) + O(1) = O(N + N \log N + O(N \cdot 1) + 1) = O(N + N \log N + N + 1) = O(N \log N + 2N + 1) = O(N \log N)$ . So the complexity class for this algorithm/function is lower than the first algorithm, the `is_unique1` function. For large  $N$  `unique2` will eventually run faster. Because we don't know the constants, we don't know which is faster for small  $N$ .

Notice that the complexity class for sorting is dominant in this code: it does most of the work. If we double the length of `alist`, this function takes a bit more than twice the amount of time. In  $N \log N$ :  $N$  doubles and  $\log N$  gets a tiny bit bigger (i.e.,  $\log 2N = 1 + \log N$ ; e.g.,  $\log 2000 = 1 + \log 1000 = 11$ , so compared to  $1000 \log 1000$ , doubling  $N$  is  $2000 \log 2000$ , which is just 2.2 times bigger, or 10% bigger than just doubling).

Looked at another way if  $T(N) = c \cdot (N \log N)$ , then  $T(2N) = c \cdot (2N \log 2N) = c \cdot 2N (\log N + 1) = c \cdot 2N \log N + c \cdot 2N = 2 \cdot T(N) + c \cdot 2N$ . Or, computing the doubling signature

$$\frac{T(2N)}{T(N)} = \frac{c \cdot 2N \log N + c \cdot 2N}{c \cdot N \log N} = \frac{c \cdot 2N \log N}{c \cdot N \log N} + \frac{c \cdot 2N}{c \cdot N \log N} = 2 + \frac{2}{\log N}$$

So, the ratio is  $2 +$  a bit (and that bit gets smaller -very slowly- as  $N$  increases): for  $N \geq 10^3$  it is  $\leq 2.2$ ; for  $N \geq 10^6$  it is  $\leq 2.1$ ; for  $N \geq 10^9$  it is  $< 2.07$ . So, it is a bit worse than doubling each time, but much better than  $O(N^2)$  which is quadrupling each time.

We could also simplify

```
copy = list(alist)           O(N)
copy.sort()                  O(N Log N) - for fast Python sorting
```

to just

```
copy = sorted(alist)         O(N Log N) - for fast Python sorting
```

Which won't change the complexity analysis because  $O(N + N \log N) = O(N \log N)$ . Finally, Algorithm 2 works only if all the values in the list are comparable: it would fail if the list contained both integers and strings.

3) Algorithm 3: A list is unique if when we turn it into a set, its length is unchanged: if duplicate values were added to the set, its length would be smaller than the length of the list by exactly the number of duplicates in the list added to the set.

```
def is_unique3 (alist : [int]) -> bool:
    aset = set(alist)          O(N): construct set from alist values
    return len(aset) == len(alist)  O(1): 2 len (each O(1)) and == ints O(1)
```

The complexity class for executing the entire function is  $O(N) + O(1) = O(N + 1) = O(N)$ . So the complexity class for this algorithm/function is lower than both the first and second algorithms/functions. If we double the length of alist, this function takes just twice the amount of time. We could write the body of this function more simply as: `return len(set(alist)) == len(alist)`, where evaluating `set(alist)` takes  $O(N)$  and then computing the two len's and comparing them for equality are all  $O(1)$ .

Unlike Algorithm 2, it can work for lists containing both integers and strings. But, Algorithm 3 works only if all the values in the list are immutable (a requirement for storing values in a set). So, it would not work for a list of lists.

So the bottom line here is that there might be many algorithms/functions to solve some problem. If the function bodies are small, we can analyze them statically (looking at the code, not needing to run it) to determine their complexity classes. For large problem sizes, the algorithm/function with the smallest complexity class will ultimately be best, running in the least amount of time. But, for small problem sizes, complexity classes don't determine which is best: for small problem we need to take into account the CONSTANTS and lower order terms that we ignored when computing complexity classes). We can run the functions (dynamic analysis, aka empirical analysis) to test which is fastest on small problem sizes.

-----

Using a Class (implementable 3 ways) Example:

We will now look at the solution of a few problems (combining operations on a priority queue: pq) and how the complexity class of the result is affected by three different classes/implementations of priority queues.

In a priority queue, we can add values to and remove values from the data structure. A correctly working priority queue always removes the maximum value remaining in the priority queue (the one with the highest priority). Think of a line/queue outside of a Hollywood nightclub, such that whenever space opens up inside, the most famous person in line gets to go in (the "highest priority" person), no matter how long less famous people have been standing in line (contrast this with first come/first serve, which is a regular -non priority- queue; in a regular queue, whoever is first in the line -has been standing in line longest- is admitted next).

For the problems below, all we need to know is the complexity class of the "add" and "remove" operations.

	add	remove
Implementation 1	$O(1)$	$O(N)$
Implementation 2	$O(N)$	$O(1)$

Implementation 3		$O(\log N)$		$O(\log N)$	
		+-----+-----+			

Implementation 1 adds the new value into the pq by appending the value at the rear of a list or the front of a linked list: both are  $O(1)$ ; it removes the highest priority value by scanning through the list or linked list to find the highest value, which is  $O(N)$ , and then removing that value, also  $O(N)$  in the worst case (removing at the front of a list; at the rear of a linked list).

Implementation 2 adds the new value into the pq by scanning the list or linked list for the right spot to put it and putting it there, which is  $O(N)$ . Lists store their highest priority at the rear (linked lists at the front); it removes the highest priority value from the rear for lists (or the front for linked lists), which is  $O(1)$ .

So Implementations 1 and 2 swap the complexity classes in their add/remove method. Implementation 1 doesn't keep the values in order: so easy to add but hard to find/remove the maximum (must scan). Implementation 2 keeps the values in order: so hard to add (need to scan to find where it goes) but easy to find/remove the maximum (at one end).

Implementation 3, which is discussed in ICS-46, uses a binary heap tree (not a binary search tree) to implement both operations with "middle" complexity  $O(\log N)$ : this complexity class greater than  $O(1)$  but less than  $O(N)$ . Because  $\log N$  grows so slowly,  $O(\log N)$  is actually closer to  $O(1)$  than  $O(N)$  even though  $O(1)$  doesn't grow at all:  $\log N$  grows that slowly.

Problem 1: Suppose we wanted to use the priority queue to sort  $N$  values: we add  $N$  values in the pq and then remove all  $N$  values (first the highest, next the second highest, ...). Here is the complexity of these combined operations for each implementation.

Implementation 1:	$N \cdot O(1) + N \cdot O(N)$	$= O(N) + O(N^2)$	$= O(N^2)$
Implementation 2:	$N \cdot O(N) + N \cdot O(1)$	$= O(N^2) + O(N)$	$= O(N^2)$
Implementation 3:	$N \cdot O(\log N) + N \cdot O(\log N)$	$= O(N \log N) + O(N \log N)$	$= O(N \log N)$

Note  $N \cdot O(\dots)$  is the same as  $O(N) \cdot O(\dots)$  which is the same as  $O(N * \dots)$

Here, Implementation 3 has the lowest complexity class for the combined operations. Implementations 1 and 2 each do one operation quickly but the other slowly: both are done  $O(N)$  times. The slowest operation determines the complexity class, and both are equally slow. The complexity class  $O(\log N)$  is between  $O(1)$  and  $O(N)$ ; surprisingly, it is actually "closer" to  $O(1)$  than  $O(N)$ , even though it does grow -because it grows so slowly; yes,  $O(1)$  doesn't grow at all, but  $O(\log N)$  grows very slowly: the known Universe has about  $10^{90}$  particles of matter, and  $\log 10^{90} = \log(10^3)^{30} = 300$ , which isn't very big compared to  $10^{90}$  (like 86 orders of magnitude less).

Problem 2: Suppose we wanted to use the priority queue to find the 10 biggest (of  $N$ ) values: we would enqueue  $N$  values and then dequeue 10 values. Here is the complexity of these combined operations for each implementation..

Implementation 1:	$N \cdot O(1) + 10 \cdot O(N)$	$= O(N) + O(N)$	$= O(N)$
Implementation 2:	$N \cdot O(N) + 10 \cdot O(1)$	$= O(N^2) + O(1)$	$= O(N^2)$
Implementation 3:	$N \cdot O(\log N) + 10 \cdot O(\log N)$	$= O(N \log N) + O(\log N)$	$= O(N \log N)$

Here, Implementation 1 has the lowest complexity for the combined operations. That makes sense, as the operation done  $N$  times (add) is very simple (add to the end of a list/the front of a linked list, where each add is  $O(1)$ ) and the operation done a constant number of times (10, independent of  $N$ ) is the expensive operation (remove, which is  $O(N)$ ). It even beats the complexity of Implementation 3. So, as  $N$  gets bigger, implementation 1 will eventually become faster than the other two for the "find the 10 biggest" task.

So, the bottom line here is that sometimes there is NOT a "best all the time"

implementation for a data structure. We need to know what problem we are solving (the complexity classes of all the operations in various implementations and how often we must do these operations) to choose the most efficient implementation for solving the problem.

-----

Problems:

TBA