

A2 demo: if you don't have your mark posted, let me know the time and room.
A2 code: autotests have been run, the TAs should have it marked by Tuesday.

CSC148 winter 2017

(more...) efficiency
week 11

Danny Heap

heap@cs.toronto.edu / BA4270 (behind elevators)

<http://www.cdf.toronto.edu/~csc148h/winter/>
416-978-5899

March 29, 2017

Outline

big-Oh, Omega, Theta examples

hash tables

sequences

Binary digITS (BITS)

pretend that many operations (multiplication, division, comparison, etc.) are constant with respect to problem size (often called n), BUT for large enough n (man, many bits) this is not strictly true.

```
def silly(n):  
    n = 17 * n**(1/2)  constant wrt n  
    n = n + 3           addition constant wrt n  
    print("n is: {}".format(n))  constant wrt n  
  
    if n > 97:          constant...  
        print('big!')  constant  
    else:  
        print('not so big!')
```

How does the running time of silly depend on n ? constant

loops

How does the running time of this code fragment depend on n ?

```
sum = 0          linear (proportional to n)
                  Theta(n)
for i in range(n):
    sum += i      proportional to n    i = 0, 1, 2, ..., n-1
                  constant wrt n
```

How does the running time of this code fragment depend on n ?

```
sum = 0          constant wrt
for i in range(n//2):    i = 0, 1, 2, ..., n//2 - 1    n//2 values
    for j in range(n**2):    j = 0, 1, ..., n**2 - 1    n**2 values
        sum += i * j        (i, j) pairs? (n**3)//2
                              Theta(n**3)
```

Suppose it took x "steps" when $n = 10$
How many "steps" do you estimate when
 $n = 30$? --- $27x$



more loops

How does the running of this code fragment depend on n ?

```
i, sum = 0, 0    constant wrt n
while i**2 < n:  i = 0, 1, 2, ..., floor(sqrt(n) - 1) ... or sqrt(n)
    j = 0        constant wrt n
    while j**2 < n: j = 0, 1, 2, 3, ..., floor(sqrt(n) - 1) ... or sqrt(n)
        sum += i * j    constant    n (i, j)
        j += 1          constant
    i += 1              constant
```

Theta(n)

How does the running time of this code fragment depend on n ?

```
i, sum = 0, 0
while i < n * n:
    sum += i
    i += 1
```

conditions

How does the running time of this code fragment depend on n ?

```
sum = 0
if n % 2 == 0:
    for i in range(n*n):
        sum += 1
else:
    for i in range(5, n+3):
        sum += i
```

$O(n^2)$ # if n is even

$\Omega(n)$ # if n is odd

$i = 0, \dots, n^2 - 1$

condition: you need to consider both branches



halving

How does the running time of `twoness` depend on `n`?

```
def twoness(n):  
    count = 0  
    while n > 1:  
        n = n // 2  
        count = count + 1  
    return count
```

constant wrt n
 $n, n/2, n/4, n/8, \dots, 1$ $\text{floor}(\lg(n))$
constant wrt
constant
constant

x for n
 $x+1$ for $2n$, $x+2$ for $4n$, etc.

working with lg

$\lg(n)$: this is the number of times you can divide n in half before reaching 1.

- ▶ refresher: $a^b = c$ means $\log_a c = b$. definition of log
- ▶ this runtime behaviour often occurs when we “divide and conquer” a problem (e.g. binary search)
- ▶ we usually assume $\lg n$ (log base 2), but the difference is only a constant:

$$2^{\lg_2 n} = n = 10^{\lg_{10} n} \implies \lg_2 n = \lg_2 10 \times \lg_{10} n$$

... so logs of different bases are related by a constant...

- ▶ so we just say $\mathcal{O}(\lg n)$.

miscellaneous

How does the running time of this code fragment depend on n ?

constant --- n doesn't even appear in the code

```
for k in range(5000):  
    if L[k] % 2 == 0:  
        even += 1  
    else:  
        odd += 1
```

more miscellaneous

How does the running time of this code fragment depend on n and m ?

```
sum = 0
for i in range(n): i = 0, 1, ..., n-1 -> n values
    for j in range(m): j = 0, 1, ..., m-1 -> m values
        sum += (i + j)      nxm (i, j) pairs...
```

summary

sequences: `max`

loops: `count and multiply`

conditions: `consider each branch separately...`



why hash

lists are contiguous (adjacent) sequences of references to objects, so access to a list position is fast (just arithmetic)

list
id { [, , , , , , , , , ,]
list
id + 5,392,486

what if we could convert — hash — other data to a suitable integer for a list index, we'd want:

- ▶ fast no sense being slow
- ▶ deterministic: the same (or equivalent values) gets hashed to the same integer each time.
- ▶ well-distributed: We'd like a typical set of values to get hashed pretty uniformly over the available list positions.

don't want every value to hash to ... 42

you can't hash everything!

```
>>> list1 = [0]
>>> id(list1)
3069263116
>>> list2 = [0, 1]
>>> id(list2)
3069528300
>>> list1.append(1)
>>> id(list1)
3069263116
```

oops! two equivalent objects should *not* have the same hash!



hash to hash table (dictionary)...

`(key1, value1) , (key2, value2), ...`

Once you have hashed an object to a number, you can easily use part of that number as an index into a list to store the object, or something related to that object. If the list is of length n , you might store information about object o at index $\text{hash}(o) \% n$.

$$0 \leq \text{hash}(o) < n$$

collisions

328, 338, 37, 280, 181, 114, 257, 200, 250, 68, 242, 250 ... 12 "birthdays" in class led to collision!

even a well-distributed hash function will have a surprising number of collisions...

how many people do you need to poll before you find two with the same birthday (out of 366 possibilities, including leap-year)?

the mathematics is a bit counter-intuitive... the probability of a non-collision for 23 birthdays is:

$$p = \frac{366}{366} \times \frac{365}{366} \times \cdots \times \frac{344}{366} \approx 0.493$$

chaining or probing

a couple of tactics for dealing with two different keys ending up at the same index

"danny"	"david"	hash("danny") % 8
		hash("david") % 8 ...
[[], [], [], [], [], [], [], []]		

^{"jacqueline"}
chaining: keep a small (one hopes) list at that index
(sometimes called bucket)

probing: explore, in a systematic way, until the next open index

either tactic has costs, so keep collisions to a minimum by keeping the list partly empty

Python dictionaries are implemented¹ using hash tables and probing. The cost of collisions is kept small by enlarging the underlying list when necessary, and the cost of enlarging is amortized over many dictionary accesses.

The result is that access to a dictionary element is $\Theta(1)$, essentially the time it takes to access a list element.

One downside is that extra work is required to order the keys or values of a dictionary. What is their “natural” order?

¹Google “Tim Peters” but beware of obnoxious culture



Python dictionaries are implemented¹ using hash tables and probing. The cost of collisions is kept small by enlarging the underlying list when necessary, and the cost of enlarging is amortized over many dictionary accesses.

The result is that access to a dictionary element is $\Theta(1)$, essentially the time it takes to access a list element.

One downside is that extra work is required to order the keys or values of a dictionary. What is their “natural” order?

¹Google “Tim Peters” but beware of obnoxious culture

