

COMP 204

Algorithm design: Linear and Binary Search

Yue Li

based on material from Mathieu Blanchette, Christopher J.F.
Cameron and Carlos G. Oliver

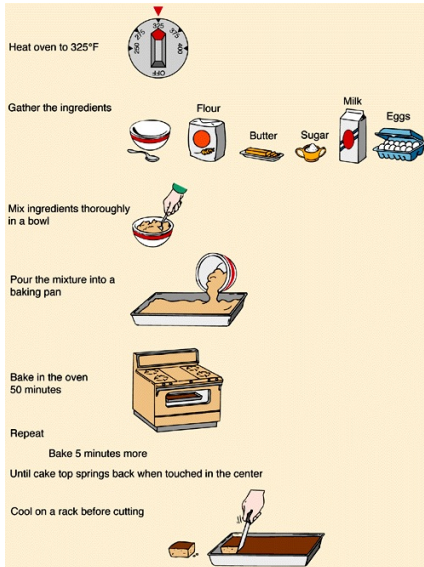
Algorithms

An **algorithm** is a predetermined series of instructions for carrying out a task in a finite number of steps

▶ or a recipe

Input \rightarrow algorithm \rightarrow output

Example algorithm: baking a cake



What is the input?

algorithm?

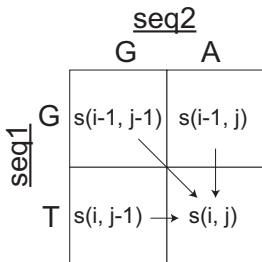
output?

Example algorithm: sequence alignment (A2)

Input: seq1, seq2

Output: alignments of seq1 and seq2

Algorithm:



$$s(i, j) = \max \begin{cases} s(i-1, j-1) + (\text{mis})\text{match} \\ s(i-1, j) + \text{gap} \\ s(i, j-1) + \text{gap} \end{cases}$$

- ▶ $s(i-1, j-1) + (\text{mis})\text{match}$: align letter seq1[i] with letter seq2[j] (match: +2, mismatch: -2)
- ▶ $s(i-1, j) + \text{gap}$: align a gap "-" from seq2 with seq1[i] (gap: -2)
- ▶ $s(i, j-1) + \text{gap}$: align a gap "-" from seq1 with seq2[j] (gap: -2)

Pseudocode

Pseudocode is a universal and informal language to describe algorithms from humans to humans

It is not a programming language (it can't be executed by a computer), but it can easily be translated by a programmer to any programming language

It uses variables, control-flow operators (while, do, for, if, else, etc.)

Example Python statements

```
1 students = ["Kris", "David", "JC", "Emmanuel"]
2 grades = [75, 90, 45, 100]
3 for student, grade in zip(students, grades):
4     if grade >= 60:
5         print(student, "has passed")
6     else:
7         print(student, "has failed")
8 #output:
9 #Kris has passed
10 #David has passed
11 #JC has failed
12 #Emmanuel has passed
```

Example pseudocode

Algorithm 1 Student assessment

```
1: for each student do  
2:   if student's grade  $\geq$  60 then  
3:     print 'student has passed'  
4:   else  
5:     print 'student has failed'  
6:   end if  
7: end for
```

Example algorithm: longest hydrophobic patch (L12)

Input:

amino acid
sequence

Output:

longest
hydrophobic
patch

Algorithm:

findLongestHydrophobicPatch(protein)

isHydrophobicPatch(sequence)?

EDAYQ**IALEGA**ASTE

outer for loop:
start position from
start = 0

inner for loop
end position from
end = start + 1

isHydrophobicPatch(sequence)?

isHydrophobic('E')
(1) first a.a.

isHydrophobic('L')
(2) last a.a.

↓ ↓
EDAYQ**IAL**

for-loop
patchLen += isHydrophobic(s[aa])
(3) length of hydrophobic amino acids (min 80%)

isHydrophobic(aa)?

aa in ["G", "A", "V", "L", "I", "P", "F", "M", "W"]?

findLongestHydrophobicPatch Python code

```
41 # This returns the longest hydrophobic patch found in a sequence
42 def findLongestHydrophobicPatch(protein):
43     longestPatch="" # the longest patch found so far
44
45     # for every possible starting point
46     for start in range(0,len(protein)):
47
48         # and every possible end point
49         for end in range(start+1,len(protein)+1):
50             # get the sequence
51             candidate = protein[start:end]
52
53             # test hydrophobicity
54             if isHydrophobicPatch(candidate):
55
56                 # if longer than longest seen so far, update
57                 if len(candidate)>len(longestPatch):
58                     longestPatch = candidate
59
60     return longestPatch
```

findLongestHydrophobicPatch pseudocode

Algorithm 2 findLongestHydrophobicPatch

```
1: while start position < protein length do
2:   end position  $\leftarrow$  start position + 1
3:   while end position < protein length do
4:     candidate  $\leftarrow$  protein substring from start to end position
5:     if candidate is hydrophobic patch then
6:       if length(candidate) > length(longestHydroPho)
7:       then longestHydroPho  $\leftarrow$  candidate
8:     end if
9:   end while
10:  end position  $\leftarrow$  end position + 1
11: end while
```

Search algorithms

Search algorithms locate an item in a data structure

Input: a list of (un)sorted items and value of item to be searched

Algorithms: linear and binary search algorithms will be covered

- ▶ images if search algorithms taken from:
http://www.tutorialspoint.com/data_structures_algorithms/

Output: if value is found in the list, return index of item

Example:

- ▶ search (key = 5, list = [3, 7, 6, 2, 5, 2, 8, 9, 2]) should return 4.
- ▶ search (key = 1, list = [3, 7, 6, 2, 5, 2, 8, 9, 2]) should return nothing.

Linear search

A very simple search algorithm

- ▶ a sequential search is made over all items one by one
- ▶ every item is checked
- ▶ if a match is found, then index is returned
- ▶ otherwise the search continues until the end of the sequence

Example: search for the item with value 33



Linear search #2

Starting with the first item in the sequence:

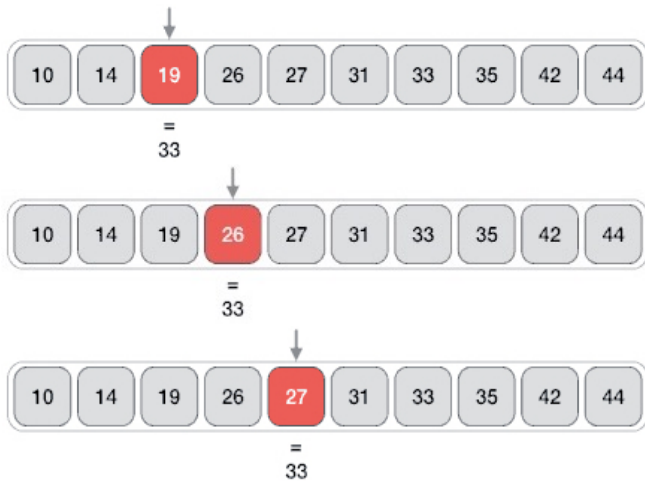


Then the next:



Linear search #3

And so on and so on...



Linear search #4

Until an item with a matching value is found:



If no item has a matching value, the search continues until the end of the sequence

Linear search: pseudocode

Algorithm 3 Linear search

```
1: procedure LINEAR_SEARCH(sequence, key)
2:   for index = 0 to length(sequence) do
3:     if sequence[index] == key then
4:       return index
5:     end if
6:   end for
7:   return None
8: end procedure
```

Linear search: Python implementation

```
1 def linear_search(sequence, key):
2     for index in range(0, len(sequence)):
3         if sequence[index] == key:
4             return index
5     return None
6
7 # import random
8 # L = random.sample(range(1,10**9),10**7)
9 # import time
10 # time_start = time.time()
11 # print(f"start: {time.asctime(time.localtime(time_start))}")
12 # index = linear_search(L, -1)
13 # time_finish = time.time()
14 # print(f"end: {time.asctime(time.localtime(time_finish))}")
15 # print("time taken (seconds):", time_finish-time_start)
```

Issues with linear search

Running time: If the sequence to be searched is very long, the function will run for a long time.

Example: The list of all medical records in Quebec contains more than 8 Million elements!

Much of computer science is about designing *efficient* algorithms, that are able to yield a solution quickly even on large data sets.

See experimentation on Wing...

Binary search

A fast search algorithm (compared to linear)

- ▶ the sequence of items must be sorted
- ▶ works on the principle of 'divide and conquer'

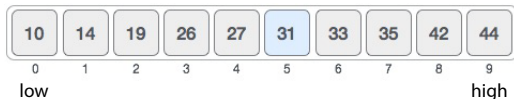
Analogy: Searching for a word (called the key) in an English dictionary.

To look for a particular word:

- ▶ Compare the word in the middle of the dictionary to the key
- ▶ If they match, you've found the word! Stop.
- ▶ If the middle word is greater than the key, then the key is searched for in the left half of the dictionary
- ▶ Otherwise, the key is searched for in the right half of the dictionary
- ▶ This repeated halves the portion of the dictionary that needs to be considered, until either the word is found, or we've narrowed it down to a portion that contains zero word, and we conclude that the key is not in the dictionary

Binary search #2

Example: let's search for the value 31 in the following sorted sequence



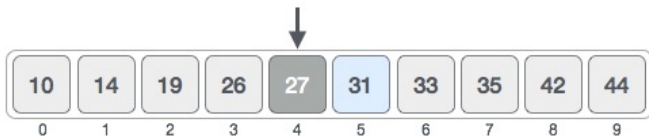
First, we need to determine the middle item:

```
1 sequence = [10, 14, 19, 26, 27, 31, 33, 35, 42, 44]
2 low = 0
3 high = len(sequence) - 1
4 mid = low + (high-low)//2      # integer division
5 print (mid) # prints: 4
```

Binary search #3

Since $index = 4$ is the midpoint of the sequence

- ▶ we compare the value stored (27)
- ▶ against the value being searched (31)



The value at index 4 is 27, which is not a match

- ▶ the value being search is greater than 27
- ▶ since we have a sorted array, we know that the target value can only be in the upper portion of the list

Binary search #4

low is changed to *mid* + 1



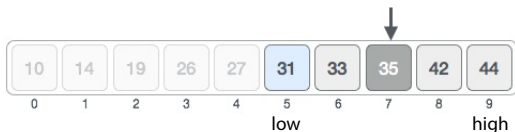
Now, we find the new *mid*

```
1 low = mid + 1    # 5
2 mid = low + (high-low)//2    # integer division
3 print (mid) # prints: 7
```

Binary search #4

mid is 7 now

- ▶ compare the value stored at index 7 with our value being searched (31)



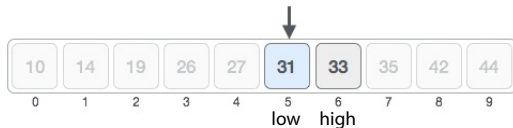
The value stored at location 7 is not a match

- ▶ 35 is greater than 31
- ▶ since it's a sorted list, the value must be in the lower half
- ▶ set *high* to *mid* - 1

Binary search #5

Calculate the mid again

- *mid* is now equal to 5



We compare the value stored at index 5 with our value being searched (31)

- It is a match!



Binary search #6

Remember,

- ▶ binary search halves the searchable items
- ▶ improves upon linear search, but...
- ▶ requires a sorted collection

Useful links

bisect - Python module that implements binary search

- ▶ <https://docs.python.org/2/library/bisect.html>

Visualization of binary search

- ▶ <http://interactivepython.org/runestone/static/pythonds/SortSearch/TheBinarySearch.html>

Binary search: pseudocode

Algorithm 4 Binary search

```
1: procedure BINARY_SEARCH(sequence, key)
2:   low = 0, high = length(sequence) - 1
3:   while low ≤ high do
4:     mid = (low + high) / 2
5:     if sequence[mid] > key then
6:       high = mid - 1
7:     else if sequence[mid] < key then
8:       low = mid + 1
9:     else
10:      return mid
11:    end if
12:  end while
13:  return 'Not found'
14: end procedure
```

Binary search: Python implementation

```
1 def binary_search(sequence, key):
2     low = 0
3     high = len(sequence) - 1
4     while low <= high:
5         mid = (low + high)//2
6         if sequence[mid] > key:
7             high = mid - 1
8         elif sequence[mid] < key:
9             low = mid + 1
10        else:
11            return mid
12    return None
```

Linear vs Binary search efficiency

Try `linear_and_binary_search.py` to see for yourself the difference in running time for large lists!

For a list of 100 Million elements, linear search takes about 3 seconds, and binary search takes about 0.001 seconds binary search is more than 3,000 times faster than linear search.

In general,

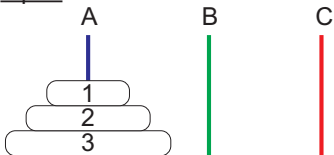
- ▶ the running time of linear search is proportional to the length of the list being searched.
- ▶ the running time of linear search is proportional to the **logarithm** of the length of the list being searched.

Binary search versus Linear search

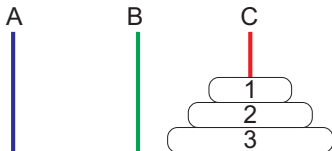
```
1 import random
2 import time
3 from decimal import Decimal
4 from linear_search import linear_search
5 from binary_search import binary_search
6
7 # generate list of 100 Million elements,
8 # where each element is a random number between 0 and 100,000,000
9 print("Generating list...")
10 n = 10**7
11 L = random.sample(range(10**9), n)
12
13 L.append(876567) # for testing purpose
14
15 print("Sorting list...")
16 L.sort()
17
18 key = int(input("Enter key for linear search: "))
19
20 # perform linear search
21 print("Starting linear search ...")
22 time_start = time.time()
23 index = linear_search(L, key)
24 time_finish = time.time()
25 linear_search_time = time_finish-time_start
26 print(f"Found at position: {index}; time taken:", "{:.2e}".format(linear_search_time), "seconds")
27
28 print("Starting binary search ...")
29 time_start = time.time()
30 index = binary_search(L, key)
31 time_finish = time.time()
32 binary_search_time = time_finish-time_start
33 print(f"Found at position: {index}; time taken:", "{:.2e}".format(binary_search_time), "seconds")
```

Example algorithm: Tower of Hanoi (Advanced)

Input:



Output:



Rules:

- ▶ Only one peg can be moved at a time
- ▶ Take the top disk from one of the stacks and place it on top of another stack or empty rod
- ▶ No larger disk may be placed on top of a smaller disk

Algorithm (recursive):

- ▶ Move $n - 1$ disks from source peg to spare peg
- ▶ Move m^{th} disk from the source to the target peg
- ▶ Move the $n - 1$ disks from spare peg to the target peg

See `tower_of_hanoi.py`

tower_of_hanoi Python code (Advanced)

```
1  def move(n, source, target, spare):
2      if n > 0:
3          # move n - 1 disks from source to spare
4          move(n - 1, source, spare, target)
5
6          # move the nth disk from source to target
7          target.append(source.pop())
8
9          # Display our progress
10         print(A, B, C, '#####', sep = '\n')
11
12         # move the n - 1 disks that we left on spare onto target
13         move(n - 1, spare, target, source)
14
15     # initiate call from source A to target C with spare B
16     A = [3, 2, 1]
17     B = []
18     C = []
19
20     move(3, A, C, B)
```

tower_of_hanoi pseudocode (Advanced)

Algorithm 5 Tower of Hanoi Mover

- 1: Move $n - 1$ disks from source peg to spare peg
 - 2: Move the n^{th} disk from the source to the target peg
 - 3: Move the $n - 1$ disks from spare peg to the target peg
 - 4: Do nothing if no disk left on source and spare peg
-