# Problem Set 1 – Introduction to Algorithms

1) Functions in order of growth:
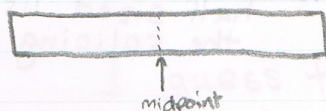
$$1/n, \; 1/5, \; 10^{100}n, \; n\log n, \; n^{100}, \; 3^{\sqrt{n}}, \; 3^n, \; 4^n, \; 2^{2n},$$
$$\log(n!), \; \binom{n}{100}$$

Equivalence Classes:

constant   linear   loglinear   polynomial

$$\boxed{1/n}, \; \boxed{1/5}, \; \boxed{10^{100}n}, \; \boxed{n\log n}, \; \boxed{n^{100}},$$

$$\boxed{3^{\sqrt{n}}}, \; \boxed{3^n}, \; \boxed{4^n, \; 2^{2n}}, \; \boxed{\log(n!)}, \; \boxed{\binom{n}{100}}$$

exponential          factorial

2) a) Iterative version runtime:



↑
midpoint

$$T(n) = \begin{cases} \Theta(1) & \text{if } n = 1 \\ T(n/2) + \Theta(1) & \text{if } n > 1 \end{cases}$$

Solve recurrence relation using substitution method:

1) Guess the order of the algorithm is:

$$\Theta(\lg n)$$

2) Show: $T(n) \leq c \cdot \lg n$, for $n \geq n_0$

    i) Assume $T(n/2) \leq c\lg(n/2)$
    ii) $T(n) \leq c\lg(n/2) + 1$
$$= c\lg(n) - c\lg(2) + 1$$
$$= c\lg(n) - c + 1$$

Try again: Guess $T(n) \leq c\lg n - b$

$$T(n) \leq c\lg(n/2) + 1$$
$$T(n) \leq c\lg n - c + 1$$

which for $b = 1 - c$

$$T(n) \leq c\lg n - b$$

$$\therefore T(n) = O(\lg n)$$

b) Recursive Version Runtime:

$$T(n) = \begin{cases} \Theta(1) & \text{if } n = 1 \\ T(n/2) + n/2 & \text{if } n > 1 \end{cases}$$

Solve subproblem of half

cost of creating new, half-sized list via the splicing operator

Solve Recurrence relation using substitution

1) Guess: $T(n) = O(n\lg n)$

2) Show: $T(n) \leq cn\lg n$ for $n > n_0$

   i) Assume: $T(n/2) \leq c(n/2)\lg(n/2)$

   ii) Substitute:

$$T(n) \leq c(n/2)\lg(n/2) + n/2$$
$$\leq cn\lg(n/2) + n$$
$$= cn\lg n - cn\lg 2 + n$$
$$= cn\lg n - cn + n$$
$$\boxed{\leq cn\lg n}$$

$$\therefore T(n) = O(n\lg n)$$

c) The extra running time of the recursive function is because the splicing operation in Python creates a new array with $n/2$ elements (which I assume has been implemented to run in $O(n)$ time.

A better implementation would be:

def binary Search (alist, item, first=0, last:

which would create no new instances of the array, like the iterative version.

3) a) Hypothesis for the order of set.intersection

$$O(\arg\min(|s|, |t|))$$

I guess that the order of the function intersection would be linear on the smaller of the two sets since any intersection command would need to run __contains__ at least once for each element in the smaller list.

b)

| time in $\mu s$ | $|s|=10^3$ | $|s|=10^4$ | $|s|=10^5$ | |
|---|---|---|---|---|
| $|t|=10^3$ | 0.003 | 0.003 | 0.003 | 0 |
| $|t|=10^4$ | 0.003 | 0.03 | 0.04 | 8 |
| $|t|=10^5$ | 0.003 | 0.04 | 0.4 | 1 |
| $|t|=10^6$ | 0.004 | 0.06 | 1.4 | 1 |

The above values (in seconds) are the total of running timeit on s.intersection(t) 100 times for varying values of $|s|$ and $|t|$

c) Based on this experimentally observed
   data the order is:

$$O(n), \text{ where } n = \min(|s|, |t|)$$

d) Yes, the new metric should run slower
   because the addition of a domain check
   was redundant, it did not change the solution
   but added an additional n-dependent operation
   for every element comparison.
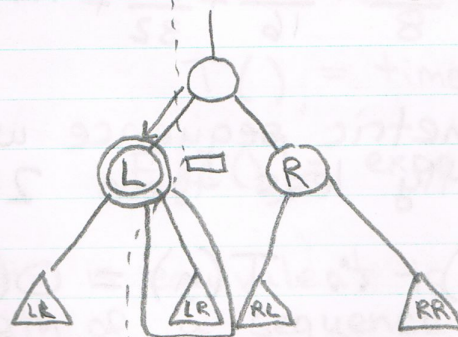
# Problem Set 2 – Introduction to Algorithms

1) Implemented in bst select.py

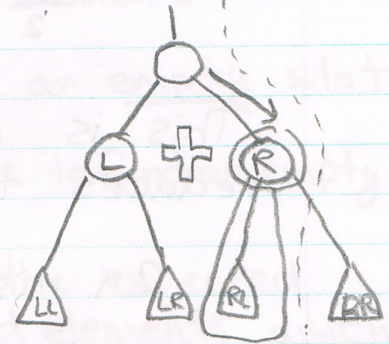   Algorithm description:

```
Set node to root
Set rank = node.left's size
while node isn't None:
    if index == rank:
        return Node
    else if index < rank:
        node = node.left
        rank = rank - node.right.size - 1
    else
        node = node.right
        rank = rank + node.left.size + 1
```
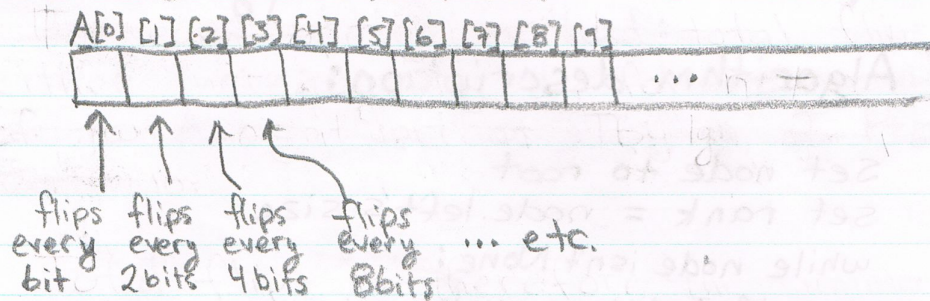


LEFT TRAVERSAL                    RIGHT TRAVERSAL

(dotted line is rank separator)

2) Amortized Analysis:

   Prove $T(n) = O(n)$ for $T(n)$:
   being the number of operations to
   increment a binary counter n times.

Bits flip depending on the k-value
for the $k^{th}$ increment:

A[0] [1] [2] [3] [4] [5] [6] [7] [8] [9]

↑   ↑   ↑   ↑
flips flips flips flips
every every every every ... etc.
bit 2bits 4bits 8bits

Therefore the total cost of k increments
is a summation of the number of flips
for each bit:

$$n + \frac{n}{2} + \frac{n}{4} + \frac{n}{8} + \frac{n}{16} + \frac{n}{32} + \dots$$

This is a geometric sequence which
reduces to strictly less than 2n.

$2n$ is $O(n)$   ∴ $T(n) = O(n)$   QED

3) a) $P(\text{chain of 3 after 3 hashes}) = \left(\frac{1}{m}\right)^3$

    Assuming the keys are independent, each
has a probability of hashing to one locatio
of $1/m$, the probability of 3 hashes occurin
is $1/m \cdot 1/m \cdot 1/m$ or $(1/m)^3$

b) This probability is also $(1/m)^3$ because
it amounts to the same thing, slot n
is hashed 3 times in a row.

c) Load Factor $= \alpha = 1 - \dfrac{1}{\log n}$

$\alpha = n/m =$ ratio of elements to slots

$T(\text{unsuccessful search}) = 1 + \underbrace{\text{average chain length}}_{\alpha}$

$$T(n) = 2 - \dfrac{1}{\lg n}$$

d)



$\alpha =$ proportion of array filled with hashed elements

$T(\ ) =$ time to find an **empty** slot

$E[T(\ )] =$ expected time to find an empty slot

This is equivalent to asking the expected length of a sequence of filled elements which has a proportion $\alpha$.

$E[T(n)] = i, \quad \underbrace{\alpha^i = 0.5}$

This is the probability where a longer and shorter subsequence are equal

$\left(1 - \dfrac{1}{\lg n}\right)^i = 0.5 \qquad i = \dfrac{0.5}{\lg\left(1 - \dfrac{1}{\lg n}\right)}$

$i \lg\left(1 - \dfrac{1}{\lg n}\right) = 0.5 \qquad E[T(n)] = \dfrac{1}{2\lg\left(1 - 1/\lg\right)}$

4)a) Ben's algorithm runs 4 nested loops:

    1) for s_start in range(0, len(s))
      ↳ This is $O(n)$

    2) for s_end in range(s-start, len(s))
      ↳ This has an amortized cost of $n/2$, also $O(n)$

    3) for t_start in range(0, len(t))
      ↳ This is $O(n)$

    4) "           "
      ↳ $n/2$, $O(n)$

The algorithm runs in $(n)(\frac{n}{2})(n)(\frac{n}{2}) = \frac{n^4}{4}$ or O⟨ time.

b) Alyssa's algorithm requires 2 nested for loops, plus a call to the Python **in** operator for lists:

1) for length in range(min(|s|, |t|))
    ↳ $O(n)$

2) for s_start ...  → $O(n)$
   for t_start ... → $O(n)$ ⎤ $+ 2*O(n) = O(n)$

3) if current in s_substrings
    ↳ I will assume this uses linear searching and therefore is $O(n)$

Therefore the total algorithm runs in $O(n^3)$ time.

c) Implemented in substring2.py

Initial (Unmodified)          $O(n^2 \lg n)$ modification
T(100) = 0.0003719.              0.0003859
T(1000) = 0.028105               0.00869703
T(10000) = 3.7917039             0.45809197
T(100000) = 204.160919           180.81293988

d) Implemented in substring4.py using Rabin-Karp:

    $O(n^2 \lg n)$ bin-search          $O(n \lg n)$ Rabin Karp

T(100)
T(1000)
T(10000)
T(100000)

# Problem Set 3 – Introduction to Algorithms

1) Implemented in heap_delete.py

Code uses first decrease_key which sets the index to $-\infty$ and pushes it to the top (maintaining minheap invariant). Then, the extract_min function is used to remove the element. Both operations run in $O(\lg n)$ so the combined function is also $O(\lg n)$.

2) a) If the MPQ were implemented with a max heap, then m operations would have a worst-case time of $O(m \lg m)$ since extract and insert would each have $O(\lg m)$ time in the worst case.

b) Here, a length k array can be generated, every insert operation would check against a stored max value, updating only if the value $\leq$ max. Updates are done by setting MPQ[value] += 1. Extract max removes one from the MPQ[max], and if this is now zero walks back down the array until it can find a new max.

The process takes $O(m+k)$ for m operations because the total cost of all walk-backs to find a new max is at most k, since the max can only decrease.

3) a) Detect collisions is $O(n^2)$ since every ball is compared against every other. (Upper triangle in a square matrix)

b) Sorting cost $= O(n)$

Collision Check $= O(mr^2)$ where $r = E[\text{\# of balls/b}$
and $m = \text{\# of bins}$

In the worst case (all balls in the same bin) thi performs the same as the original algorith

On the average case, however, $r = \dfrac{n}{(400\sqrt{n})^2 / 256^2}$

$$r = \dfrac{n}{\dfrac{160\,000\,n}{65536}} \doteq \dfrac{1}{2.44}$$

Since $m \doteq 2.44n$, this means the collision check should take $O(n)$ time to process, albeit with higher constant coefficients.

c) Implemented in detection.py and teste

PROBLEM SET 4: Introduction to Algorithms

1) Give an $O(V+E)$ time algorithm for separating a graph into connected components.

    First, create a set of all vertices, ordered by some easy process (such as position in array or matrix). Next, go through every edge. Compare the two nodes incident to the edge, keep the smaller of the two and delete (or mark) the larger one. In the end, the graph will have only $c$ unmarked nodes which are the smallest nodes in each connected component.
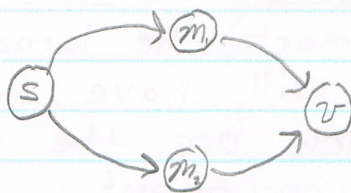
    If the purpose was to partition the graph, then instead of marking, nodes could be added to a list containing all connected nodes and stored by the "min" node.

    Since only one pass through each node and each edge is required, this algorithm requires $O(V+E)$ time.

2) a) $v$ is part of a cycle if a path connects it with itself. To implement this algorithm in $O(V+E)$ time, run a BFS starting at $v$, never color the node $v$, so future edge checks can loop back. Then if $v$ would be colored a second time, we can stop and return true. If we reach the end of all nodes, we can return false.

b) Ben's algorithm works for undirected
   graphs because going backwards from
   an edge is always possible, therefore
   the only graphs without cycles are trees
   which never repeat nodes in a BFS.

   In a digraph this process would produce
   false positives. Take the simple example belo



   There is no cycle here because $v$ has
   no outgoing edges, yet it _is_ seen twice
   from $m_1$ and $m_2$.

c) First, we conduct a DFS, but also
   give a second mark to each node (*)
   when it is discovered and remove this
   mark whenever we backtrack from a
   black node. This (*) has the property
   of saying which nodes lie in a direct
   path from s, the starting node and $v$,
   the current vertex.
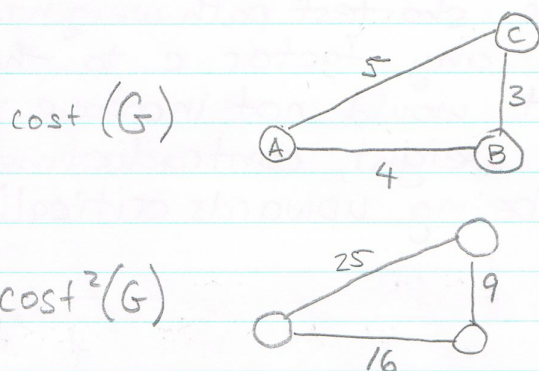
   Cycles can then be detected as occurring
   whenever a node finds a (*) node
   during its search.

3) Implemented in level.py and solver.py

# PROBLEM SET 5: INTRODUCTION TO ALGORITHMS

1) a) False. This is only true if the graph contains no negative edge weight cycles, as this will leave the shortest path undefined, but will return an answer nonetheless with this modification.

b) False. Consider the following graph:

cost (G)



In this example, the path AC is strictly short. than ABC, yet they have the same length with all edge lengths squared.

cost²(G)

c) True. This algorithm will return the largest simple path, in an acyclic graph. The relate problem of finding a longest simple path with cycles is not solvable with this algorithm and is in NP-complete.

2) a) Proof by Contradiction:

Suppose edge $(u, v)$ were downward critica and it was not on a shortest path betwee S and $v$. This means the total path weight including $(u, v)$ is strictly more than the shortest path weight by some factor c. We could then subtract some amount ac, $0 < a < 1$ from $(u, v)$ and it would still be longer than the shortest path which contradicts the definition of being downward critical.

b) Claim: All upward critical edges $(u, v)$ are on the shortest path between $s$ and $v$.

Proof by Contradiction:

Again, suppose we had an upward critical path that weren't on the shortest path. This means the path including $(u, v)$ is strictly larger than the shortest path weight. Here, we could add any factor $c$ to the weight of $(u, v)$ and it would not increase the shortest path weight, contradicting the principle of being upwards critical.

c)

3) Implemented in dijkstra.py

## PROBLEM SET 6: Introduction to Algorithms

1) Implemented in fib.py

2) a) The set of subproblems is the ways of making change with $(1, 2, 3, \ldots, C-1)$ as the amounts of change.

b) Let $SCL(c)$ represent the function for the list of the shortest list of change, given an amount of change, $c$.

$$SCL(c) = \min \begin{pmatrix} SCL(c-d_1), \\ SCL(c-d_2), \\ \cdots \\ SCL(c-d_m) \end{pmatrix} + 1$$

c) The running time of this algorithm is

$O(mC)$ as there are $C$ subproblem each requiring $m$ (# of denomina calculations.

d) Implemented in change.py

3) a) The set of subproblems are the longest
increasing subsequences, ending at index $i \in 0, 1, \ldots$
The LIS of the first $i$ elements is necessarily
$1 +$ the maximum of the LIS which end at $j$, which
is such that $j < i$, $x[i] < x[j]$, given $X$ as the sequence

b) Let $q_i$ be the LIS which ends at element $i$,
the recurrence relation is therefore:

$$q_i = \max \text{length} \left( q_j \mid j < i, x_j < x_i \right) + 1$$

From this the true LIS is simply:

$$LIS = \max \text{length} \left( q_i \mid i \in \{0, 1, 2, \ldots, n\} \right)$$

c) This algorithm requires $O(n^2)$ time because
there are $n$ problems each having $n$ subproblems

d) Implemented in progress.py.

4) Implemented in ResizeableImage.py

 ↳ Unfortunately, I couldn't get PIL working
so this code is untested and may contain
bugs.