

Problem Set 0

Harvard SEAS - Fall 2022

Due: Wed Sep. 7, 2022 (11:59PM)

Collaboration Statement: Done in collaboration with Avi Gulati, Liya Jin, Alexandra Topic, Eric Elliott, and Jacob Cremers.

The purpose of this problem set is to reactivate your skills in proofs and programming from CS20 and CS32/CS50. For those of you who haven't taken one or both those courses, the problem set can also help you assess whether you have acquired sufficient skills to enter CS120 in other ways and can fill in any missing gaps through self-study. Even for students with all of the recommended background, this problem set may still require a significant amount of thought and effort, so do not be discouraged if that is the case and do take advantage of the staff support in section and office hours.

For those of you who are wondering whether you should wait and take CS20 before taking CS120, we encourage you to also complete [the CS20 Placement Self-Assessment](#). Some problems there that are of particular relevance to CS120 and are complementary to what is covered below are Problems 2 (counting), 4 (comparing growth rates), 9 (quantificational logic), and 12 (graph theory).

Written answers must be submitted in pdf format on Gradescope. Although L^AT_EX is not required, it is strongly encouraged. You may handwrite solutions so long as they are fully legible. The `ps0` directory, which contains your code for problems 1a and 1c, must be submitted separately to an autograder on Gradescope. Be sure to pull the starter code from the [cs120 GitHub repository](#).

1. (Binary Trees) In the [cs120 GitHub repository](#), we have given you a Python implementation of a binary tree data structure, as well as a collection of test trees built using this data structure. We specify a binary tree by giving a pointer to its *root*, which is a special *vertex* (a.k.a. *node*), and giving every vertex pointers to its *children* vertices and its *parent* vertex as well as an identifying *key*:

```
class BinaryTree:
    def __init__(self, root):
        self.root: BTvertex = root

class BTvertex:
    def __init__(self, key):
        self.parent: BTvertex = None
        self.left: BTvertex = None
        self.right: BTvertex = None
        self.key: int = key
        self.size: int = None
```

In CS50, the concept of a Python `class` was not covered. Here, with `BinaryTree` and `BTvertex`, we are using them in the same way as a `struct` in C. An object `v` of the `BTvertex`

class contains five attributes, which we list with the type of the object we expect to be named by each attribute (using the Python type annotation syntax). These attributes can be accessed as `v.parent`, `v.left`, `v.right`, `v.key`, and `v.size`. For example, `v.left.key` is the key associated with `v`'s left child. An object of the `BinaryTree` class contains only one attribute, which is the `BTvertex` object that is the root of our binary tree. You can create a `BinaryTree` object as follows:

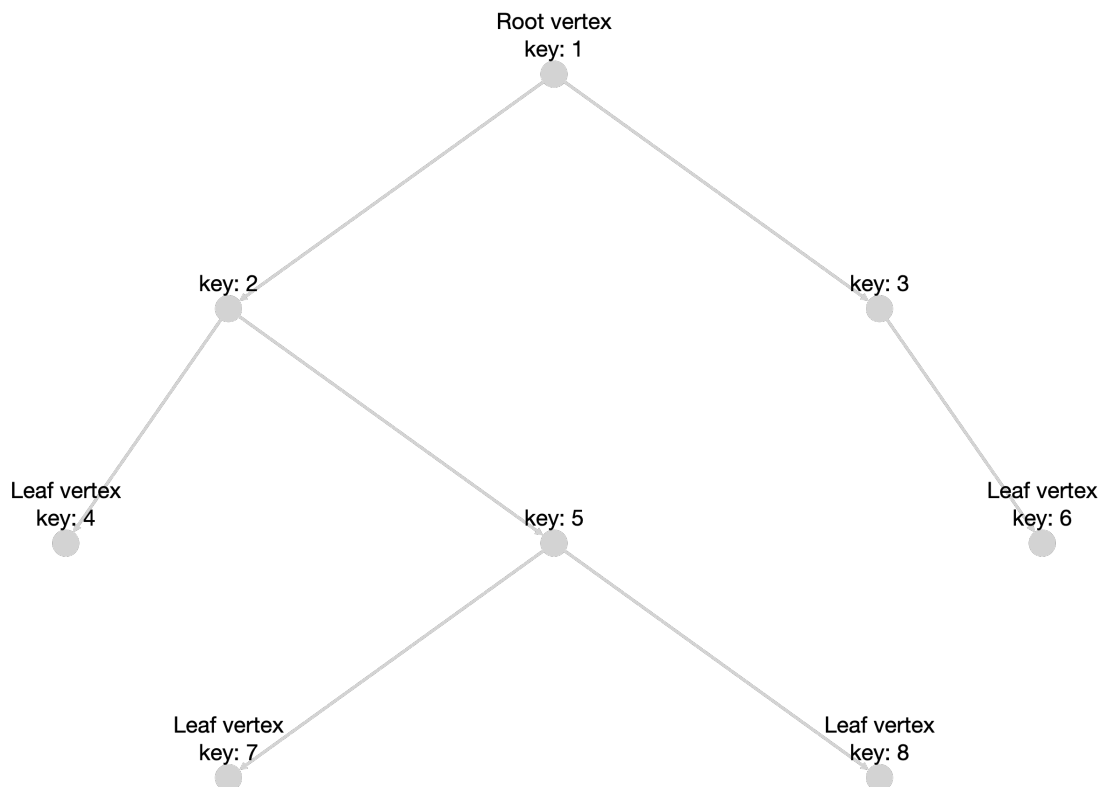
```
root = BTvertex(120)
tree = BinaryTree(root)
tree.root.left = BTvertex(121)
tree.root.right = BTvertex(124)
```

You can then print attributes of the newly created `BinaryTree` object:

```
print(tree.root.key)
>> 120
print(tree.root.left.key)
>> 121
```

Classes are more general than structs because they can also have private attributes and methods that operate on the attributes, allowing for object-oriented programming. However, you won't need that generality in this problem set.

Here is an instance `T` of `BinaryTree`:



A `BinaryTree` `T` contains only a pointer to its root vertex, `T.root`, which is required to satisfy `T.root.parent==None`. In the above example, the root is the vertex with key 1 (i.e. `T.root.key==1`). A binary tree vertex `v` can have zero, one, or two children, determined by which of `v.left` and `v.right` are equal to `None`. In the above example, the vertex `v` with key 3 has `v.left==None` but `v.right` is the vertex with key 6. A *leaf* is a vertex with zero children, i.e. `v.left==v.right==None`.

A vertex `w` is *descendent* of a vertex `v` if there is a sequence of vertices v_0, v_1, \dots, v_k , $k \in \mathbb{N}$ such that $v_0 = v$, $v_k = w$, and $v_i \in \{v_{i-1}.left, v_{i-1}.right\}$ for $i = 1, \dots, k$.¹ In the above example, the vertex with key 5 is a descendent of the root (with a path of length 2), but is not a descendent of the vertex with key 3. The sequence v_0, v_1, \dots, v_k is called a *path* from `v` to `w` and k is the *distance* from `v` to `w`. Taking $k = 0$, we see that `v` is a descendent of itself.

The *vertex set* of a binary tree `T` consists of all of the descendents of `T.root`. The *size* of `T` is its number of vertices. The *height* of `T` is the largest distance from the root to a leaf. The above example has size 8 and height 3.

Given any vertex `v` in a tree, the *subtree* rooted at `v` consists of all of `v`'s descendents. Note that we can remove a subtree and turn it into a new tree `S` by setting `S.root=v` and `v.parent=None`.

For now, the **key** attribute serves to distinguish vertices from each other in our tests and help illustrate what the algorithms are doing. The `BTvertex` class also has a **size** attribute, which is initialized to `None` in all of the test instances; it will be filled in by the program you write in Part 1a.

An instance `T` `BinaryTree` is *valid* if it satisfies the following constraints:

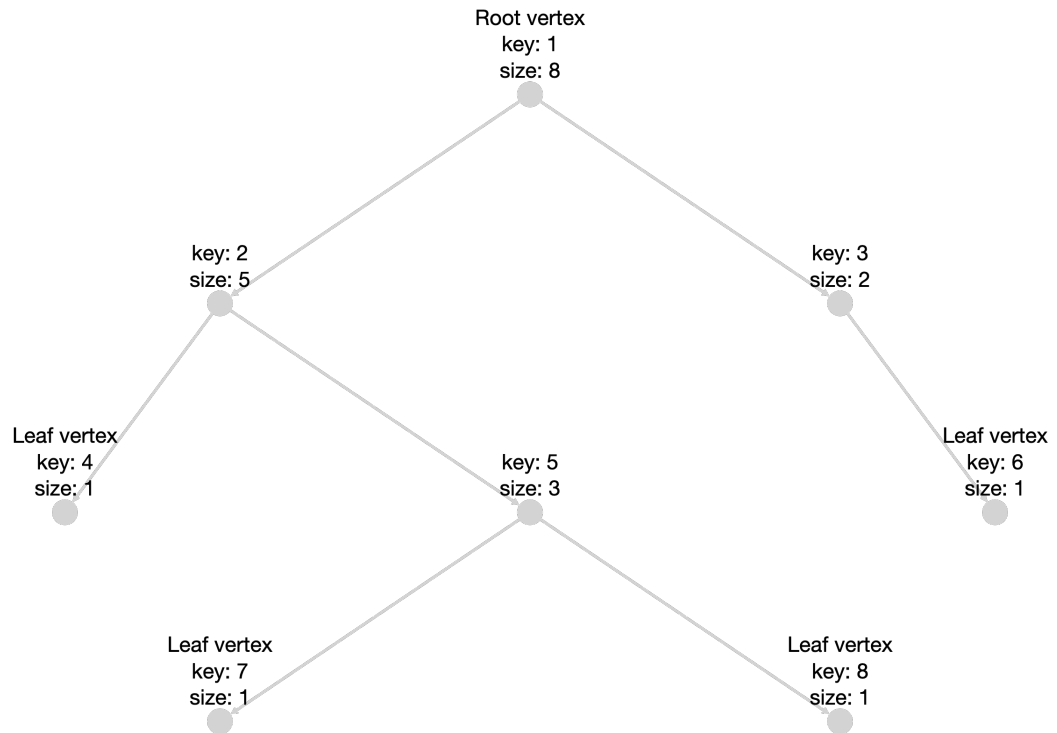
- `T.root.parent==None`
- `T` has finitely many vertices.
- No two vertices `v`, `w` of `T` share a child, i.e. $\{v.left, v.right\} \cap \{w.left, w.right\} = \emptyset$.

All of the test instances we provide are valid, and furthermore have the property that all of the vertices have distinct keys (which is something we often want, but not always).

- (a) (recursive programming) Write a recursive program `calculate_sizes` that given a vertex `v` of a binary tree `T`, calculates the sizes of all of the subtrees rooted at descendents of `v`. After running your program on `T.root`, every vertex `v` in `T` should have `v.size` set to the size of the subtree rooted at `v`. (Recall that the size attributes are initialized to `None`.) We call the resulting tree a *size-augmented* tree.

For example, if `T` is the tree shown above, then calling `calculate_sizes(T.root)` should modify `T` to be the following size-augmented tree:

¹ \mathbb{N} denotes the natural numbers $\{0, 1, 2, 3, \dots\}$. Since we are computer scientists, we start counting at 0.

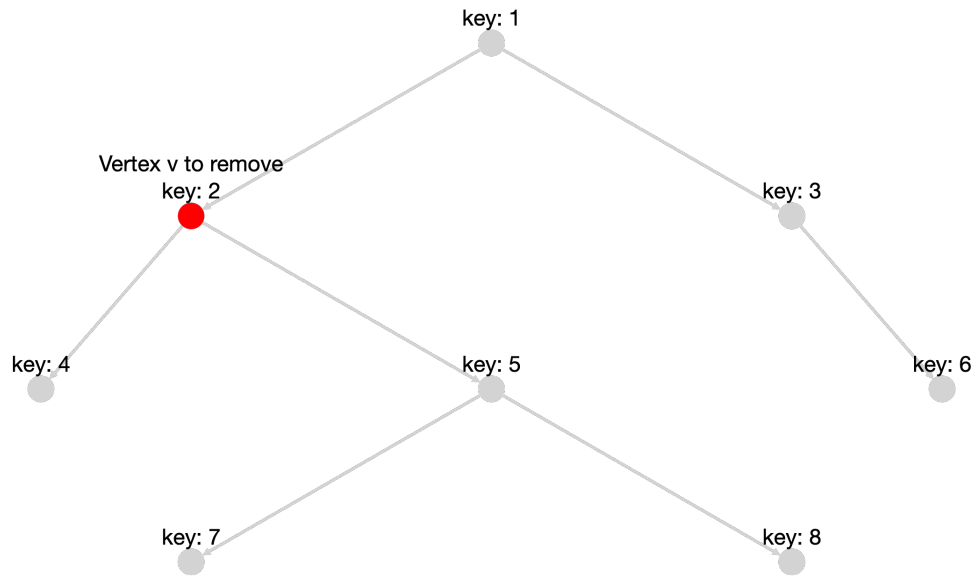


Your program should run in time $O(n)$ when given the root of a tree with n vertices. In a sentence or two, informally justify why your program has such a runtime.

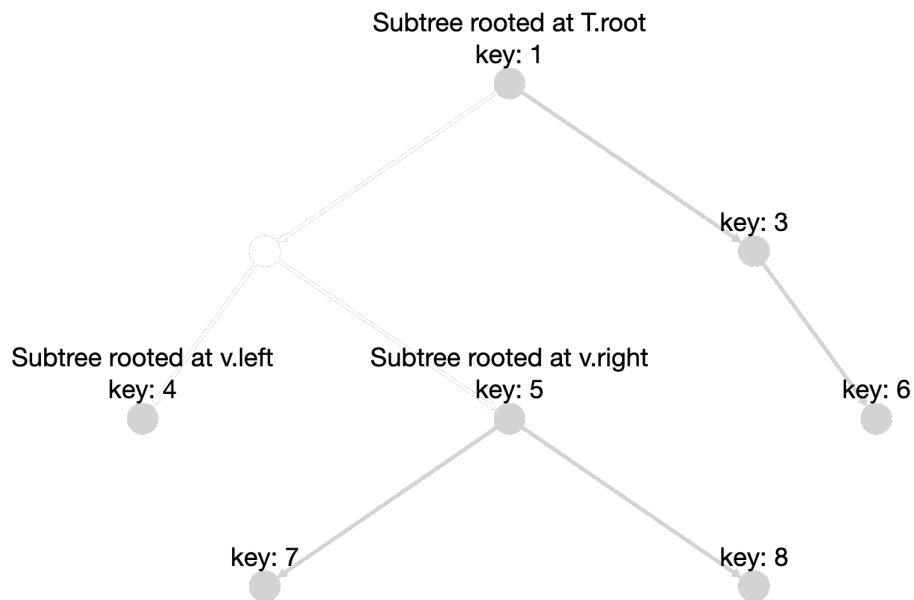
Answer: Reasoning provided within Python document.

- (b) (proofs by contradiction) Removing a vertex v from a tree T yields up to three disjoint trees: the subtree rooted at $v.\text{left}$ (unless $v.\text{left}==\text{None}$), the subtree rooted at $v.\text{right}$ (unless $v.\text{right}==\text{None}$), and a tree rooted at $T.\text{root}$ consisting of all non-descendants of v (unless $T.\text{root}==v$). For example:

Before:



After:



Prove that in every tree T of size n , there exists a vertex v such that removing v from T results in disjoint trees that all have size at most $n/2$.

You may prove this however you like, but a recommended approach is to define a “potential function” ϕ on the vertices of the tree, by setting $\phi(v)$ to equal the size of the largest tree created by removing v . Let v^* be a vertex that minimizes the value of ϕ , i.e. v^* is a vertex such that $\phi(v^*) \leq \phi(v)$ for all other vertices v . Then we want to prove that $\phi(v^*) \leq n/2$. Prove this by contradiction. (Hint: try to show that either the parent

or one of the children will have smaller potential. If you're feeling stuck, try drawing some pictures!)

Proof:

1. Using the first suggestion provided, we create a potential function, where $\phi(\mathbf{v})$ represents the size of the largest tree created when the node \mathbf{v} is removed.

2. Then, we turn our claim into a logical statement using the potential function. Thus, we get:

$$\exists \mathbf{v}. \phi(\mathbf{v}) \leq \frac{n}{2} \quad (1)$$

3. Using proof by contradiction, we negate the statement to get:

$$\forall \mathbf{v}. \phi(\mathbf{v}) > \frac{n}{2} \quad (2)$$

4. Given inequality #2, we can conclude that:

$$\phi(\mathbf{v}^*) > \frac{n}{2} \quad (3)$$

Now, we'll be working under this new assumption, and trying to disprove it.

5. Given a BinaryTree T , let us find a vertex \mathbf{v}^* in T such that the inequality #3 is true. Having found \mathbf{v}^* , we know that \mathbf{v}^* is connected to the largest sub-tree whose size would be $> \frac{n}{2}$ if \mathbf{v}^* is removed.

6. Moving in the direction of that largest sub-tree, let the first node we encounter next (either the parent or a child of \mathbf{v}^*) be \mathbf{v}' . We need to prove that removing \mathbf{v}' will lead to a smaller potential function.

7. Upon removing \mathbf{v}' , we will encounter two cases: (1) a case where the removal of \mathbf{v}' will not affect the largest sub-tree (i.e. the previous largest sub-tree created when removing \mathbf{v}^* is still the largest sub-tree created when removing \mathbf{v}'), and (2) a case where the removal of \mathbf{v}' will heavily affect the largest sub-tree (i.e. the previous largest sub-tree created when removing \mathbf{v}^* is split into two smaller sub-trees).

7a. Let's address case 1: If we remove \mathbf{v}' and the largest sub-tree created remains, then its size is now $> \frac{n}{2} - 1$, which is less than $> \frac{n}{2}$ (the potential upon removing \mathbf{v}^*). This is a contradiction because now we have found a smaller largest sub-tree.

7b. Let's address case 2: If we remove \mathbf{v}' and the largest sub-tree created splits, then the largest sub-tree is now the sub-tree rooted at \mathbf{v}^* and its size should be $> \frac{n}{2}$. However, in the beginning of the proof we assumed that any other sub-tree rooted at \mathbf{v}^* that wasn't the largest sub-tree has the size of $< \frac{n}{2}$. But now, here, we're claiming that that same sub-tree is $> \frac{n}{2}$. That is a contradiction.

8. Therefore, both cases are accounted for and the original claim is proven true using contradiction.

Q.E.D.

- (c) (from proofs to algorithms) Turn your proof from Part 1b into a Python program that given a root vertex r of a *size-augmented* tree T with n vertices finds a vertex v with $\phi(v) \leq n/2$. Your program should run in time $O(h)$ on all size-augmented trees of height h ; again informally justify why your program has such a runtime. (Hint: try to repeatedly reduce the potential function by moving to children. Why don't we need to try moving to parents as in the previous proof?)

Answer: Reasoning provided within Python document.

2. (matchings and induction) Later in the course, we will study matching algorithms that are used in practice to match kidney donors to patients. The challenge in general is that some donors are incompatible with some patients (i.e. the patient's body is likely to reject the donated kidney). Suppose we are very lucky and have n donors and n patients where each donor d is incompatible with exactly one patient, denoted $incomp(d)$, and each patient p is incompatible with exactly one donor $incomp(p)$. (Incompatibility is symmetric, so $incomp(d) = p$ iff $incomp(p) = d$.) Let $f(n)$ be the number of ways, under these circumstances, of matching donors to patients so that each donor donates exactly one kidney to a compatible patient and each patient receives exactly one kidney from a compatible donor.

- (a) Show that $f(1) = 0$, $f(2) = 1$, and for all $n \geq 3$, we have

$$f(n) = (n - 1) \cdot (f(n - 1) + f(n - 2)).$$

(Hint: let d be one of the donors, and consider all possible patients p with whom d could be matched. Then consider cases according to whether $incomp(p)$ is matched with $incomp(d)$ or not.)

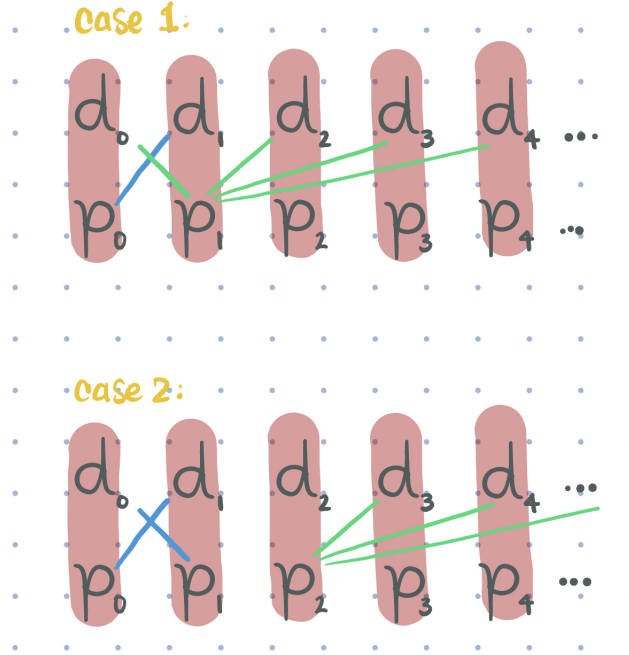
Proof:

1. Let us prove this claim using the story-proof method, and deconstruct the function $f(n)$ given.

2. Let us, for the purposes of making the proof more concrete, assume that $n = 5$, such that each of the 5 donors and 5 patients are placed in front of each other. Each donor-patient pair that is sitting right in front of each other is incompatible. That means that each of the 5 donors has 4 patients they are compatible with, and each of the 5 patients has 4 donors they are compatible with. This explains the first term of $f(n)$, $(n - 1)$.

3. Within the parenthesis, we have the second and third terms that we want to justify. Let us begin with the first term, $f(n - 1)$:

3a. (See attached image) Going back to our concrete example, let's assume that we match a patient, p_0 , to a valid donor, d_1 , such that $incomp(d_0) = incomp(p_1)$. Locking



this pair in place, this means that the next patient, p_1 , can match with any other $(n - 1)$ donors. Thus, together, there could be $f(n - 1)$ possible pairings.

3b. (See attached image) Now, let's consider another scenario, where on top of the pairing from 3a, we match another patient, p_1 , with another valid donor, d_0 , such that $incomp(d_1) = incomp(p_0)$. Having two locked pairs, the next patient, p_3 , can match with any other $(n - 2)$ donors. Thus, together, there could be $f(n - 2)$ possible pairings. We could continue this process of locking pairs and seeing the other possible pairings this creates, but this would not be necessary because we would be repeating the scenarios – either the scenario where the number of patients equals the number of donors, or the scenario where one number is greater than the other.

4. In this way, we've justified every term within the equation $f(n)$ and this claim is proved true.
Q.E.D.

(b) Prove by strong induction that for all $n \geq 2$, we have

$$\frac{n!}{3} \leq f(n) \leq \frac{n!}{2}.$$

Proof:

Base Cases:

$n = 3$:

$$\frac{3!}{3} \leq (2) \cdot (f(2) + f(1)) \leq \frac{3!}{2} \quad (4)$$

$$2 \leq 2 \leq 3 \quad (5)$$

$n = 2$:

$$\frac{2!}{3} \leq 1 \leq \frac{2!}{2} \quad (6)$$

$$\frac{2}{3} \leq 1 \leq 1 \quad (7)$$

Inductive Hypotheses:

For any k and $k - 1$, the following hold:

$$\frac{k!}{3} \leq f(k) \leq \frac{k!}{2} \quad (8)$$

$$\frac{k-1!}{3} \leq f(k-1) \leq \frac{k-1!}{2} \quad (9)$$

Inductive Step:

WTS that $\frac{(k+1)!}{3} \leq f(k+1) \leq \frac{(k+1)!}{2}$ is true for all k , knowing that $f(k+1) = k(f(k) + f(k-1))$.

Let us add our two hypotheses together:

$$\frac{k! + (k-1)!}{3} \leq f(k) + f(k-1) \leq \frac{k! + (k-1)!}{2} \quad (10)$$

Now, let us multiply everything by k and simplify:

$$\frac{k(k! + (k-1)!)}{3} \leq k(f(k) + f(k-1)) \leq \frac{k(k! + (k-1)!)}{2} \quad (11)$$

$$\frac{k * k! + k!}{3} \leq k(f(k) + f(k-1)) \leq \frac{k * k! + k!}{2} \quad (12)$$

$$\frac{k!(k+1)}{3} \leq k(f(k) + f(k-1)) \leq \frac{k!(k+1)}{2} \quad (13)$$

$$\frac{(k+1)!}{3} \leq f(k+1) \leq \frac{(k+1)!}{2} \quad (14)$$

By this step, we are done!

Q.E.D.