

Practice Problems for the Final Exam

Our final exam will be held on **Fri, May 13, 4–6pm** in **Tydings Hall (TYD) 0130**. It is close-book, closed-notes, but you will be allowed three sheets of notes, front and back.

Disclaimer: The exam will be comprehensive, emphasizing material in the latter half of the semester. **These practice problems reflect the just material since the second midterm, but you should expect coverage of other topics as well.** They have been extracted from old homework assignments and exams. Material changes from semester to semester. These do **not** necessarily reflect the actual coverage, difficulty, or length of the midterm exam.

Problem 0. Since the exam is comprehensive, please look back over the previous homework assignments, the two midterm exams, and the practice problems for both midterms. You should expect at least one problem that involves tracing through an algorithm or construction given in class.

Problem 1. Short answer questions. Except where noted, explanations are not required but may be given for partial credit.

- (a) Let T be extended binary search tree (that is, one having internal and external nodes). You visit the nodes of T according to one of the standard traversals (preorder, postorder, or inorder). Which of the following statements is necessarily true? (Select all that apply.)
 - (i) In a *postorder traversal*, all the external nodes appear in the order *before* any of the internal nodes
 - (ii) In a *preorder traversal*, all the internal nodes appear in the order *after* any of the external nodes
 - (iii) In an *inorder traversal*, internal and external node *alternate* with each other
 - (iv) None of the above is true
- (b) You have an AVL tree containing n keys, and you insert a new key. As a function of n , what is the maximum number of rotations that might be needed as part of this operation? (A double rotation is counted as two rotations.) Explain briefly.
- (c) Repeat (b) in the case of deletion. (Give your answer as an asymptotic function of n .)
- (d) Splay trees are known to support efficient finger search queries. What is a “finger search query”?
- (e) In class, we mentioned that when using double hashing, it is important that the second hash function $g(x)$ should not share any common divisors with the table size m . What might go wrong if this were not the case?
- (f) Consider the following dictionary structures: (1) Unbalanced binary search trees, (2) AVL trees, (3) AA-trees, (4) quake-heaps, (5) treaps, (6) splay trees, (7) scapegoat trees. Suppose we insert one key into such a data structure that contains n keys. For which of these data structures can we assert that the worst-case number of structural changes to the tree is $O(\log n)$? (A *structure change* is any local alteration of the structure: creating/modifying node contents, rotation, node split, etc.)

- (g) A scapegoat tree containing n keys has height $O(\log n)$... (select one):
- (i) Always—the height is guaranteed
 - (ii) In expectation, over the algorithm’s random choices
 - (iii) In expectation, assuming that keys are inserted in random order
 - (iv) In the amortized sense—the average height will be $O(\log n)$ over a long sequence of operations
 - (v) Maybe yes, maybe no—there is just no way of knowing
- (h) Between the classical dynamic storage allocation algorithm (with arbitrary-sized blocks) or the buddy system (with blocks of size power of 2) which is more susceptible to *internal fragmentation*? Explain briefly.
- (i) (Check out Problem 1(g) from the Practice Problems for Midterm 2 on hashing.)

Problem 2. This problem involves an input which is a binary search tree having n nodes of height $O(\log n)$. You may assume that each node p has a field $p.size$ that stores the number of nodes in its subtree (including p itself). Here is the node structure:

```
class Node {
    int key;
    Node left;
    Node right;
    int size; // number of nodes in this subtree
}
```

- (a) Present pseudocode for a function `printMaxK(int k)`, which is given $0 \leq k \leq n$, and prints the values of the k largest keys in the binary search tree. (See, for example, Fig. 1.)

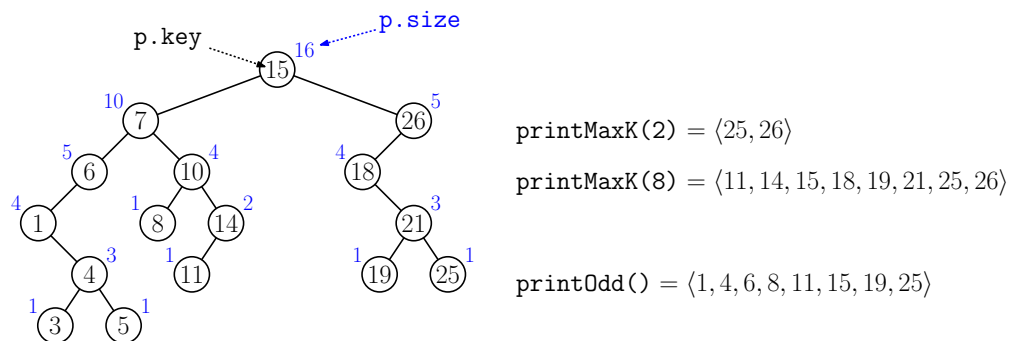


Figure 1: The functions `printMaxK` and `printOdd`.

You should do this by traversing the tree. You are not allowed to “cheat” but storing an auxiliary list of sorted nodes.

For fullest credit, the keys should be printed in *ascending order*, and your algorithm should run in time $O(k + \log n)$ (see part (4.2) below). (Partial credit will be given otherwise, but an $O(n)$ time algorithm is not worth anything.)

You may assume that $0 \leq k \leq n$, where n is the total number of nodes in the tree. Briefly explain your algorithm.

Hint: I would suggest using the helper function `printMaxK(Node p, int k)`, where `k` is the number of keys to print from the subtree rooted at `p`.

- (b) Derive the running time of your algorithm in (a).
- (c) Give pseudocode for a function `printOdd()`, which does the following. Let $\langle x_1, x_2, \dots, x_n \rangle$ denote the keys of the tree in ascending order, this function prints every other key, namely $\langle x_1, x_3, x_5, \dots, x_n \rangle$, if n is odd, and $\langle x_1, x_3, x_5, \dots, x_{n-1} \rangle$, if n is even.

Beware: We are not printing the “odd-valued” keys, rather we are printing the odd numbered positions in the sorted order (see Fig. 1.)

Again, you should do this by traversing the tree. You are not allowed to “cheat” by storing auxiliary lists or using global variables. Your program should run in time $O(n)$. Briefly explain your algorithm.

Problem 3. Throughout this problem, assume that you are given a standard kd-tree storing a set P of n points in \mathbb{R}^2 (see Fig. 2(a)). Assume that the cutting dimension alternates between x and y . You may also assume that the tree stores a bounding box `bbox`, which is a 2-dimensional rectangle containing all the points of P . You may also assume that any geometric computations on primitive objects (distances, disjointness, containment, etc.) can be computed in constant time, without explanation.

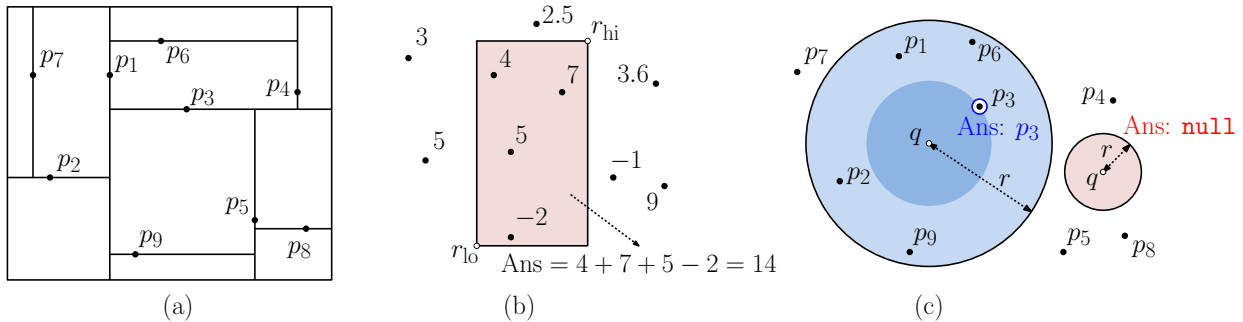


Figure 2: Queries on kd-trees.

- (a) In a standard range-counting query, we want to count the number of points in the range. Suppose that each point $p_i \in P$ has an associated real-valued weight w_i . In a *weighted orthogonal range query*, we are given a query rectangle R , given by its lower-left corner r_{lo} and upper-right corner r_{hi} , and the answer is the sum of the weights of the points that lie within R (see Fig. 2(b)). If there are no points in the range, the answer is 0. Explain how to modify the kd-tree (by adding additional fields to the nodes, if you like) so that weighted orthogonal range queries can be answered efficiently. Based on your modified data structure, present an efficient algorithm in pseudo-code for answering these queries and explain. (For full credit, your algorithm should run in $O(\sqrt{n})$ time). You may handle the edge cases (e.g., points lying on the boundary of R) however you like. **Hint:** You may use whatever helper function(s) you like, but I would suggest using:

`double weightedRange(Rectangle R, KDNODE p, Rectangle cell)`

where `p` is the current node in the kd-tree, `cell` is the associated cell.

- (b) Briefly analyze the running time of your algorithm, assuming that the tree is balanced. (You may apply/modify results proved in class.)
- (c) In a *fixed-radius nearest neighbor query*, we are given a point $q \in \mathbb{R}^d$ and a radius $r > 0$. Let C denote the circular disk centered at q whose radius is r . If no points of P lie within this disk, the answer to the query is `null`. Otherwise, it returns the point of P within the disk that is closest to q . Present (in pseudo-code) an efficient kd-tree algorithm that answers such a query.

You may handle the edge cases (e.g., multiple points at the same distance or points lying on the boundary of C) however you like. **Hint:** You may use whatever helper function(s) you like, but I would suggest using:

`Point frnn(Point q, double r, KDNODE p, Rectangle cell, Point best)`

where `p` is the current node in the kd-tree, `cell` is the associated cell, and `best` is the best point seen so far.

Briefly explain your algorithm, but you *do not* need to derive its running time.

Problem 4. In this problem we will build a suffix tree for the string $S = \text{baabaabababaa}\$$.

- (a) List the substring identifiers for the 14 suffixes of S . For the sake of uniformity, list them in order (either back to front or front to back). For example, you could start with “\$” and end with the substring identifier for the entire string.
- (b) List the substring identifiers again together with their indices (0 through 13), but this time in alphabetical order (where “a” < “b” < “\$”).
- (c) Draw a picture of the suffix tree for S . For the sake of uniformity, when drawing your tree, use the convention of Fig. 7 in the Lecture 17 LaTeX lecture notes. In particular, label edges of the final tree with substrings, index the suffixes from 0 to 13, and order subtrees in ascending lexicographical order.

Problem 5. In this problem, we will consider how to use/modify range trees to answer two related queries. While the answer should be based on range trees, you may need to make modifications including possibly transforming the points and even adding additional coordinates. In each case, describe the points that are stored in the range tree and how the search process works. An English explanation (as opposed to pseudocode) is sufficient. Justify your algorithm’s correctness and derive its running time.

- (a) Assume you are given an n -element point set P in \mathbb{R}^2 (see Fig. 3(a)). In addition to its coordinates (p_x, p_y) , each point $p \in P$ is associated with a numeric *rating*, p_z . In an *orthogonal top- k query*, you are given an axis-aligned query rectangle R (given, say, by its lower-left and upper-right corners) and a positive integer k . The query returns a list of the (up to) k points of P that lie within R having the highest ratings (see Fig. 3(b)). (As an application, imagine you are searching for the k highest rated restaurants in a rectangular region of some city.)

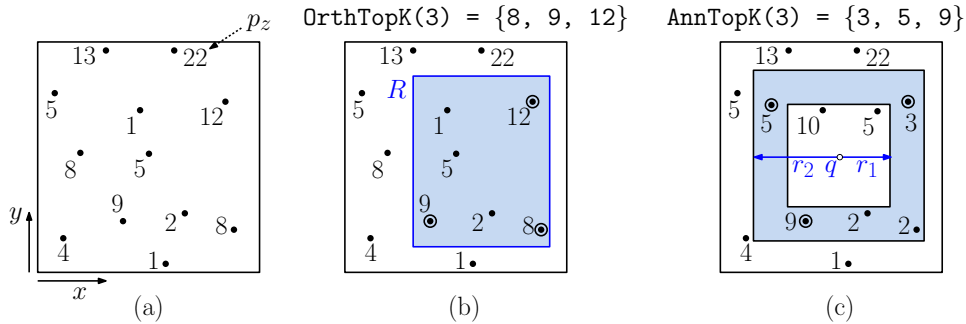


Figure 3: Orthogonal top- k queries and annulus top- k queries.

Describe how to preprocess the point set P into a data structure that can efficiently answer any orthogonal top- k query (R, k) . Your data structure should use $O(n \log^2 n)$ storage and answer queries in at most $O(k \log^2 n)$ time. (I don't care how you handle edge cases, such as points lying on the boundary of the rectangle or points having the same rating.) If there are k points or fewer in the query region, the list will contain them all.

- (b) In an *annulus top- k query* a query is given by a query point $q \in \mathbb{R}^2$ and two positive radii $r_1 < r_2$. Let $S_1 = S(q, r_1)$ be the square centered at q whose half side length is r_1 and define S_2 similarly for q and r_2 . The square annulus $A(q, r_1, r_2)$ is defined to be the region between these two squares. The query returns a list of the (up to) k points of P that lie within the annulus $A(q, r_1, r_2)$ that have the highest ratings (see Fig. 3(c)).

Problem 6. Suppose you have a large span of memory, which starts at some address `start` and ends at address `end-1` (see Fig. 4). (The variables `start` and `end` are generic pointers of type `void*`.) As the dynamic memory allocation method of Lecture 15, this span is subdivided into blocks. The block starting at address `p` is associated with the following information:

- `p.inUse` is 1 if this block is in-use (allocated) and 0 otherwise (available)
- `p.prevInUse` is 1 if the block immediately preceding this block in memory is in-use. (It should be 1 for the first block.)
- `p.size` is the number of words in this block (including all header fields)
- `p.size2` each available block has a copy of the size stored in its last word, which is located at address `p + p.size - 1`.

(For this problem, we will ignore the available-list pointers `p.prev` and `p.next`.)

In class, we said that in real memory-allocation systems, blocks cannot be moved, because they may contain pointers. Suppose, however, that the blocks are movable. Present pseudo-code for a function that compacts memory by copying all the allocated blocks to a single contiguous span of blocks at the start of the memory span (see Fig. 4). Your function `compress(void* start, void* end)` should return a pointer to the head of the available block at the end. Following these blocks is a single available block that covers the rest of the memory's span.

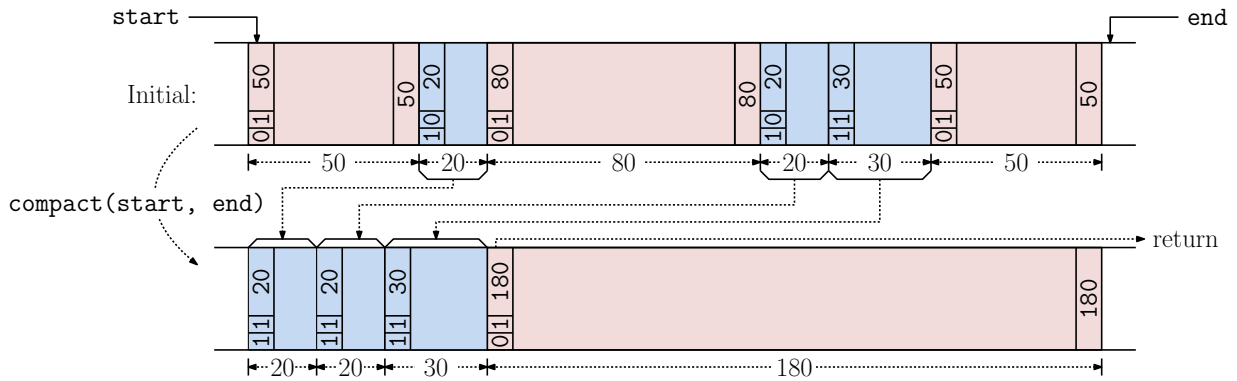


Figure 4: Memory compactor.

To help copy blocks of memory around, you may assume that you have access to a function `void* memcpy(void* dest, void* source, int num)`, which copies `num` words of memory from the address `source` to the address `dest`.

Problem 7. Recall the buddy system of allocating blocks of memory (see Fig. 5). Throughout this problem you may use the following standard bit-wise operators:

<code>&</code>	bit-wise “and”	<code> </code>	bit-wise “or”
<code>^</code>	bit-wise “exclusive-or”	<code>~</code>	bit-wise “complement”
<code><<</code>	left shift (filling with zeros)	<code>>></code>	right shift (filling with zeros)

You may also assume that you have access to a function `bitMask(k)`, which returns a binary number whose k lowest-order bits are all 1’s. For example `bitMask(3) = 1112 = 7`.

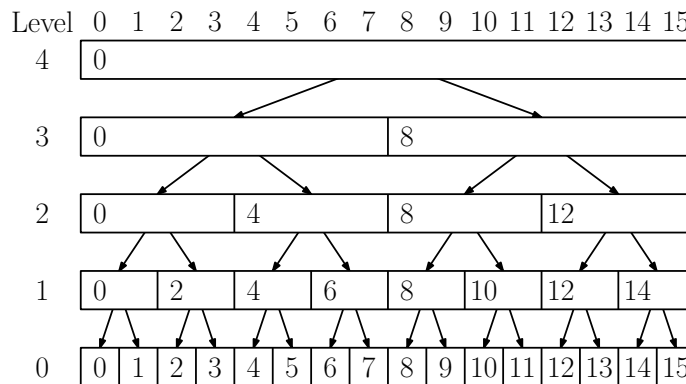


Figure 5: Buddy relatives.

Present a short (one-line) expression for each of the following functions in terms of the above bit-wise functions:

- (a) `boolean isValid(int k, int x)`: True if and only if $x \geq 0$ a valid starting address for a buddy block at level $k \geq 0$.

- (b) `int sibling(int k, int x)`: Given a valid buddy block of level $k \geq 0$ starting at address x , returns the starting address of its *sibling* (that is, its “buddy”).
- (c) `int parent(int k, int x)`: Given a valid buddy block of level $k \geq 0$ starting at address x , returns the starting address of its *parent* at level $k + 1$.
- (d) `int left(int k, int x)`: Given a valid buddy block of level $k \geq 1$ starting at address x , returns the starting address of its *left child* at level $k - 1$.
- (e) `int right(int k, int x)`: Given a valid buddy block of level $k \geq 1$ starting at address x , returns the starting address of its *right child* at level $k - 1$.

For example, given the tree shown in the figure, we have

```

isValid(2, 12) = isValid(2, 01100) = True
isValid(2, 10) = isValid(2, 01010) = False
sibling(2, 12) = sibling(2, 01100) = 8 = 01000
parent(2, 12) = parent(2, 01100) = 8 = 01000
  left(2, 12) = left(2, 01100) = 12 = 01100
  right(2, 12) = right(2, 01100) = 14 = 01110

```

Problem 8. This problem involves a data structure called an *erasable stack*. This data structure is just a stack with an additional operation that allows us to “erase” any element that is currently in the stack. Whenever we pop the stack, we skip over the erased elements, returning the topmost “unerased” element. The pseudocode below provides more details to be implemented.

```

class EStack {      // erasable stack of Objects
    int top          // index of stack top
    Object A[HUGE]  // array is so big, we will never overflow
    Object ERASED   // special object which indicates an element is erased

    EStack() { top = -1 } // initialize

    void push(Object x) { // push
        A[++top] = x
    }

    void erase(int i) {   // erase (assume 0 <= i <= top)
        A[i] = ERASED
    }

    Object pop() {       // pop (skipping erased items)
        while (top >= 0 && A[top] == ERASED) top--
        if (top >= 0) return A[top--]
        else return null
    }
}

```

Let $n = \text{top} + 1$ denote the current number of entries in the stack (including the ERASED entries). Define the *actual cost* of operations as follows: `push` and `erase` both run in 1 unit of time and `pop` takes $k + 1$ units of time where k is the number of ERASED elements that were skipped over.

- (a) As a function of n , what is the *worst-case running time* of the `pop` operation? (For fullest credit, make your bound as tight as possible.) Justify your answer.
- (b) Starting with an empty stack, we perform a sequence of m `push`, `erase`, and `pop` operations. Give an upper bound on the *amortized running time* of such a sequence. You may assume that all the operations are valid and the array never overflows. (For fullest credit, make your bound as tight as possible.) Justify your answer.
- (c) Given two (large) integers k and m , where $k \leq m/2$, we start from an empty stack, push m elements, and then erase k elements *at random*, finally we perform a single `pop` operation. What is the *expected running time* of the final `pop` operation. You may express your answer asymptotically as a function of k and m .

In each case, state your answer first, and then provide your justification.

Problem 9. (Check out Problem 11 from the Practice Problems for Midterm 2 on deletion in open-addressing hashing.)