

Think Algorithms

Algorithms and Data Structures in Java

Version 0.1.0

Think Algorithms

Algorithms and Data Structures in Java

Version 0.1.0

Allen B. Downey

Green Tea Press

Needham, Massachusetts

Copyright © 2016 Allen B. Downey.

Green Tea Press
9 Washburn Ave
Needham, MA 02492

Permission is granted to copy, distribute, and/or modify this work under the terms of the Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License, which is available at <http://creativecommons.org/licenses/by-nc-sa/3.0/>.

The original form of this book is L^AT_EX source code. Compiling this code has the effect of generating a device-independent representation of the book, which can be converted to other formats and printed.

The L^AT_EX source for this book is available from [???](#).

DRAFT: Not ready for distribution!

Contents

Preface	xiii
1 Interfaces	1
1.1 Prerequisites	2
1.2 Why are there two kinds of List?	3
1.3 interfaces in Java	3
1.4 The List interface	4
1.5 Exercise 1	6
2 Analysis of algorithms	9
2.1 What You Should Know Before Reading This	9
2.2 Overview	9
2.3 Objectives	9
2.4 Algorithm analysis	10
2.5 Big O notation	13
2.6 Resources	15
3 ArrayList	17
3.1 Objectives	17
3.2 Implementing an Array-backed List	17

3.3	Instructions	17
3.4	MyArrayList code	18
3.5	Instructions	20
4	Analyzing ArrayList	23
4.1	Overview	23
4.2	Objectives	23
4.3	Classifying MyArrayList methods	23
4.4	Classifying add	25
4.5	Classifying removeAll	28
4.6	Resources	29
5	LinkedList	31
5.1	Objectives	31
5.2	Overview	31
5.3	Linked data structures	31
5.4	Getting set up	33
5.5	MyLinkedList code	34
5.6	Adding nodes	35
5.7	Instructions	36
6	A note on garbage collection	39
6.1	Resources	40
7	Analysis of LinkedList	41
7.1	Overview	41
7.2	Objectives	41
7.3	Classifying MyLinkedList methods	41

7.4	Comparing <code>MyLinkedList</code> and <code>MyArrayList</code>	44
8	Performance profiling	47
8.1	Objectives	47
8.2	Overview	47
8.3	Profiler	47
8.4	Interpreting results	50
8.5	Instructions	51
8.6	Resources	52
9	Doubly-linked list	53
9.1	Overview	53
9.2	Objectives	53
9.3	Performance profiling results	53
9.4	Profiling <code>LinkedList</code> methods	55
9.5	Adding to the end of a <code>LinkedList</code>	56
9.6	Doubly-linked list	58
9.7	Resources	60
10	Tree traversal	61
10.1	Learning goals	61
10.2	Overview	61
10.3	The road ahead	61
10.4	Parsing HTML	63
10.5	Using jsoup	64
10.6	Iterating through the DOM	65
10.7	Depth-first search	66

10.8	Stacks in Java	67
10.9	Iterative DFS	69
10.10	Resources	70
11	Getting to Philosophy	71
11.1	Objectives	71
11.2	Overview	71
11.3	Getting started	71
11.4	Iterables and Iterators	72
11.5	WikiFetcher	74
11.6	Filling in WikiPhilosophy	76
11.7	Resources	78
12	Indexer	79
12.1	Learning goals	79
12.2	Overview	79
12.3	Data structure selection	80
12.4	TermCounter	81
12.5	Finishing off TermCounter	84
12.6	Finishing off Index	85
13	The Map interface	89
13.1	Learning goals	89
13.2	Overview	89
13.3	Implementing MyLinearMap	89
13.4	Instructions	91
14	Hashing	93

14.1	Learning goals	93
14.2	Overview	93
14.3	Analyzing <code>MyLinearMap</code>	93
14.4	Hashing	96
14.5	How does hashing work?	98
14.6	Hashing and mutation	100
14.7	Resources	102
15	Hash table implementation	103
15.1	Learning goals	103
15.2	Overview	103
15.3	Finishing <code>MyBetterMap</code>	103
15.4	Implementing <code>MyHashMap</code>	105
15.5	Analyzing <code>MyHashMap</code>	106
15.6	The tradeoffs	107
15.7	Profiling <code>MyHashMap</code>	108
15.8	Resources	109
16	Fixing the hash table	111
16.1	Learning goals	111
16.2	Overview	111
16.3	Fixing <code>MyHashMap</code>	111
16.4	UML class diagrams	114
16.5	Resources	115
17	TreeMap	117
17.1	Learning goals	117

17.2	Overview	117
17.3	Binary search tree	118
17.4	Implementing a tree-backed Map	120
17.5	Instructions	121
17.6	Resources	124
18	Binary search tree	125
18.1	Learning goals	125
18.2	Overview	125
18.3	Our version of MyTreeMap	125
18.4	Searching for values	127
18.5	Implementing put	128
18.6	In-order traversal	130
18.7	The logarithmic methods	131
18.8	Self-balancing trees	134
18.9	One more exercise	134
18.10	Resources	135
19	Redis	137
19.1	Learning goals	137
19.2	Overview	137
19.3	Persistence	138
19.4	Redis clients and servers	139
19.5	Hello, Jedis	140
19.6	Redis data types	143
19.7	Making a Jedis-backed index	145
19.8	More suggestions if you want them	147

19.9	A few design hints	148
19.10	Resources	149
20	Crawling Wikipedia	151
20.1	Learning goals	151
20.2	Overview	151
20.3	Our Redis-backed indexer	151
20.4	Analysis of lookup	154
21	Analysis of indexing	157
21.1	Graph traversal	158
21.2	Making a crawler	159
21.3	Resources	162
22	Boolean search	163
22.1	Learning goals	163
22.2	Overview	163
22.3	Crawler solution	163
22.4	Information retrieval	166
22.5	Boolean search	167
22.6	Relevance scores	168
22.7	Implementing boolean operators	169
22.8	Comparable and Comparator	171
22.9	Extensions	173
22.10	Resources	174
23	Sorting	175
23.1	Learning goals	175

23.2	Overview	175
23.3	Insertion sort	176
23.4	Merge sort	178
23.5	Analysis of merge sort	180
23.6	Radix sort	181
23.7	Heap sort	183
23.8	Bounded heap	184
23.9	Space complexity	185
23.10	Resources	186
Index		189

Preface

The philosophy behind the book

Here are the guiding principles that make the book the way it is:

- *One concept at a time.* We break down topics that give beginners trouble into a series of small steps, so that they can exercise each new concept in isolation before continuing.
- *Balance of Java and concepts.* The book is not primarily about Java; it uses code examples to demonstrate computer science. Most chapters start with language features and end with concepts.
- *Conciseness.* An important goal of the book is to be small enough so that students can read and understand the entire text in a one-semester college or AP course.
- *Emphasis on vocabulary.* We try to introduce the minimum number of terms and define them carefully when they are first used. We also organize them in glossaries at the end of each chapter.
- *Program development.* There are many strategies for writing programs, including bottom-up, top-down, and others. We demonstrate multiple program development techniques, allowing readers to choose methods that work best for them.
- *Multiple learning curves.* To write a program, you have to understand the algorithm, know the programming language, and be able to debug errors. We discuss these and other aspects throughout the book, and include an appendix that summarizes our advice.

Working with the code

Most of the code examples in this book are available from a Git repository at <https://github.com/AllenDowney/ThinkJavaCode>. Git is a “version control system” that allows you to keep track of the files that make up a project. A collection of files under Git’s control is called a “repository”.

GitHub is a hosting service that provides storage for Git repositories and a convenient web interface. It provides several ways to work with the code:

- You can create a copy of the repository on GitHub by pressing the **Fork** button. If you don’t already have a GitHub account, you’ll need to create one. After forking, you’ll have your own repository on GitHub that you can use to keep track of code you write. Then you can “clone” the repository, which downloads a copy of the files to your computer.
- Alternatively, you could clone the repository without forking. If you choose this option, you don’t need a GitHub account, but you won’t be able to save your changes on GitHub.
- If you don’t want to use Git at all, you can download the code in a ZIP archive using the **Download ZIP** button on the GitHub page, or this link: <http://tinyurl.com/ThinkJavaCodeZip>.

After you clone the repository or unzip the ZIP file, you should have a directory called **ThinkJavaCode** with a subdirectory for each chapter in the book.

All examples in this book were developed and tested using Java SE Development Kit 8. If you are using a more recent version, the examples in this book should still work. If you are using an older version, some of them may not.

Contributors

- Flatiron?

Additional contributors who found one or more typos:

If you have additional comments or ideas about the text, please send them to: feedback@greenteapress.com.

Chapter 1

Interfaces

This book presents three topics:

- Data structures: Starting with the structures in the Java Collections Framework (JCF), you will learn how to use data structures like lists and maps, and you will see how they work.
- Analysis of algorithms: I will present techniques for analyzing code and predicting how fast it will run and how much space (memory) it will require.
- Information retrieval: To motivate the first two topics, and to make the exercises more interesting, we will use data structures and algorithms to build a simple Web search engine.

Here's an outline of the order of topics:

- We'll start with the `List` interface and you will write classes that implement this interface two different ways. Then we'll compare your implementations with the Java classes `ArrayList` and `LinkedList`.
- Next I'll introduce tree-shaped data structures and you will work on the first application, a program that reads pages from Wikipedia, parses the contents, and navigates the resulting tree to find links and other features. We'll use these tools to test the "Getting to Philosophy" conjecture (you can get a preview by reading https://en.wikipedia.org/wiki/Wikipedia:Getting_to_Philosophy).

- We'll learn about the `Map` interface using Java's `HashMap` implementation. Then you'll write classes that implement this interface using a hash table and a binary search tree.
- Finally, you will use these classes (and a few others I'll present along the way) to implement the pieces of a Web search engine: a `Crawler` that finds and reads pages, an `Indexer` that stores the contents of Web pages in a form that can be searched efficiently, and a `Retriever` that takes queries from a user and returns relevant results.

But first, a few prerequisites.

1.1 Prerequisites

Before you start this book, you should know Java pretty well; in particular, you should know how to define a new class that extends an existing class or implements an `interface`. If your Java is rusty, here are two books you might start with:

- Downey and Mayfield, *Think Java* (O'Reilly Media 2016), which is intended for people who have never programmed before.
- Sierra and Bates, *Head First Java* (O'Reilly Media, 2005), which might be more appropriate for people who know another programming language.

If you are not familiar with Interfaces in Java, you might want to work through the tutorial called “What Is an Interface?” at <https://docs.oracle.com/javase/tutorial/java/concepts/interface.html>.

One vocabulary note: the word “interface” can be confusing. In the context of an **application programming interface** (API), it refers to a set of classes and methods that provide certain capabilities.

In the context of Java, it also refers to a language feature, similar to a class, that specifies a set of methods. To help avoid confusion, we'll use “interface” in the normal typeface for the general idea of an interface, and `interface` in the code typeface for the Java language feature.

You should also be familiar with type parameters and generic types. If not, you can read about them at <https://docs.oracle.com/javase/tutorial/java/generics/types.html>.

Finally, I will assume you are familiar with the Java Collections Framework (JCF), which you can read about at <https://docs.oracle.com/javase/7/docs/api/java/util/Collections.html>. In particular, you should know about the `List` **interface** and the classes `ArrayList` and `LinkedList`.

1.2 Why are there two kinds of List?

When people start working with the JCF, they are sometimes confused about `ArrayList` and `LinkedList`. Why does Java provide two implementations of the `List` **interface**? And how should we choose which one to use? We will answer these questions in the next few chapters.

We'll start by reviewing **interfaces** and classes that implement them, and we'll present the idea of “programming to an interface”. In the first two exercises, you'll implement classes similar to `ArrayList` and `LinkedList`, so you'll know how they work, and we'll see that each of them has pros and cons. Some operations are faster or use less space with `ArrayList`; others are faster or smaller with `LinkedList`. Which one is better for a particular application depends on which operations it performs most often.

1.3 **interfaces** in Java

A Java **interface** specifies a set of methods; any class that implements this **interface** has to provide these methods. For example, here is the source code for `Comparable`, which is an **interface** defined in the package `java.lang`:

```
public interface Comparable<T> {  
    public int compareTo(T o);  
}
```

This **interface** definition uses a type parameter, `T`, which makes `Comparable` a **generic type**.

In order to implement this **interface**, a class has to

- Specify the type `T` refers to, and
- Provide a method named `compareTo` that takes an object as a parameter and returns an `int`.

For example, here's the source code for `java.lang.Integer`:

```
public final class Integer extends Number implements Comparable<Integer> {  
  
    public int compareTo(Integer anotherInteger) {  
        int thisVal = this.value;  
        int anotherVal = anotherInteger.value;  
        return (thisVal < anotherVal ? -1 : (thisVal == anotherVal ? 0 : 1));  
    }  
  
    // other methods omitted  
}
```

This class extends `Number`, so it inherits the methods and instance variables of `Number`; and it implements `Comparable<Integer>`, so it provides a method named `compareTo` that takes an `Integer` and returns an 'int'.

When a class declares that it implements an `interface`, the compiler checks that it provides implementations of all methods defined by the `interface`.

1.4 The List `interface`

The Java collections framework (JCF) defines an interface called `List` and provides two implementations, `ArrayList` and `LinkedList`.

The interface defines what it means to be a `List`; any class that implements this `interface` has to provide a particular set of methods, including `add`, `get`, `remove`, and about 20 more.

`ArrayList` and `LinkedList` provide these methods, so they can be used interchangeably. A method written to work with a `List` will work with an `ArrayList`, `LinkedList`, or any other object that implements `List`.

Here's a contrived example that demonstrates this point:

```
public class ListClientExample {
    private List list;

    public ListClientExample() {
        list = new LinkedList();
    }

    private List getList() {
        return list;
    }

    public static void main(String[] args) {
        ListClientExample lce = new ListClientExample();
        List list = lce.getList();
        System.out.println(list);
    }
}
```

`ListClientExample` doesn't do anything useful, but it has the essential elements of a class that **encapsulates** a `List`; that is, it contains a `List` as an instance variable. We'll use this class to make a point, and then you'll work with it in the first exercise.

The `ListClientExample` constructor initializes `list` by **instantiating** (that is, creating) a new `LinkedList`; the getter method called `getList` returns the internal `List` object; and `main` contains a few lines of code to test these methods.

The important thing about this example is that it uses `List` whenever possible and avoids specifying `LinkedList` or `ArrayList` unless it is necessary. For example, the instance variable is declared to be a `List`, and `getList` returns a `List`, but neither specifies which kind of list.

If you change your mind and decide to use an `ArrayList`, you only have to change the constructor; you don't have to make any other changes. You'll see how this works in the first exercise.

This style is called **interface-based programming**, or more casually, "programming to an interface" (see https://en.wikipedia.org/wiki/interface-based_programming).

(Notice that now we are talking about the general idea of an interface, not a Java `interface`.)

When you use a library, your code should only depend on the interface the library provides; it should not depend on details of the implementation. That way, if the implementation of the library changes in the future, your code will still work.

On the other hand, if the interface changes, the code that depends on it has to change, too. That's why library developers avoid changing interfaces unless absolutely necessary.

1.5 Exercise 1

Since this is the first exercise, we'll keep it simple. You will take the code from the previous section and **swap the implementation**; that is, you will replace the `LinkedList` with an `ArrayList`. Because the code programs to an interface, you will be able to swap the implementation by changing a single line (and adding an `import` statement).

- Set up your development environment. For this Track you will need to be able to compile and run Java code. We developed the examples using the Java SE Development Kit 7. If you are using a more recent version, everything should still work. If you are using an older version, you might find some incompatibilities.

We recommend using an IDE like Eclipse that provides syntax-checking, auto-completion, and source code refactoring. These features help you avoid errors or find them quickly. However, if you are preparing for a technical interview, remember that you will not have these tools during the interview, so you might also want to practice writing code without them.

Either way, we'll assume that you know how to compile and run Java code in your environment of choice.

- When you check out the repository for this exercise, you should find these files:

- CONTRIBUTING.md contains information about how you can notify us if you find a problem in a lesson.
 - javacs-lab01 is a directory that contains the source code for this exercise.
 - LICENSE.md contains license information about these materials.
 - README.md contains the text you are reading now.
- If you open `javacs-lab01`, you'll find these files:
 - `bin` contains compiled `class` files.
 - `build.xml` is an Ant file that makes it easier to compile and run the code.
 - `lib` contains the libraries you'll need (for this exercise, just JUnit).
 - `src` contains the source code.
 - And if you navigate into `src/com/flatironschool/javacs`, you'll find the source code for this exercise:
 - `ListClientExample.java` contains the code from the previous README.
 - `ListClientExampleTest.java` contains a JUnit test for `ListClientExample`.
 - Review `ListClientExample` and make sure you understand what it does. Then compile and run it. If you use Ant, you can navigate to `javacs-lab01` and run `ant ListClientExample`. <https://ant.apache.org/manual/tutorial-HelloWorldWithAnt.html> You can read about Ant here.

NOTE: You might get a warning like `List` is a raw type. References to generic type `List` should be in the form of `List<Integer>`. To keep the example simple, we didn't bother to specify the type of the elements in the `List`. If this warning bothers you, you can suppress it by replacing each `List` or `LinkedList` with `List<Integer>` or `LinkedList<Integer>`.

- Review `ListClientExampleTest`. It runs one test, which creates a `ListClientExample` invokes `getList`, and then checks whether the result is an `ArrayList`. Initially, this test will fail because the result is a `LinkedList`, not an `ArrayList`. Run this test and confirm that it fails.

NOTE: This test makes sense for this exercise, but it is not a good example of a test. Good tests should check whether the class under test

satisfies the requirements of the *interface*; they should not depend on the details of the *implementation*.

- In the `ListClientExample`, replace `LinkedList` with `ArrayList`. You might have to add an import statement. Compile and run `ListClientExample`. Then run the test again. With this change, the test should now pass.
- To make this test pass, you only had to replace `LinkedList` with `ArrayList` in the constructor. You did not have to change any of the places where `List` appears. What happens if you do? Go ahead and replace one or more appearances of `List` with `ArrayList`. The program should still work correctly, but now it is “overspecified”. If you change your mind in the future and want to swap the interface again, you would have to change more code.
- In the `ListClientExample` constructor, what happens if you replace `ArrayList` with `List`? Why can’t you instantiate a `List`?

DRAFT: Not ready for distribution

Chapter 2

Analysis of algorithms

2.1 What You Should Know Before Reading This

- <http://www.codejava.net/java-core/collections/java-list-collection-tutorial>
List interface
- <http://www.go4expert.com/articles/selection-sort-algorithm-absolute-t2788>
Sort
- http://www.tutorialspoint.com/java/java_arraylist_class.htmJava
ArrayList Usage

2.2 Overview

This README reviews the basic ideas of algorithm analysis. After this lesson, you should be able to compare algorithms written in Java and predict which will be more efficient for large problems.

2.3 Objectives

1. Describe the goals of algorithm analysis.

2. Define constant time, linear time, and quadratic time algorithms.
3. Analyze a Java function and classify its run time.
4. Read and write Big O notation.
5. Analyze the selection sort algorithm.

2.4 Algorithm analysis

If you have worked with the Java Collections Framework (JCF), you might have wondered why Java provides two implementations of the `List` interface. The short answer is that for some methods `LinkedList` is faster, and for other methods `ArrayList` is faster. Which one is better depends on how you use it.

To decide which one is better for a particular application, one approach is to try them both and see how long they take. This approach, which is called “profiling” has a few problems:

1. Before you compare the algorithms, you have to implement them both.
2. The results might depend on what kind of computer you use. One algorithm might be better on one machine; the other might be better on a different machine.
3. The results might depend on the size of the problem or the data provided as input.

We can address some of these problems using http://en.wikipedia.org/wiki/Analysis_of_algorithms analysis of algorithms. When it works, we can use algorithm analysis to compare algorithms without having to implement them. But we have to make some assumptions:

1. To avoid dealing with the details of computer hardware, we usually identify the basic operations that make up an algorithm — like addition, multiplication, and comparison of numbers — and count the number of operations each algorithm requires.
2. To avoid dealing with the details of the input data, we try to analyze the average performance for the inputs we expect. If that’s not possible, a common alternative is to analyze the worst case scenario, invoking the principle that we should hope for the best, but prepare for the worst.

3. Finally, we have to deal with the possibility that one algorithm works best for small problems and another for big ones. In that case, we usually focus on the big ones, because for big problems a good algorithm is often much faster than a bad one.

This kind of analysis lends itself to simple classification of algorithms. For example, if we know that the run time of Algorithm A tends to be directly proportional to the size of the input, n , and Algorithm B tends to be directly proportional to n^2 , we expect A to be faster than B, at least for large values of n .

Most simple algorithms fall into just a few categories.

- Constant time: An algorithm is “constant time” if the run time does not depend on the size of the input. For example, if you have an array of n elements and you use the bracket operator (`[]`) to access one of the elements, this operation takes pretty much the same amount of time regardless of how big the array is.
- Linear: An algorithm is “linear” if the run time is directly proportional to the size of the input. For example, if you add up the elements of an array, you have to access n elements and perform $n-1$ additions. The total number of operations (element accesses and additions) is $2n-1$, which is directly proportional to n .
- Quadratic: An algorithm is “quadratic” if the run time is directly proportional to n^2 . For example, if the input has 2 elements, it might require 4 operations; with 3 elements, it might require 9, and so on.

For example, here’s an implementation of a simple algorithm called https://en.wikipedia.org/wiki/Selection_sort selection sort:

```
public class SelectionSort {  
  
    /**  
     * Swaps the elements at indexes i and j.  
     */  
    public static void swapElements(int[] array, int i, int j) {  
        int temp = array[i];
```

```
        array[i] = array[j];
        array[j] = temp;
    }

    /**
     * Finds the index of the lowest value
     * starting from the index at 'start' (inclusive)
     * and going to the end of the array.
     */
    public static int indexLowest(int[] array, int start) {
        int lowIndex = start;
        for (int i = start; i < array.length; i++) {
            if (array[i] < array[lowIndex]) {
                lowIndex = i;
            }
        }
        return lowIndex;
    }

    /**
     * Sorts the elements (in place) using selection sort.
     */
    public static void selectionSort(int[] array) {
        for (int i = 0; i < array.length; i++) {
            int j = indexLowest(array, i);
            swapElements(array, i, j);
        }
    }
}
```

The first method, `swapElements`, swaps two elements of the array. Reading and writing elements are constant time operations, regardless of the size of the array. That works because if we know where the beginning of the array is, we can compute the location of element `i` or `j` with one multiplication and one addition. And those are constant time operations. Since everything in `swapElements` is constant time, the whole method is constant time.

The second method, `indexLowest` finds the index of the smallest element of the array starting at a given index, `start`. Each time through the loop, it

accesses two elements of the array and performs one comparison. Since these are all constant time operations, it doesn't really matter which ones we count. To keep it simple, let's count the number of comparisons.

1. If `start` is 0, `indexLowest` traverses the entire array, and the total number of comparisons is `n`.
2. If `start` is 1, the number of comparisons is `n-1`.
3. In general, the number of comparisons is `n-start`.

So we normally consider `indexLowest` to be linear time, except for special values of `start`.

The third method, `selectionSort`, sorts the array. It loops from 0 to `n-1`, so the loop executes `n` times. Each time, it performs a constant time operation, `swapElements`, and then calls `indexLowest`.

The first time `selectionSort` calls `indexLowest`, it performs `n` comparisons. The second time it performs `n-1` comparisons, and so on. The total number of comparisons is

$$n + n-1 + n-2 + \dots + 1 + 0$$

The sum of this series is $n(n+1)/2$, which is directly proportional to n^2 ; and that means that `selectionSort` is quadratic.

To get to the same result a different way, we can think of `indexLowest` as a nested loop. Each time we call `indexLowest`, the number of operations is directly proportional to `n`. We call it `n` times, so the total number of operations is directly proportional to n^2 .

2.5 Big O notation

All constant time algorithms belong to a set called $O(1)$. So another way to say that an algorithm is constant time is to say that it is in $O(1)$. Similarly, all linear algorithms belong to $O(n)$, and all quadratic algorithms belong to $O(n^2)$. This way of classifying algorithms is called "big O notation".

NOTE: I am providing a casual definition of big O notation. For a more mathematical treatment, see https://en.wikipedia.org/wiki/Big_O_notation Big O notation.

This notation provides a convenient way to write general rules about how algorithms behave when we compose them. For example, if you perform a linear time algorithm followed by a constant algorithm, the total run time is linear. Using \mathcal{O} to mean “is a member of”:

If $f \in \mathcal{O}(n)$ and $g \in \mathcal{O}(1)$, $f+g \in \mathcal{O}(n)$.

If you perform two linear operations, the total is still linear:

If $f \in \mathcal{O}(n)$ and $g \in \mathcal{O}(n)$, $f+g \in \mathcal{O}(n)$.

In fact, if you perform a linear operation any number of times, k , the total is linear, as long as k is a constant that does not depend on n .

If $f \in \mathcal{O}(n)$ and k is a constant, $kf \in \mathcal{O}(n)$.

But if you perform a linear operation n times, the result is quadratic:

If $f \in \mathcal{O}(n)$, $nf \in \mathcal{O}(n^2)$.

In general, we only care about the largest exponent of n . So if the total number of operations is $2n + 1$, it belongs to $\mathcal{O}(n)$. The leading constant, 2, and the additive term, 1, are not important for this kind of analysis. Similarly, $n^2 + 100n + 100$ is in $\mathcal{O}(n^2)$.

One other piece of vocabulary you should know: an “order of growth” is a set of algorithms whose runtimes grow in the same way as problem size increases; for example, all linear algorithms belong to the same order of growth because their runtimes increase linearly with problem size.

NOTE: In this context, an “order” is a group, like the *Order of the Knights of the Round Table*, which is a group of knights, not a way of lining them up. So you can imagine the *Order of Linear Algorithms* as a set of brave, chivalrous, and particularly efficient knights.

2.6 Resources

http://en.wikipedia.org/wiki/Analysis_of_algorithms Analysis of algorithms

https://en.wikipedia.org/wiki/Selection_sort Selection sort

https://en.wikipedia.org/wiki/Big_O_notation Big O notation # cs-implementing-an-arraylist-lab

DRAFT: Not ready for distribution!

DRAFT: Not ready for distribution!

Chapter 3

ArrayList

3.1 Objectives

1. Write an implementation of an ArrayList.

3.2 Implementing an Array-backed List

For this lesson we provide a partial implementation of an ArrayList that uses a Java array to store the elements. We left four of the methods incomplete; your job is to fill them in. We provide JUnit tests you can use to check your work.

3.3 Instructions

When you check out the repository for this exercise, you should find a file structure similar to what you saw in the previous exercise. The top level directory contains CONTRIBUTING.md, LICENSE.md, README.md, and the directory that contains the code for this lab, javacs-lab02.

In the subdirectory `javacs-lab02/src/com/flatironschool/javacs` you'll find the source files you need for this exercise:

- * 'MyArrayList.java' contains a partial implementation of the 'List' interface using
- * 'MyArrayListTest.java' contains JUnit tests for 'MyArrayList'.

In javacs-lab02, you'll find the Ant build file `build.xml`. If you are in this directory, you should be able to run `ant MyArrayList` to run `MyArrayList.java`, which contains a few simple tests. Or you can run `ant MyArrayListTest` to run the JUnit test.

When you run the tests, several of them should fail. If you examine the source code, you'll find four `TODO` comments indicating which methods you will fill in.

3.4 MyArrayList code

Before you start filling in the missing methods, let's walk through some of the code. Here are the instance variables and the constructor.

```
public class MyArrayList<E> implements List<E> {
    int size;                      // keeps track of the number of elements
    private E[] array;             // stores the elements

    public MyArrayList() {
        array = (E[]) new Object[10];
        size = 0;
    }
}
```

As the comments indicate, `size` keeps track of how many elements are in `MyArrayList`, and `array` is the array that actually contains the elements.

The constructor creates an array of 10 elements, which are initially `null`, and sets `size` to 0. Most of the time, the length of the array is bigger than `size`, so there are unused slots in the array.

One detail about Java: You can't instantiate an array of `T[]`, so you have to instantiate an array of `Object` and then typecast it. You can <http://www.ibm.com/developerworks/java/library/j-jtp01255/index.html> read more about this issue here.

Next we'll look at the method that adds elements to the list. Here's our implementation of `add`:


```
public boolean add(E element) {
    if (size >= array.length) {
        // make a bigger array and copy over the elements
        E[] bigger = (E[]) new Object[array.length * 2];
        System.arraycopy(array, 0, bigger, 0, array.length);
        array = bigger;
    }
    array[size] = element;
    size++;
    return true;
}
```

If there are no unused spaces in the array, we have to create a bigger array and copy over the elements. Then we can store the element in the array and increment `size`.

It might not be obvious why this method returns a boolean, since it seems like it always returns `true`. As always, [https://docs.oracle.com/javase/7/docs/api/java/util/Collection.html#add\(E\)](https://docs.oracle.com/javase/7/docs/api/java/util/Collection.html#add(E)) you can find the answer in the documentation.

It's also not obvious how to analyze the performance of this method. In the normal case, it's constant time, but if we have to resize the array, it's linear. In the next README we'll explain how we can handle this.

We'll look at `get` next, and then you can fill in `set`. Actually, `get` is pretty simple:

```
public T get(int index) {
    if (index < 0 || index >= size) {
        throw new IndexOutOfBoundsException();
    }
    return array[index];
}
```

If the index is out of bounds, it throws an exception; otherwise it reads and returns an element of the array. Notice that it checks whether the index is less than `size`, not `array.length`, so it's not possible to access the unused elements of the array.

3.5 Instructions

- In `MyArrayList.java`, you'll find a stub for `set` that looks like this:

```
public T set(int index, T element) {  
    // TODO: fill in this method.  
    return null;  
}
```

Read [https://docs.oracle.com/javase/7/docs/api/java/util/List.html#set\(int,%20E\)](https://docs.oracle.com/javase/7/docs/api/java/util/List.html#set(int,%20E)) the documentation of `set`, then fill in the body of this method. If you run `MyArrayListTest` again, `testSet` should pass.

HINT: Try to avoid repeating the index-checking code.

- Your next mission is to fill in `indexOf`. As usual, you should [https://docs.oracle.com/javase/7/docs/api/java/util/List.html#indexOf\(java.lang.Object\)](https://docs.oracle.com/javase/7/docs/api/java/util/List.html#indexOf(java.lang.Object)) read the documentation so you know what it's supposed to do. In particular, notice how it is supposed to handle `null`.

To make things a little easier, we've provided a helper method called `equals` that compares an element from the array to a target value and returns `true` if they are equal (and it handles `null` correctly). Notice that this method is `private` because it is used inside this class but it is not part of the `List` interface.

When you are done, run `MyArrayListTest` again; `testIndexOf` should pass now, as well as the other tests that depend on it.

- Only two more methods to go, and you'll be done with this exercise. The next one is an overloaded version of `add` that takes an index and stores the new value at the given index, shifting the other elements to make room, if necessary.

Again, [https://docs.oracle.com/javase/7/docs/api/java/util/List.html#add\(int,%20E\)](https://docs.oracle.com/javase/7/docs/api/java/util/List.html#add(int,%20E)) read the documentation, write an implementation, and run the tests for confirmation.

HINT: Avoid repeating the code that makes the array bigger.

- Last one: fill in the body of `remove`. [https://docs.oracle.com/javase/7/docs/api/java/util/List.html#remove\(int\)](https://docs.oracle.com/javase/7/docs/api/java/util/List.html#remove(int)) The documentation is here. When you finish this one, all tests should pass.

Once you have your implementation working, compare it to mine, which you can find by checking out the solutions branch of the repo, or <https://github.com/learn-co-students/cs-implementing-an-arraylist-lab-codeU/tree/solution> you can read it on GitHub.

View Implementing an Array List on Learn.co and start learning to code for free.

DRAFT: Not ready for distribution.

DRAFT: Not ready for distribution!

Chapter 4

Analyzing ArrayList

4.1 Overview

This README kills two birds with one stone: we present solutions to the previous exercise, and also use them to demonstrate analysis of algorithms.

4.2 Objectives

1. Read and understand solutions to **Implementing an ArrayList** exercise.
2. Classify the methods in `MyArrayList`.
3. Use amortized analysis to classify appropriate algorithms.

4.3 Classifying MyArrayList methods

For many methods, we can identify the order of growth by examining the code. For example, here's the implementation of `get` from `MyArrayList`:

```
public E get(int index) {
    if (index < 0 || index >= size) {
        throw new IndexOutOfBoundsException();
    }
}
```

```
    }  
    return array[index];  
}
```

Everything in `get` is constant time, so `get` is constant time. No problem.

Now that we've classified `get`, we can classify `set`, which uses it. Here is our implementation of `set` from the previous exercise:

```
public E set(int index, E element) {  
    E old = get(index);  
    array[index] = element;  
    return old;  
}
```

One slightly clever part of this solution is that it does not check the bounds of the array explicitly; it takes advantage of `get`, which raises an exception if the index is invalid.

Everything in `set`, including the invocation of `get`, is constant time, so `set` is also constant time.

Next we'll look at some linear methods. For example, here's my implementation of `indexOf`:

```
public int indexOf(Object target) {  
    for (int i = 0; i < size; i++) {  
        if (equals(target, array[i])) {  
            return i;  
        }  
    }  
    return -1;  
}
```

Each time through the loop, `indexOf` invokes `equals`, so we have to classify `equals` first. Here it is:

```
private boolean equals(Object target, Object element) {  
    if (target == null) {  
        return element == null;  
    } else {  
        return target.equals(element);  
    }  
}
```

This method invokes `target.equals`; the runtime of this method might depend on the size of `target` or `element`, but it probably doesn't depend on the size of the array, so we consider it constant time for purposes of analyzing `indexOf`.

Getting back to `indexOf`, everything inside the loop is constant time, so the next question we have to consider is: how many times does the loop execute?

If we get lucky, we might find the target object right away and return after testing only one element. If we are unlucky, we might have to test all of the elements. On average, we expect to test half of the elements, so this method is considered linear except in the (unlikely) case that we know the target element is at the beginning of the array.

The analysis of `remove` is similar. Here's our implementation:

```
public E remove(int index) {
    E element = get(index);
    for (int i=index; i<size-1; i++) {
        array[i] = array[i+1];
    }
    size--;
    return element;
}
```

It uses `get`, which is constant time, and then loops through the array, starting from `index`. If we remove the element at the end of the list, the loop never runs and this method is constant time. If we remove the first element, we loop through all of the remaining elements, which is linear. So, again, this method is considered linear except in the special case where we know the element is at the end.

4.4 Classifying add

In the previous exercise, you wrote a version of `add` that takes an index and an element as parameters. Here's our solution:

```
public void add(int index, E element) {
    if (index < 0 || index > size) {
        throw new IndexOutOfBoundsException();
    }
    // add the element to get the resizing
    add(element);

    // shift the other elements
    for (int i=size-1; i>index; i--) {
        array[i] = array[i-1];
    }
    // put the new one in the right place
    array[index] = element;
}
```

This two-parameter version, which we'll call `add(int, E)`, uses the one-parameter version, which we'll call `add(E)`, to put the new element at the end. Then it shifts the other elements to the right, and puts the new element in the correct place.

Before we can classify the two-parameter `add(int, E)`, we have to classify the one-parameter `add(E)`:

```
public boolean add(E element) {
    if (size >= array.length) {
        // make a bigger array and copy over the elements
        E[] bigger = (E[]) new Object[array.length * 2];
        System.arraycopy(array, 0, bigger, 0, array.length);
        array = bigger;
    }
    array[size] = element;
    size++;
    return true;
}
```

This one-parameter version turns out to be hard to analyze. If there is an unused space in the array, it is constant time, but if we have to resize the array, it's linear. So which is it?

We can classify this method by thinking about the average number of operations per add over a series of n adds. For simplicity, assume we start with an array that has room for 2 elements.

- The first time we call add, it finds unused space in the array, so it stores 1 element.
- The second time, it finds unused space in the array, so it stores 1 element.
- The third time, we have to resize the array, copy 2 elements, and store 1 element. Now the size of the array is 4.
- The fourth time stores 1 element.
- The fifth time resizes the array, copies 4 elements, and stores 1 element. Now the size of the array is 8.
- The next 3 adds store 3 elements.
- The next one copies 8 and stores 1. Now the size is 16
- The next 7 adds store 7 elements.

And so on. Adding things up:

- After 4 adds, we've stored 4 elements and copied 2.
- After 8 adds, we've stored 8 elements and copied 6.
- After 16 adds, we've stored 16 elements and copied 14.

By now you should see the pattern: to do n adds, we have to store n elements and copy $n-2$. So the total number of operations is $n + n - 2$, which is $2n-2$.

To get the average number of operations per add, we divide the total by n ; the result is $2 - 2/n$. As n gets big, the second term, $2/n$, gets small. Invoking the principle that we only care about the largest exponent of n , we can think of add as constant time.

It might seem strange that an algorithm that is sometimes linear can be constant time, on average. The key is that we double the length of the array each time it gets resized. That limits the number of times each element gets copied. Otherwise — if we add a fixed amount to the length of the array, rather than multiplying by a fixed amount — the analysis doesn't work.

This way of classifying an algorithm, by computing the average time in a series of invocations, is called [https://en.wikipedia.org/wiki/Amortized_](https://en.wikipedia.org/wiki/Amortized_time)

analysis amortized analysis. The idea is that the extra cost of copying the array is spread, or “amortized”, over a series of invocations.

Now, if `add(E)` is constant time, what about `add(int, E)`? After calling `add(E)`, it loops through part of the array and shifts elements. This loop is linear, except in the special case where we are adding at the end of the list. So `add(int, E)` is linear.

This is an example of a tricky API, where two methods that seem similar have qualitatively different performance.

4.5 Classifying `removeAll`

The last example we’ll consider is `removeAll`; here’s the implementation in `MyArrayList`:

```
public boolean removeAll(Collection<?> collection) {
    boolean flag = true;
    for (Object obj: collection) {
        flag &= remove(obj);
    }
    return flag;
}
```

Each time through the loop, `removeAll` invokes `remove`, which is linear. So it is tempting to think that `removeAll` is quadratic. But that’s not necessarily the case.

In this method, the loop runs once for each element in the `Collection c`. If `c` contains `m` elements and the list we are removing from contains `n` elements, this method is in $O(nm)$. If the size of the collection can be considered constant, `removeAll` is linear. But if the size of the collection is proportional to `n`, it’s quadratic. For example, if `c` always contains 100 or fewer elements, `removeAll` is linear. But if `c` generally contains 1% of the elements in the list, `removeAll` is quadratic.

Sometimes when we talk about “problem size” we have to be careful about which size, or sizes, we are talking about. This example demonstrates a pitfall

of algorithm analysis, which is the tendency to count loops. If there is one loop, the algorithm is *often* linear. If there are two loops (one nested inside the other), the algorithm is *often* quadratic. But be careful! You have to think about how many times each loop runs. If the number of iterations is proportional to n , you can get away with just counting the loops. But if, as in this example, the number of iterations does not depend on n , you have to give it more thought.

4.6 Resources

https://en.wikipedia.org/wiki/Amortized_analysis Amortized analysis: Wikipedia page.

View Analyzing Our Array List on Learn.co and start learning to code for free.

DRAFT: Not ready for distribution!

Chapter 5

LinkedList

5.1 Objectives

1. Implement a linked list in Java.
2. Write an implementation of the List interface.

5.2 Overview

For this exercise we provide a partial implementation of the List interface that uses a linked list to store the elements. We left three of the methods incomplete; your job is to fill them in. We provide JUnit tests you can use to check your work.

If you are not familiar with a linked list, you might want to read https://en.wikipedia.org/wiki/Linked_list the Wikipedia page about it, but we'll also start with a brief introduction.

5.3 Linked data structures

A data structure is “linked” if it is made up of objects, often called “nodes”, that contain references to other nodes. In a linked *list*, each node contains a

reference to the next node in the list. Other linked structures include trees and graphs, in which nodes can contain references to more than one other node.

A linked list is used to store a sequence of elements, so each node contains a reference to an element, or sometimes to a collection of elements. The element part of the node is sometimes called “cargo”, so you can think of nodes as rail cars, where each car contains cargo and the cars are connected together.

Let’s see what this looks like in code. Here’s a class definition for a simple Node:

```
public class ListNode {  
  
    public Object cargo;  
    public ListNode next;  
  
    public ListNode() {  
        this.cargo = null;  
        this.next = null;  
    }  
  
    public ListNode(Object cargo) {  
        this.cargo = cargo;  
        this.next = null;  
    }  
  
    public ListNode(Object cargo, ListNode next) {  
        this.cargo = cargo;  
        this.next = next;  
    }  
  
    public String toString() {  
        return "ListNode(" + cargo.toString() + ")";  
    }  
}
```

The `ListNode` object has two instance variables: `cargo` is a reference to some kind of `Object`, and `next` is a reference to the next node in the list. In the last node in the list, by convention, `next` is `null`.

ListNode provides several constructors, allowing you to provide values for `cargo` and `next`, or initialize them to the default value, `null`.

You can think of each ListNode as a list with a single element, but more generally, a list can contain any number of nodes. There are several ways to make a new list. A simple way is to create a set of ListNode objects, like this:

```
ListNode node1 = new ListNode(1);
ListNode node2 = new ListNode(2);
ListNode node3 = new ListNode(3);
```

And then link them up, like this:

```
node1.next = node2;
node2.next = node3;
node3.next = null;
```

Alternatively, you can create a Node and link it at the same time. For example, if you want to add a new Node at the beginning of a list, you can do it like this:

```
ListNode node0 = new ListNode(0, node1);
```

After this sequence of instructions, we have four Nodes containing the Integers 0, 1, 2, and 3 as cargo, linked up in increasing order. In the code that has 3 as cargo, the `next` field is `null`.

The following diagram shows these variables and the objects they refer to:

5.4 Getting set up

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory that contains the code for this lab, `javacs-lab03`.

In the subdirectory `javacs-lab03/src/com/flatironschool/javacs` you'll find the source files you need for this exercise:

- * 'MyLinkedList.java' contains a partial implementation of the 'List' interface using
- * 'MyLinkedListTest.java' contains JUnit tests for 'MyLinkedList'.

In `javacs-lab03`, you'll find the Ant build file `build.xml`. If you are in this directory, you should be able to run `ant MyArrayList` to run `MyArrayList.java`, which contains a few simple tests. Or you can run `ant test` to run the JUnit test.

When you run the tests, several of them should fail. If you examine the source code, you'll find three `TODO` comments indicating which methods you will fill in.

5.5 MyLinkedList code

Before you start filling in the missing methods, let's walk through some of the code. Here are the instance variables and the constructor for `MyLinkedList`:

```
public class MyLinkedList<E> implements List<E> {  
  
    private int size;           // keeps track of the number of elements  
    private Node head;         // reference to the first node  
  
    public MyLinkedList() {  
        head = null;  
        size = 0;  
    }  
}
```

As the comments indicate, `size` keeps track of how many elements are in `MyLinkedList`; `head` is a reference to the first `Node` in the list, or `null` if the list is empty.

Storing the number of elements is not necessary, and in general it is risky to keep redundant information, because if it's not updated correctly, it creates opportunities for error. It also takes a little bit of extra space.

But if we store `size` explicitly, we can implement the `size` method in constant time; otherwise, we would have to traverse the list and count the elements, which requires linear time.

Because we store `size` explicitly, we have to update it each time we add or remove an element, so that slows down those methods a little, but it doesn't change their order of growth, so it's probably worth it.

The constructor sets `head` to `null`, which indicates an empty list, and sets `size` to 0.

This class uses the type parameter `E` for the type of the elements. If you are not familiar with type parameters, you might want to read <https://docs.oracle.com/javase/tutorial/java/generics/types.html> this tutorial.

The type parameter also appears in the definition of `Node`, which is nested inside `MyLinkedList`:

```
private class Node {
    public E cargo;
    public Node next;

    public Node(E cargo, Node next) {
        this.cargo = cargo;
        this.next = next;
    }
}
```

Other than that, `Node` is similar to `ListNode` above.

5.6 Adding nodes

At this point you should have a general idea of how this implementation of linked lists works. But before you get started on the exercises, we'll look at one more method. Here's my implementation of `add`:

```
public boolean add(E element) {
    if (head == null) {
        head = new Node(element);
    } else {
        Node node = head;
        // loop until the last node
```

```
        for ( ; node.next != null; node = node.next) {}
        node.next = new Node(element);
    }
    size++;
    return true;
}
```

This example demonstrates two patterns you'll need for your solutions:

1. For many methods, we have to handle the first element of the list as a special case. In this example, if we are adding the first element of a list, we have to modify `head`. Otherwise, we traverse the list, find the end, and add the new node.
2. This method shows how to use a `for` loop to traverse the nodes in a list. In your solutions, you will probably write several variations on this loop. Notice that we have to declare `node` before the loop so we can access it after the loop.

5.7 Instructions

Now it's your turn:

- Fill in the body of `indexOf`. As usual, you should [https://docs.oracle.com/javase/7/docs/api/java/util/List.html#indexOf\(java.lang.Object\)](https://docs.oracle.com/javase/7/docs/api/java/util/List.html#indexOf(java.lang.Object)) read the documentation so you know what it's supposed to do. In particular, notice how it is supposed to handle `null`.

To make things a little easier, I've provided a helper method called `equals` that compares an element from the array to a target value and checks whether they are equal — and it handles `null` correctly. This method is private because it is used inside this class but it is not part of the `List` interface.

When you are done, run `ant test` again; `testIndexOf` should pass now, as well as the other tests that depend on it.

- Next, you should fill in the two-parameter version of `add`, which takes an index and stores the new value at the given index. Again, [https://docs.oracle.com/javase/7/docs/api/java/util/List.html#add\(int,%20E\)](https://docs.oracle.com/javase/7/docs/api/java/util/List.html#add(int,%20E)) read the documentation, write an implementation, and run the tests for confirmation.
- Last one: fill in the body of `remove`. [https://docs.oracle.com/javase/7/docs/api/java/util/List.html#remove\(int\)](https://docs.oracle.com/javase/7/docs/api/java/util/List.html#remove(int)) The documentation is here. When you finish this one, all tests should pass.
- Once you have your implementation working, compare it to mine, which you can find by checking out the solutions branch of the repo, or https://TODO:%20add_this_later you can read it on GitHub.

DRAFT: Not ready for distribution!

Chapter 6

A note on garbage collection

In `MyArrayList` from the previous exercise, the array grows if necessary, but it never shrinks. The array never gets garbage collected, and the elements don't get garbage collected until the list itself is destroyed.

One advantage of the linked list implementation is that it shrinks when elements are removed, and the unused `Nodes` can get garbage collected immediately.

Here is my implementation of the `clear` method:

```
public void clear() {  
    head = null;  
    size = 0;  
}
```

When we set `head` to `null`, we remove a reference to the first `Node`. If there are no other references to that `Node` (and there shouldn't be), it will get garbage collected. At that point, the reference to the second `Node` is removed, so it gets garbage collected, too. This process continues until all `Nodes` are collected.

So how should we classify `clear`? The method itself contains two constant time operations, so it sure looks like it's constant time. But when you invoke it, you make the garbage collector do work that's proportional to the number of elements. So there's an argument that we should consider it linear!

This is a subtle example of what is sometimes called a “performance bug”: a program that is correct in the sense that it does the right thing, but it doesn’t belong to the order of growth we expected. In languages like Java that do a lot of work, like garbage collection, behind the scenes, this kind of error can be hard to find.

6.1 Resources

https://en.wikipedia.org/wiki/Linked_list Linked list Wikipedia page
<https://docs.oracle.com/javase/tutorial/java/generics/types.html> Generic types # cs-analyzing-our-linkedlist-readme

DRAFT: Not ready for distribution!

Chapter 7

Analysis of LinkedList

7.1 Overview

This README kills two birds with one stone: we present solutions to the previous exercise, and continue the discussion of analysis of algorithms.

7.2 Objectives

1. Read and understand solutions to the previous exercise.
2. Classify the methods in `MyLinkedList`.

7.3 Classifying `MyLinkedList` methods

My implementation of `indexOf` is below. Read through it and see if you can classify it before you read the explanation.

```
public int indexOf(Object target) {
    Node node = head;
    for (int i=0; i<size; i++) {
        if (equals(target, node.cargo)) {
            return i;
        }
    }
}
```

```
        }
        node = node.next;
    }
    return -1;
}
```

Initially `node` gets a copy of `head`, so they both refer to the same Node. The loop variable, `i`, counts from 0 to `size-1`. Each time through the loop, we use `equals` to see if we've found the target. If so, we return `i` immediately. Otherwise we advance to the next Node in the list. Normally we would check to make sure the next Node is not `null`, but in this case it is safe because the loop ends when we get to the end of the list (as long as `size` is consistent with the actual number of Nodes in the list).

If we get through the loop without finding the target, we return `-1`.

So what is the order of growth for this method?

1. Each time through the loop we invoke `equals`, which is constant time (it might depend on the size of `target` or `cargo`, but it doesn't depend on the size of the list). The other operations in the loop are also constant time.
2. The loop might run `n` times, because in the worse case, we might have to traverse the whole list.

So this method is linear.

Up next, here is my implementation of the two-parameter `add` method. Again, you should try to classify it before you read the explanation.

```
public void add(int index, E element) {
    if (index == 0) {
        head = new Node(element, head);
    } else {
        Node node = getNode(index-1);
        node.next = new Node(element, node.next);
    }
    size++;
}
```


If `index==0`, we're adding the new `Node` at the beginning, so we handle that as a special case. Otherwise, we have to traverse the list to find the element at `index-1`. We use the helper method `getNode`:

```
private Node getNode(int index) {
    if (index < 0 || index >= size) {
        throw new IndexOutOfBoundsException();
    }
    Node node = head;
    for (int i=0; i<index; i++) {
        node = node.next;
    }
    return node;
}
```

`getNode` checks whether `index` is out of bounds; if so, it throws an exception. Otherwise it traverses the list and returns the requested `Node`.

Jumping back to `add`, once we find the right `Node`, we create the new `Node` and put it between `node` and `node.next`. You might find it helpful to draw a diagram of this operation to make sure you understand it.

So, what's the order of growth for `add`?

1. Starting with `getNode`, it is very similar to `indexOf`, and it is linear for the same reason.
2. In `add`, everything before and after `getNode` is constant time.

So all together, `add` is linear.

Finally, let's look at `remove`:

```
public E remove(int index) {
    E element = get(index);
    if (index == 0) {
        head = head.next;
    } else {
        Node node = getNode(index-1);
        node.next = node.next.next;
    }
}
```

```

    }
    size--;
    return element;
}

```

`remove` uses `get` to find and store the element at `index`. Then it removes the Node that contained it.

If `index==0`, we handle that as a special case again. Otherwise we find the node at `index-1` and modify it to skip over `node.next` and link directly to `node.next.next`. This effectively removes `node.next` from the list, and it can be garbage collected.

Finally, we decrement `size` and return the element we retrieved at the beginning.

So, what's the order of growth for `remove`? Everything in `remove` is constant time except `get` and `getNode`, which are linear. So `remove` is linear.

When people see two linear operations, they sometimes think the result is quadratic, but that only applies if one operation is nested inside the other. If you invoke one operation after the other, the runtimes add. If they are both in $O(n)$, the sum is also in $O(n)$.

7.4 Comparing MyLinkedList and MyArrayList

The following table summarizes the differences between `MyLinkedList` and `MyArrayList`, where `1` means $O(1)$ or constant time and `n` means $O(n)$ or linear.

	MyArrayList	MyLinkedList
add (at the end)	1	n
add (at the beginning)	n	1
add (in general)	n	n

	MyArrayList	MyLinkedList
get / set	1	n
indexOf / lastIndexOf	n	n
isEmpty / size	1	1
remove (from the end)	1	n
remove (from the beginning)	n	1
remove (in general)	n	n

The operations where **MyArrayList** has an advantage are adding at the end, removing from the end, getting and setting.

The operations where **MyLinkedList** has an advantage are adding at the beginning and removing from the beginning.

For the other operations, the two implementations are in the same order of growth.

Which implementation is better? It depends on which operations you are likely to use the most. And that's why Java provides more than one implementation, because it depends. #

DRAFT: Not ready for distribution!

Chapter 8

Performance profiling

8.1 Objectives

1. Classify algorithms by measuring runtime for a range of problem sizes.

8.2 Overview

For this exercise we provide a class called `Profiler` that contains code for running code with a range of problem sizes, measuring runtimes, then plotting and analyzing the results.

You will use `Profiler` to measure and classify the performance of the `add` method for the Java implementations of `ArrayList` and `LinkedList`.

You'll need to “test” that you've got the right output for this lesson by making sure that your results conform to the runtime expectations.

8.3 Profiler

Here's an example that shows how to use the profiler:

```
public static void profileArrayListAddEnd() {
    Timeable timeable = new Timeable() {
        List<String> list;

        public void setup(int n) {
            list = new ArrayList<String>();
        }

        public void timeMe(int n) {
            for (int i=0; i<n; i++) {
                list.add("a string");
            }
        }
    };

    String title = "ArrayList add end";
    Profiler profiler = new Profiler(title, timeable);

    int startN = 4000;
    int endMillis = 1000;
    XYSeries series = profiler.timingLoop(startN, endMillis);
    profiler.plotResults(series);
}
```

This method measures the time it takes to run `add` on an `ArrayList`, which adds the new element at the end. We'll explain the code and then show the results.

In order to use `Profiler`, we need to create a `Timeable` object that provides two methods: `setup` and `timeMe`. The `setup` method does whatever needs to be done before we start the clock; in this case it creates an empty list. Then `timeMe` does whatever operation we are trying to measure; in this case it adds `n` elements to the list.

The code that creates `timeable` is an anonymous class that defines a new implementation of the `Timeable` interface and creates an instance of the new class at the same time. If you are not familiar with anonymous classes, you can <https://docs.oracle.com/javase/tutorial/java/java00/anonymousclasses.html> read about them here. But you really don't need to know much for this

exercise; even if you are not comfortable with anonymous classes, you can do the exercises by copying and modifying the example code.

The next step is to create the `Profiler` object, passing the `Timeable` object and a title as parameters.

The `Profiler` provides `timingLoop` which uses the `Timeable` object stored as an attribute. It invokes the `timeMe` method on the `Timeable` object several times with a range of values of `n`. `timingLoop` takes two parameters:

- `startN` is the value of `n` the timing loop should start at.
- `endMillis` is a threshold in milliseconds. As `timingLoop` increases the problem size, the runtime increases; when the runtime exceeds this threshold, `timingLoop` stops.

When you run these experiments, you might have to adjust these parameters. If `startN` is too low, the runtime might be too short to measure accurately. If `endMillis` is too low, you might not get enough data to see a clear relationship between problem size and run time.

This code is in `ProfileAdd.java`, which you'll run in the next lab. When we ran it, we got this output:

```
4000, 3
8000, 0
16000, 1
32000, 2
64000, 3
128000, 6
256000, 18
512000, 30
1024000, 88
2048000, 185
4096000, 242
8192000, 544
16384000, 1325
```

The first column is problem size, `n`; the second column is runtime in milliseconds. The first few measurements are pretty noisy, it might have been better to set `startN` around 64000.

The result from `timingLoop` is an `XYSeries` that contains this data. If you pass this series to `plotResults`, it generates a plot like this:

The next section explains how to interpret it.

8.4 Interpreting results

Based on our understanding of how `ArrayList` works, we expect the `add` method to take constant time when we add elements to the end. So the total time to add `n` elements should be linear.

To test that theory, we could plot total runtime versus problem size, and we should see a straight line, at least for problem sizes that are big enough to measure accurately. Mathematically, we can write the function for that line:

$$\text{runtime} = a + b * n$$

Where `a` is the intercept of the line and `b` is the slope.

On the other hand, if `add` is linear, the total time for `n` adds would be quadratic. If we plot runtime versus problem size, we expect to see a parabola. Or mathematically, something like:

$$\text{runtime} = a + b * n + c * n^2$$

With perfect data, we might be able to tell the difference between a straight line and a parabola, but if the measurements are noisy, it can be hard to tell. A better way to interpret noisy measurements is to plot runtime and problem size on a **log-log** scale.

Why? Let's suppose that runtime is proportional to n^k , but we don't know what the exponent `k` is. We can write the relationship like this:

$$\text{runtime} = a + b * n + \dots + c * n^k$$

For large values of `n`, the term with the largest exponent is the most important, so:

$$\text{runtime} \approx c * n^k$$

Where \approx is the symbol for “approximately equal”. Now, if we take the log of both sides of this equation:

$$\log(\text{runtime}) \approx \log(c) + k * \log(n)$$

This equation implies that if we plot runtime versus n on a log-log scale, we expect to see a straight line with intercept $\log(c)$ and slope k . We don’t care much about the intercept, but the slope indicates the order of growth: if $k=1$, the algorithm is linear; if $k=2$, it’s quadratic.

Looking at the figure in the previous section, you can estimate the slope by eye. But when you call `plotResults` it computes a least squares fit to the data and prints the estimated slope. In this example:

```
Estimated slope= 1.06194352346708
```

Which is close to 1; and that suggests that the total time for n adds is linear, so each add is constant time, as expected.

One important point: if you see a straight line on a graph like this, that does **not** mean that the algorithm is linear. If the run time is proportional to n^k for any exponent k , we expect to see a straight line with slope k . If the slope is close to 1, that suggests the algorithm is linear. If it is close to 2, it’s probably quadratic.

8.5 Instructions

When you check out the respository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory that contains the code for this lab, `javacs-lab04`.

In the subdirectory `javacs-lab04/src/com/flatiron-school/javacs` you’ll find the source files you need for this exercise:

1. `Profiler.java` contains the implementation of the `Profiler` class described above. You will use this class, but you don’t have to know how it works. But feel free to read the source.

2. `ProfileAdd.java` contains starter code for this exercise, including the example, above, which profiles `ArrayList.add`. You will modify this file to profile a few other methods.

Also, in `javacs-lab04`, you'll find the Ant build file `build.xml`.

- In `javacs-lab04` run `ant ProfileAdd` to run `ProfileAdd.java`. You should get results similar to ours, but you might have to adjust `startN` or `endMillis`. The estimated slope should be close to 1, indicating that performing `n` add operations takes time proportional to `n` raised to the exponent 1; that is, it is $O(n)$.
- In `ProfileAdd.java`, you'll find an empty method named `profileArrayListAddBeginning`. Fill in the body of this method with code that tests `ArrayList.add`, always putting the new element at the beginning. If you start with a copy of `profileArrayListAddEnd`, you should only have to make a few changes. Add a line in `main` to invoke this method.

Run `ant ProfileAdd` again and interpret the results. Based on our understanding of how `ArrayList` works, we expect each add operation to be linear, so the total time for `n` adds should be quadratic. If so, the estimated slope of the line, on a log-log scale, should be near 2. Is it?

- Now let's compare that to the performance of `LinkedList`. Fill in the body of `profileLinkedListAddBeginning` and use it to classify `LinkedList.add` when we put the new element at the beginning. What performance do you expect? Are the results consistent with your expectations?
- Finally, fill in the body of `profileLinkedListAddEnd` and use it to classify `LinkedList.add` when we put the new element at the end. What performance do you expect? Are the results consistent with your expectations?

We'll present our results and answer these questions in the next README.

8.6 Resources

<https://docs.oracle.com/javase/tutorial/java/java00/anonymousclasses.html> Anonymous classes: Java tutorial.

Chapter 9

Doubly-linked list

9.1 Overview

This README reviews results from the previous exercise and introduces yet another implementation of the `List` interface, the doubly-linked list.

9.2 Objectives

1. Interpret results from the previous exercise.
2. Understand the implementation of the doubly-linked list.
3. Analyze the performance of doubly-linked list operations.

9.3 Performance profiling results

In the previous exercise, we used `Profiler.java` to run various `ArrayList` and `LinkedList` operations with a range of problem sizes. We plotted the runtime versus problem size on a log-log scale and estimated the slope of the resulting curve, which indicates the leading exponent of the relationship between run time and problem size. For example, when we used the `add` method to add elements to the end of an `ArrayList`, we found that the total time to perform n adds was proportional to n ; that is, the estimated slope was close to 1. We

concluded that performing n adds is in $O(n)$, so on average the time for a single add is constant time, or $O(1)$, which is what we expected based on algorithm analysis.

The exercise asked you to fill in the body of `profileArrayListAddBeginning`, which tests the performance of adding new elements at the beginning of an `ArrayList`. Based on our analysis, we expect each add to be linear, because it has to shift the other elements to the right; so we expect n adds to be quadratic.

Here's our solution to the exercise, which you can see and run by checking out the `solutions` branch of the exercise:

```
public static void profileArrayListAddBeginning() {
    Timeable timeable = new Timeable() {
        List<String> list;

        public void setup(int n) {
            list = new ArrayList<String>();
        }

        public void timeMe(int n) {
            for (int i=0; i<n; i++) {
                list.add(0, "a string");
            }
        }
    };
    int startN = 4000;
    int endMillis = 10000;
    runProfiler("ArrayList add beginning", timeable, startN, endMillis);
}
```

This method is almost identical to `profileArrayListAddEnd`. The only difference is in `timeMe`, which uses the two-parameter version of `add` to put the new element at index 0. Also, we increased `endMillis` to get one additional data point.

Here are the timing results (problem size on the left, run time in milliseconds on the right):

4000, 14
8000, 35
16000, 150
32000, 604
64000, 2518
128000, 11555

Here's the graph of runtime versus problem size:

Remember that a straight line on this graph does **not** mean that the algorithm is linear. Rather, if the runtime is proportional to nk for any exponent, k , we expect to see a straight line with slope k . In this case, we expect the total time for n adds to be proportional to n^2 , so we expect a straight line with slope 2. In fact, the estimated slope is 1.992, which is so close we would be afraid to fake data this good.

9.4 Profiling LinkedList methods

The next exercise asked you to profile the performance of adding new elements at the beginning of a `LinkedList`. Based on our analysis, we expect each `add` to take constant time, because in a linked list, we don't have to shift the existing elements; we can just add a new node at the beginning. So we expect the total time for n adds to be linear.

Here's our solution:

```
public static void profileLinkedListAddBeginning() {
    Timeable timeable = new Timeable() {
        List<String> list;

        public void setup(int n) {
            list = new LinkedList<String>();
        }

        public void timeMe(int n) {
            for (int i=0; i<n; i++) {
                list.add(0, "a string");
            }
        }
    };
}
```

```
        }  
    };  
    int startN = 128000;  
    int endMillis = 2000;  
    runProfiler("LinkedList add beginning", timeable, startN, endMillis);  
}
```

Again, we only had to make a few small changes, replacing `ArrayList` with `LinkedList` and adjusting `startN` and `endMillis` to get a good range of data. We found that these measurements were noisier than the previous batch; here are the results:

```
128000, 16  
256000, 19  
512000, 28  
1024000, 77  
2048000, 330  
4096000, 892  
8192000, 1047  
16384000, 4755
```

And here's the graph:

It's not a very straight line, and the slope is not exactly 1; the slope of the least squares fit is 1.23. But these results indicate that the total time for n adds is at least approximately $O(n)$, so each add is constant time.

9.5 Adding to the end of a LinkedList

Adding elements at the beginning is one of the operations where we expect `LinkedList` to be faster than `ArrayList`. But for adding elements at the end, we expect `LinkedList` to be slower. In our implementation, we have to traverse the entire list to add an element to the end, which is linear. So we expect the total time for n adds to be quadratic.

Well, it's not.

Here's the code:

```
public static void profileLinkedListAddEnd() {
    Timeable timeable = new Timeable() {
        List<String> list;

        public void setup(int n) {
            list = new LinkedList<String>();
        }

        public void timeMe(int n) {
            for (int i=0; i<n; i++) {
                list.add("a string");
            }
        }
    };
    int startN = 64000;
    int endMillis = 1000;
    runProfiler("LinkedList add end", timeable, startN, endMillis);
}
```

Here are the results:

```
64000, 7
128000, 5
256000, 21
512000, 23
1024000, 73
2048000, 228
4096000, 956
8192000, 911
16384000, 4539
```

And here's the graph:

Again, the measurements are noisy and the line is not perfectly straight, but the estimated slope is 1.24, which is almost exactly what we got adding elements at the beginning, and not very close to 2, which is what we expected based on our analysis. In fact, it is closer to 1, which suggests that adding elements at the end is at least approximately linear. What's going on?

9.6 Doubly-linked list

Our implementation of a linked list, `MyLinkedList`, uses a singly-linked list; that is, each element contains a link to the next, and the `MyArrayList` object itself has a link to the first node.

But if you read <https://docs.oracle.com/javase/7/docs/api/java/util/LinkedList.html> the documentation of `LinkedList`, it says

Doubly-linked list implementation of the List and Deque interfaces. [...] All of the operations perform as could be expected for a doubly-linked list. Operations that index into the list will traverse the list from the beginning or the end, whichever is closer to the specified index.

If you are not familiar with doubly-linked lists, you can https://en.wikipedia.org/wiki/Doubly_linked_list read more about them here, but the short version is:

- Each node contains a link to the next node and a link to the previous node.
- The `LinkedList` object contains links to the first and last elements of the list.

So we can start at either end of the list and traverse it in either direction. As a result, we can add and remove elements from the beginning and the end of the list in constant time!

The following table summarizes the performance we expect from `ArrayList`, `MyLinkedList` (singly-linked), and `LinkedList` (doubly-linked):

	MyArrayList	MyLinkedList	LinkedList
add (at the end)	1	n	1
add (at the beginning)	n	1	1
add (in general)	n	n	n

	MyArrayList	MyLinkedList	LinkedList
get / set	1	n	n
indexOf / lastIndexOf	n	n	n
isEmpty / size	1	1	1
remove (from the end)	1	n	1
remove (from the beginning)	n	1	1
remove (in general)	n	n	n

The doubly-linked implementation is better than **ArrayList** for adding and removing at the beginning, and just as good as **ArrayList** for adding and removing at the end. So the only advantage of **ArrayList** is for **get** and **set**, which require linear time in a linked list, even if it is doubly-linked.

If you know that the runtime of your application depends on the time it takes to **get** and **set** elements, an **ArrayList** might be the better choice. If the runtime depends on adding and removing elements near the beginning or the end, **LinkedList** might be better.

But remember that these recommendations are based on the order of growth for large problems. There are other factors to consider:

- If these operations don't take up a substantial fraction of the runtime for your application – that is, if your applications spends most of its time doing other things – then your choice of a **List** implementation won't matter very much.
- If the lists you are working with are not very big, you might not get the performance you expect. For small problems, an quadratic algorithm might be faster than a linear algorithm, or linear might be faster than constant time. And for small problems, the difference probably doesn't matter.

- Also, don't forget about space. So far we have focused on runtime, but different implementations require different amounts of space. In an `ArrayList`, the elements are stored side-by-side in a single chunk of memory, so there is very little wasted space, and computer hardware is often faster with contiguous chunks. In a linked list, each element requires a node with one or two links. The links take up space (sometimes more than the cargo!), and with nodes scattered around in memory, the hardware might be less efficient.

In summary, analysis of algorithms provides some guidance for choosing data structures, but only if

1. The runtime of your application is important,
2. The runtime of your application depends on your choice of data structure, and
3. The problem size is large enough that the order of growth actually predicts which data structure is better.

You could have a long career as a software engineer without ever finding yourself in this situation.

9.7 Resources

https://en.wikipedia.org/wiki/Doubly_linked_list Doubly linked list at Wikipedia.

View Interpreting Performance Profiles on Learn.co and start learning to code for free.

Chapter 10

Tree traversal

10.1 Learning goals

1. Parse HTML using jsoup.
2. Traverse a DOM tree.
3. Use the `Deque` interface and its implementations.

10.2 Overview

This README introduces the application we will develop during the upcoming units. We describe the elements of a Web search engine and introduce the first application, a Web crawler that downloads and parses pages from Wikipedia. We present a recursive implementation of depth-first search and an iterative implementation that uses a Java `Deque` to implement a “last in, first out” stack.

10.3 The road ahead

At the end of each unit, you will have a chance to work on an application that’s a bit more substantial than an exercise. You will write more code, and you’ll have to make more design decisions, rather than just filling in methods.

By the time we get to the end of this track, you will have built a simple **Web search engine**, which is a tool, like Google Search and Bing, that takes a set of “search terms” and returns a list of web pages that are relevant to those terms (we’ll discuss what “relevant” means later). If you are not familiar with search engines, https://en.wikipedia.org/wiki/Web_search_engine you can read more here, but we’ll explain what you need as we go along.

The essential pieces of a search engine are:

- **Crawling:** We’ll need a program that can download a web page, parse it, and extract the text and any links to other pages.
- **Indexing:** We’ll need an index that makes it possible to look up a search term and find the pages that contain it.
- **Retrieval:** And we’ll need a way to collect results from the Index and identify pages that are most relevant to the search terms.

We’ll start with the Crawler; the goal of a Crawler is to discover and download a set of web pages. For search engines like Google and Bing, the goal is to find *all* web pages, but often Crawlers are limited to a smaller domain. In our case, we will only read pages from Wikipedia.

As a first step, we’ll build a Crawler that reads a Wikipedia page, finds the first link, follows the link to another page, and repeats. We will use this Crawler to test the “Getting to Philosophy” conjecture, which states:

Clicking on the first lowercase link in the main text of a Wikipedia article, and then repeating the process for subsequent articles, usually eventually gets one to the Philosophy article.

This conjecture is https://en.wikipedia.org/wiki/Wikipedia:Getting_to_Philosophy stated on this Wikipedia page, and you can read its history there.

Testing the conjecture will allow us to build the basic pieces of a Crawler without having to crawl the entire web, or even all of Wikipedia. And we think the exercise is kind of fun!

In the next unit, we’ll work on the Indexer, and then we’ll get to the Retriever.

10.4 Parsing HTML

When you download a web page, the contents are written in <https://en.wikipedia.org/wiki/HTML> HyperText Markup Language, aka HTML. For example, here is a minimal HTML document:

```
<!DOCTYPE html>
<html>
  <head>
    <title>This is a title</title>
  </head>
  <body>
    <p>Hello world!</p>
  </body>
</html>
```

The phrases “This is a title” and “Hello world!” are the text that actually appears on the page; the other elements are **tags** that indicate how the text should be displayed.

When our Crawler downloads a page, it will need to parse the HTML in order to extract the text and find the links. To do that, we’ll use **jsoup**, which is an open-source Java library that downloads and parses HTML.

The result of parsing HTML is a Document Object Model tree or **DOM tree**, that contains the elements of the document, including text and tags. The tree is a linked data structure made up of nodes; the nodes represent text, tags, and other document elements.

The relationships between the nodes are determined by the structure of the document. In the example above, the first node, called the **root**, is the **<html>** tag, which contains links to the two nodes it contains, **<head>** and **<body>**; these nodes are the **children** of the root node.

The **<head>** node has one child, **<title>**, and the **<body>** node has one child, **<p>** (which stands for “paragraph”). The following figure represents this tree graphically:

Each node contains links to its children; in addition, each node contains a link to its **parent**, so from any node it is possible to navigate up and down

the tree. The DOM tree for real pages is usually more complicated than this example.

Most web browsers provide tools for inspecting the DOM of the page you are viewing. In Chrome, you can right-click on any part of a web page and select “Inspect” from the menu that pops up. In Firefox, you can right-click and select “Inspect Element” from the menu. Safari provides a tool called Web Inspector, https://developer.apple.com/library/mac/documentation/AppleApplications/Conceptual/Safari_Developer_Guide/GettingStarted/GettingStarted.html#//apple_ref/doc/uid/TP40007874-CH2-SW1 which you can read about here. For Internet Explorer, <http://support.janova.us/entries/20181293-How-to-inspect-an-Element> can read the instructions here.

Here is a screenshot of the DOM for [https://en.wikipedia.org/wiki/Java_\(programming_language\)](https://en.wikipedia.org/wiki/Java_(programming_language)) the Wikipedia page on Java.

The element that’s highlighted is the first paragraph of the main text of the article, which is contained in a `<div>` element with `id="mw-content-text"`. We’ll use this element id to identify the main text of each article we download.

10.5 Using jsoup

Jsoup makes it easy to download and parse web pages, and to navigate the DOM tree. Here’s an example:

```
String url = "https://en.wikipedia.org/wiki/Java_(programming_language)";

// download and parse the document
Connection conn = Jsoup.connect(url);
Document doc = conn.get();

// select the content text and pull out the paragraphs.
Element content = doc.getElementById("mw-content-text");

// TODO: avoid selecting paragraphs from sidebars and boxouts
Elements paras = content.select("p");
```

`Jsoup.connect` takes a URL as a `String` and makes a connection to the web server; the `get` method downloads the HTML, parses it, and returns a `Document` object, which represents the DOM.

`Document` provides lots of helpful methods for navigating the tree and selecting nodes. In fact, it provides so many methods, it can be confusing. This example demonstrates two ways to select nodes:

- `getElementById` takes a `String` and searches the tree for an element that has a matching “id” field. Here it selects the node `<div id="mw-content-text" lang="en">` which appears on every Wikipedia page to identify the `<div>` element that contains the main text of the page, as opposed to the navigation sidebar and other elements.

The return value from `getElementById` is an `Element` object that represents this `<div>` and contains the elements in the `<div>` as children, grandchildren, etc.

- `select` takes a `String`, traverses the tree, and returns all the elements with tags that match the `String`. In this example, it returns all paragraph tags that appear in `content`. The return value is an `Elements` object, which is a `Collection` that contains multiple elements.

Before you go on, you should skim the documentation of these classes so you know what they can do. The most important classes are: <http://jsoup.org/apidocs/org/jsoup/nodes/Element.html> `Element`, <http://jsoup.org/apidocs/org/jsoup/select/Elements.html> `Elements`, and <http://jsoup.org/apidocs/org/jsoup/nodes/Node.html> `Node`.

It’s easy to get `Node` and `Element` confused: `Node` is the superclass of `Element`, so every `Element` is a `Node`, but not every `Node` is an `Element`. Other kinds of `Nodes` include `TextNode` (which we will use soon), `DataNode`, and `Comment`.

10.6 Iterating through the DOM

To make your life easier, we’ve provided a class called `WikiNodeIterable` that lets you iterate through the nodes in a DOM tree. Here’s an example that shows how to use it:

```
Elements paras = content.select("p");
Element firstPara = paras.get(0);

Iterable<Node> iter = new WikiNodeIterable(firstPara);
for (Node node: iter) {
    if (node instanceof TextNode) {
        System.out.print(node);
    }
}
```

This example picks up where the previous one leaves off. It selects the first paragraph in `paras` and creates a `WikiNodeIterable`, which implements `Iterable<Node>`. The `for` loop iterates the nodes in the tree using a “depth first search”, which produces the nodes in the order they would appear on the page.

In this example, we print a `Node` only if it is a `TextNode` and ignore other types of `Node`, specifically the `Element` objects that represent tags. The result is the plain text of the HTML paragraph without any markup. The output is:

Java is a general-purpose computer programming language that is concurrent, class-based, object-oriented,[13] and specifically designed...

At this point, you know what you need for the next exercise, so you can skip ahead. But if you are not familiar with depth-first search, you should keep reading.

10.7 Depth-first search

There are several ways you might reasonably traverse a tree, each with different applications. We’ll start with “depth-first search,” or DFS. DFS starts at the root of the tree and selects the first child. If the child has children, it selects the first child again. When it gets to a node with no children, it backtracks, moving up the tree to the parent node, where it selects the next child if there is one; otherwise it backtracks again. When it has explored the last child of the root, it’s done.

There are two common ways to implement DFS, recursively and iteratively. The recursive implementation is simple and elegant:

```
private static void recursiveDFS(Node node) {
    if (node instanceof TextNode) {
        System.out.print(node);
    }
    for (Node child: node.childNodes()) {
        recursiveDFS(child);
    }
}
```

This method gets invoked on every `Node` in the tree, starting with the root. If the `Node` it gets is a `TextNode`, it prints the contents. If the `Node` has any children, it invokes `recursiveDFS` on each one of them in order.

In this example, we print the contents of each `TextNode` before traversing the children, so this is an example of a “pre-order” traversal. https://en.wikipedia.org/wiki/Tree_traversal You can read about “pre-order”, “post-order”, and “in-order” traversals here. For this application, the traversal order doesn’t matter.

By making recursive calls, `recursiveDFS` uses the https://en.wikipedia.org/wiki/Call_stack call stack to keep track of the child nodes and process them in the right order. As an alternative, we can use a stack data structure to keep track of the nodes ourselves; if we do that, we can avoid the recursion and traverse the tree iteratively.

10.8 Stacks in Java

Before we can explain the iterative version of DFS, we have to explain the stack data structure. We’ll start with the general concept of a stack, which we’ll call a “stack” with a lowercase “s”. Then we’ll talk about two Java implementations of a stack, which are called **Stack** and **Deque**.

A stack is a data structure that is similar to a list: it is a collection that maintains the order of the elements. The primary difference between a stack and a list is that the stack provides fewer methods. In the usual convention, it provides:

- **push**: which adds an element to the top of the stack.
- **pop**: which removes the top-most element from the stack.
- **peek**: which returns the top-most element without modifying the stack.
- **isEmpty**: which indicates whether the stack is empty.

Because **pop** always returns the top-most element, a stack is also called a “LIFO”, which stands for “last in, first out”. An alternative to a stack is a “queue”, which returns elements in the same order they are added; that is, “first in, first out”, or FIFO.

It might not be obvious why stacks and queues are useful: they don’t provide any capabilities that aren’t provided by lists; in fact, they provide fewer capabilities. So why not use lists for everything? There are two reasons:

1. If you limit yourself to a small set of methods — that is, a small API — your code will be more readable and less error-prone. For example, if you use a list to represent a stack, you might accidentally remove an element in the wrong order. With the stack API, this kind of mistake is literally impossible. And the best way to avoid errors is to make them impossible.
2. If a data structure provides a small API, it is easier to implement efficiently. For example, a simple way to implement a stack is a singly-linked list. When we push an element onto the stack, we add it to the beginning of the list; when we pop an element, we remove it from the beginning. For a linked list, adding and removing from the beginning are constant time operations, so this implementation is efficient. Conversely, big APIs are harder to implement efficiently.

To implement a stack in Java, you have three options:

1. Go ahead and use **ArrayList** or **LinkedList**. If you use **ArrayList**, be sure to add and remove from the *end*, which is a constant time operation. And be careful not to add elements in the wrong place or remove them in the wrong order.
2. Java provides a class called **Stack** that provides the standard set of stack methods. But this class is an old part of Java: it is not consistent with the Java Collections Framework, which came later, and it does not handle synchronization as well as more recent classes.

3. Probably the best choice is to use one of the implementations of the `Deque` interface, like `ArrayDeque`.

“Deque” stands for “double-ended queue”; it’s supposed to be pronounced “deck”, but some people say “deek”. In Java, the `Deque` interface provides `push`, `pop`, `peek`, and `isEmpty`, so you can use a `Deque` as a stack. It provides other methods <https://docs.oracle.com/javase/7/docs/api/java/util/Deque.html> you can read about here, but we won’t use them for now.

10.9 Iterative DFS

Here is an iterative version of DFS that uses an `ArrayDeque` to represent a stack of `Node` objects:

```
private static void iterativeDFS(Node root) {
    Deque<Node> stack = new ArrayDeque<Node>();
    stack.push(root);

    while (!stack.isEmpty()) {
        Node node = stack.pop();
        if (node instanceof TextNode) {
            System.out.print(node);
        }

        List<Node> nodes = new ArrayList<Node>(node.childNodes());
        Collections.reverse(nodes);

        for (Node child: nodes) {
            stack.push(child);
        }
    }
}
```

The parameter, `root`, is the root of the tree we want to traverse, so we start by creating the stack and pushing the root onto it.

The loop continues until the stack is empty. Each time through, it pops a `Node` off the stack. If it gets a `TextNode`, it prints the contents. Then it pushes the

children onto the stack. In order to process the children in the right order, we have to push them onto the stack in reverse order; we do that by copying the children into an `ArrayList`, reversing the elements in place, and then iterating through the reversed `ArrayList`.

One advantage of the iterative version of DFS is that it is easier to implement as a Java `Iterator`; you'll find out how in the next lab. But first, one last note about the `Deque` interface: in addition to `ArrayDeque`, Java provides another implementation of `Deque`, our old friend `LinkedList`. `LinkedList` implements both interfaces, `List` and `Deque`. Which interface you get depends on how you use it. For example, if you assign a `LinkedList` object to a `Deque` variable, like this:

```
Deque<Node> deque = new LinkedList<Node>();
```

You can use the methods in the `Deque` interface, but not all methods in the `List` interface. If you assign it to a `List` variable, like this:

```
List<Node> deque = new LinkedList<Node>();
```

You can use `List` methods but not all `Deque` methods. And if you assign it like this:

```
LinkedList<Node> deque = new LinkedList<Node>();
```

You can use *all* the methods. But if you combine methods from different interfaces, your code will be less readable and more error-prone.

10.10 Resources

https://en.wikipedia.org/wiki/Web_search_engine Web search engine: Wikipedia.

<https://en.wikipedia.org/wiki/HTML> HyperText Markup Language: Wikipedia

<http://jsoup.org/apidocs/org/jsoup/nodes/Element.html> `Element`: jsoup documentation
<http://jsoup.org/apidocs/org/jsoup/select/Elements.html> `Elements`: jsoup documentation
<http://jsoup.org/apidocs/org/jsoup/nodes/Node.html> `Node`: jsoup documentation

https://en.wikipedia.org/wiki/Tree_traversal Tree traversal: Wikipedia

https://en.wikipedia.org/wiki/Call_stack Call stack: Wikipedia

<https://docs.oracle.com/javase/7/docs/api/java/util/Deque.html> `Deque` interface: Java documentation # cs-wikipedia-philosophy-lab

Chapter 11

Getting to Philosophy

11.1 Objectives

1. Select data structures for an application.
2. Use Java data structures to write a simple Web crawler.

11.2 Overview

The goal of this exercise is to write a Web crawler that tests the “Getting to Philosophy” conjecture, which we presented in the previous exercise. We’ll provide some code to help you get started, but you will write more code for this exercise than for the previous ones.

11.3 Getting started

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory that contains the code for this lab, `javacs-lab05`.

In the subdirectory `javacs-lab05/src/com/flatironschool/javacs` you’ll find the source files you need for this exercise:

1. `WikiNodeExample.java` contains the code from the previous README, demonstrating recursive and iterative implementations of depth-first search (DFS) in a DOM tree.
2. `WikiNodeIterable.java` contains an `Iterable` class for traversing a DOM tree. We'll explain this code in the next section.
3. `WikiFetcher.java` contains a utility class that uses jsoup to download pages from Wikipedia. To help you comply with Wikipedia's terms of service, this class limits how fast you can download pages; if you request more than one page per second, it sleeps before downloading the next page.
4. `WikiPhilosophy.java` contains an outline of the code you will write for this exercise. We'll walk you through it below.

Also, in `javacs-lab05`, you'll find the Ant build file `build.xml`. If you run `ant WikiPhilosophy`, it will run a simple bit of starter code.

11.4 Iterables and Iterators

In the previous README, we presented an iterative depth-first search (DFS), and suggested that an advantage of the iterative version, compared to the recursive version, is that it is easier to wrap in an `Iterator` object. In this section we'll see how to do that.

If you are not familiar with the `Iterator` and `Iterable` interfaces, you can read about them <https://docs.oracle.com/javase/7/docs/api/java/util/Iterator.html> here and <https://docs.oracle.com/javase/7/docs/api/java/lang/Iterable.html> here.

Take a look at the contents of `WikiNodeIterable.java`. The outer class, `WikiNodeIterable` implements the `Iterable<Node>` interface, so we can use it in an “enhanced for loop” like this:

```
Node root = ...
Iterable<Node> iter = new WikiNodeIterable(root);
for (Node node: iter) {
    visit(node);
}
```

Where `root` is the root of the tree we want to traverse and `visit` is a method that does whatever we want when we “visit” a `Node`.

The implementation of `WikiNodeIterable` follows a conventional formula:

1. The constructor takes and stores a reference to the root `Node`.
2. The `iterator` method creates and returns an `Iterator` object.

Here’s what it looks like:

```
public class WikiNodeIterable implements Iterable<Node> {  
  
    private Node root;  
  
    public WikiNodeIterable(Node root) {  
        this.root = root;  
    }  
  
    @Override  
    public Iterator<Node> iterator() {  
        return new WikiNodeIterator(root);  
    }  
}
```

The inner class, `WikiNodeIterator`, does all the real work:

```
private class WikiNodeIterator implements Iterator<Node> {  
  
    Deque<Node> stack;  
  
    public WikiNodeIterator(Node node) {  
        stack = new ArrayDeque<Node>();  
        stack.push(root);  
    }  
  
    @Override  
    public boolean hasNext() {  
        return !stack.isEmpty();  
    }  
}
```

```
@Override
public Node next() {
    if (stack.isEmpty()) {
        throw new NoSuchElementException();
    }

    Node node = stack.pop();
    List<Node> nodes = new ArrayList<Node>(node.childNodes());
    Collections.reverse(nodes);
    for (Node child: nodes) {
        stack.push(child);
    }
    return node;
}
```

This code is almost identical to the iterative version of DFS, but now it's split into three methods:

1. The constructor initializes the stack (which is implemented using an `ArrayDeque`) and pushes the root node onto it.
2. `isEmpty` checks whether the stack is empty.
3. `next` pops the next `Node` off the stack, pushes its children in reverse order, and returns the `Node` it popped. If someone invokes `next` on an empty `Iterator`, it throws an exception.

11.5 WikiFetcher

When you write a Web crawler, it is easy to download too many pages too fast, which might violate the terms of service of the server you are downloading from. To help you avoid that, we provide a class called `WikiFetcher` that does two things

1. It encapsulates the code we demonstrated in the previous README for downloading pages from Wikipedia, parsing the HTML, and selecting the content text.

2. It measures the time between requests and, if we don't leave enough time between requests, it sleeps until a reasonable interval has elapsed. By default, the interval is one second.

Here's the definition of WikiFetcher

```
public class WikiFetcher {
    private long lastRequestTime = -1;
    private long minInterval = 1000;

    /**
     * Fetches and parses a URL string, returning a list of paragraph elements
     *
     * @param url
     * @return
     * @throws IOException
     */
    public Elements fetchWikipedia(String url) throws IOException {
        sleepIfNeeded();

        Connection conn = Jsoup.connect(url);
        Document doc = conn.get();
        Element content = doc.getElementById("mw-content-text");
        Elements paras = content.select("p");
        return paras;
    }

    private void sleepIfNeeded() {
        if (lastRequestTime != -1) {
            long currentTime = System.currentTimeMillis();
            long nextRequestTime = lastRequestTime + minInterval;
            if (currentTime < nextRequestTime) {
                try {
                    Thread.sleep(nextRequestTime - currentTime);
                } catch (InterruptedException e) {
                    System.err.println("Warning: sleep interrupted in fetchWiki");
                }
            }
        }
    }
}
```

```
        lastRequestTime = System.currentTimeMillis();  
    }  
}
```

The only public method is `fetchWikipedia`, which takes a URL as a `String` and returns an `Elements` collection that contains one DOM elements for each paragraph in the content text. This code should look familiar.

The new code is in `sleepIfNeeded`, which checks the time since the last request and sleeps if the elapsed time is less than `minInterval`, which is in milliseconds.

That's all there is to `WikiFetcher`. Here's an example that demonstrates how it's used:

```
WikiFetcher wf = new WikiFetcher();  
  
for (String url: urlList) {  
    Elements paragraphs = wf.fetchWikipedia(url);  
    processParagraphs(paragraphs);  
}
```

In this example, we assume that `urlList` is a collection of `Strings`, and `processParagraphs` is a method that does something with the `Elements` object returned by `fetchWikipedia`.

This example demonstrates something important: you should create one `WikiFetcher` object and use it to handle all requests. If you have multiple instances of `WikiFetcher`, they won't enforce the minimum interval between requests.

NOTE: My implementation of `WikiFetcher` is simple, but it would be easy for someone to mis-use it by creating multiple instances. You could avoid this problem by making `WikiFetcher` a “singleton”, https://en.wikipedia.org/wiki/Singleton_pattern which you can read about here.

11.6 Filling in WikiPhilosophy

In `WikiPhilosophy.java` you'll find a simple `main` method that shows how to use some of these pieces. Starting with this code, your job is to write a crawler that:

1. Takes a URL for a Wikipedia page, downloads it, and parses it.
2. It should traverse the resulting DOM tree to find the first *valid* link. We'll explain what "valid" means below.
3. If the page has no links, or if the first link is a page we have already seen, the program should indicate failure and exit.
4. If the link matches the URL of the Wikipedia page on philosophy, the program should indicate success and exit.
5. Otherwise it should go back to Step 1.

The program should build a **List** of the URLs it visits and display the results at the end (whether it succeeds or fails).

So what should we consider a "valid" link? You have some choices here. Various versions of the "Getting to Philosophy" conjecture use slightly different rules, but here are some options:

1. The link should be in the content text of the page, not in a sidebar or boxout.
2. It should not be in italics or in parentheses.
3. You should skip external links, links to the current page, and red links.
4. In some versions, you should skip a link if the text starts with an uppercase letter.

You don't have to enforce all of these rules, but we recommend that you at least handle parentheses, italics, and links to the current page.

If you feel like you have enough information to get started, go ahead. Or you might want to read these hints:

1. As you traverse the tree, the two kinds of **Node** you will need to deal with are **TextNode** and **Element**. If you find an **Element**, you will probably have to typecast it to access the tag and other information.

2. When you find an `Element` that contains a link, you can check whether it is in italics by following parent links up the tree. If there is an `i` or `em` tag in the parent chain, the link is in italics.
3. To check whether a link is in parentheses, you will have to scan through the text as you traverse the tree and keep track of opening and closing parentheses (ideally your solution should be able to handle nested parentheses (like this)).
4. If you start from [https://en.wikipedia.org/wiki/Java_\(programming_language\)](https://en.wikipedia.org/wiki/Java_(programming_language)) the Java page, you should get to <https://en.wikipedia.org/wiki/Philosophy> after following seven links (unless something has changed since we ran the code).

Ok, that's all the help you're going to get from us. Now it's up to you. Have fun!

11.7 Resources

<https://docs.oracle.com/javase/7/docs/api/java/util/Iterator.html> Iterator: Java documentation.

<https://docs.oracle.com/javase/7/docs/api/java/lang/Iterable.html> Iterable: Java documentation.

https://en.wikipedia.org/wiki/Singleton_pattern Singleton pattern: Wikipedia.

Chapter 12

Indexer

12.1 Learning goals

1. Understand the design of a Web indexer and fill in missing methods.
2. Choose appropriate data structures from the Java Collections Framework (JCF).
3. Use Java's `Map` interface and the `HashMap` implementation.
4. Use the `Set` interface and the `HashSet` implementation.

12.2 Overview

At this point we have built a basic Web crawler; the next piece we will work on is the **index**. In the context of Web search, an index is a data structure that makes it possible to look up a search term and find the pages where that term appears. In addition, we would like to know how many times the search term appears on each page, which will help identify the pages most relevant to the term.

For example, if a user submits the search terms “Java” and “programming”, we would look up both search terms and get two sets of pages. Pages with the word “Java” would include pages about the island of Java, the nickname for coffee, and the programming language. Pages with the word “programming”

would include pages about different programming languages, as well as other uses of the word. By selecting pages with both terms, we hope to eliminate irrelevant pages and find the ones about Java programming.

Now that we understand what the index is and what operations it performs, we can design a data structure to represent it.

12.3 Data structure selection

The fundamental operation of the index is a **lookup**; specifically, we need the ability to look up a term and find all pages that contain it. The simplest implementation would be a collection of pages. Given a search term, we could iterate through the contents of the pages and select the ones that contain the search term. But the runtime would be proportional to the total number of words on all the pages, which would be way too slow.

A better alternative is a **map**, which is a data structure that represents a collection of **key-value pairs** and provides a fast way to look up a **key** and find the corresponding **value**. For example, the first map we'll construct is a **TermCounter**, which maps from each search term to the number of times it appears in a page. The keys are the search terms and the values are the counts (also called "frequencies").

Java provides an interface called **Map** that specifies the methods a map should provide; the most important are:

- **get(key)**: This method looks up a key and returns the corresponding value.
- **put(key, value)**: This method adds a new key-value pair to the **Map**, or if the key is already in the map, it replaces the value associated with **key**.

Java provides several implementations of **Map**, including the two we will focus on, **HashMap** and **TreeMap**. In upcoming labs, we'll look at these implementations and analyze their performance.

In addition to the **TermCounter**, which maps from search terms to counts, we will define a class called **Index**, which maps from a search term to a collection

of pages where it appears. And that raises the next question, which is how to represent a collection of pages. Again, if we think about the operations we want to perform, that guides our decision.

In this case, we'll need to combine two or more collections and find the pages that appear in all of them. You might recognize this operation as **set intersection**: the intersection of two sets is the set of elements that appear in both.

As you might expect by now, Java provides a **Set** interface that defines the operations a set should perform. It doesn't actually provide set intersection, but it provides methods that make it possible to implement intersection and other set operations efficiently. The core **Set** methods are:

- **add(element)**: This method adds an element to a set; if the element is already in the set, it has no effect.
- **contains(element)**: This method checks whether the given element is in the set.

Again, Java provides several implementation of **Set**; the ones we will focus on are **HashSet** and **TreeSet**. We'll come back to them later.

Now that we've designed our data structures from the top down, we'll implement them from the inside out, starting with **TermCounter**.

12.4 TermCounter

TermCounter is a class that represents a mapping from search terms to the number of times they appear in a page. Here is the first part of the class definition:

```
public class TermCounter {  
  
    private Map<String, Integer> map;  
    private String exerciseel;  
  
    public TermCounter(String exerciseel) {
```

```
        this.label = exerciseel;
        this.map = new HashMap<String, Integer>();
    }
}
```

The instance variables are `map`, which contains the mapping from terms to counts, and `label`, which identifies the document the terms came from; we'll use it to store URLs.

To implement the mapping, I chose `HashMap`, which is the most commonly-used `Map`. Coming up in a few lessons, you will see how it works and why it is a common choice.

`TermCounter` provides `put` and `get`, which are defined like this:

```
public void put(String term, int count) {
    map.put(term, count);
}

public Integer get(String term) {
    Integer count = map.get(term);
    return count == null ? 0 : count;
}
```

`put` is just a **wrapper method**; when you call `put` on a `TermCounter`, it calls `put` on the embedded `map`.

On the other hand, `get` actually does some work. When you call `get` on a `TermCounter`, it calls `get` on the `map`, and then it checks the result. If the term does not appear in the `map`, `TermCounter.get` returns 0. Defining `get` this way makes it easier to write `incrementTermCount`, which takes a term and increases by one the counter associated with that term.

```
public void incrementTermCount(String term) {
    put(term, get(term) + 1);
}
```

If the term has not been seen before, `get` returns 0; we add 1, then use `put` to add a new key-value pair to the `map`. If the term is already in the `map`, we get the old count, add 1, and then store the new count, which replaces the old value.

In addition, `TermCounter` provides these other methods to help with indexing Web pages:

```
public void processElements(Elements paragraphs) {
    for (Node node: paragraphs) {
        processTree(node);
    }
}

public void processTree(Node root) {
    for (Node node: new WikiNodeIterable(root)) {
        if (node instanceof TextNode) {
            processText(((TextNode) node).text());
        }
    }
}

public void processText(String text) {
    String[] array = text.replaceAll("\\pP", " ").toLowerCase().split("\\s+");

    for (int i=0; i<array.length; i++) {
        String term = array[i];
        incrementTermCount(term);
    }
}
```

- `processElements` takes an `Elements` object, which is a collection of `jsoup Element` objects. It iterates through the collection and calls `processTree` on each.
- `processTree` takes a `jsoup Node` that represents the root of a DOM tree. It iterates through the tree to find the nodes that contain text; then it extracts the text and passes it to `processText`.
- `processText` takes a `String` that contains words, spaces, punctuation, etc. It removes punctuation characters by replacing them with spaces, converts the remaining letters to lowercase, then splits the text into words. Then it loops through the words it found and calls `incrementTermCount` on each.

Finally, here's an example that demonstrates how `TermCounter` is used:

```
String url = "https://en.wikipedia.org/wiki/Java_(programming_language)";
WikiFetcher wf = new WikiFetcher();
Elements paragraphs = wf.fetchWikipedia(url);

TermCounter counter = new TermCounter(url);
counter.processElements(paragraphs);
counter.printCounts();
```

This example uses a `WikiFetcher` to download a page from Wikipedia and parse the main text. Then it creates a `TermCounter` and uses it to count the words in the page.

In the next section, you'll have a chance to run this code and test your understanding by filling in a missing method.

12.5 Finishing off `TermCounter`

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory with the code for this exercise, `javacs-lab06`.

In the subdirectory `javacs-lab06/src/com/flatiron/school/javacs` you'll find the source files for this exercise:

- * `'TermCounter.java'` contains the code from the previous section.
- * `'Index.java'` contains the class definition for the next part of this exercise.
- * `'WikiFetcher.java'` contains the class we used in the previous exercise to download
- * `'WikiNodeIterable.java'` contains the class we used to traverse the nodes in a DOM t

Also, in `javacs-lab06`, you'll find the Ant build file `build.xml`.

- In `javacs-lab06`, run `ant build` to compile the source files. Then run `ant TermCounter`, it should run the code from the previous section and print a list of terms and their counts. The output should look something like this

```
genericservlet, 2
configurations, 1
claimed, 1
servletresponse, 2
occur, 2
Total of all counts = -1
```

When you run it, the order of the terms might be different.

- The last line is supposed to print the total of the term counts, but it returns `-1` because the method `size` is incomplete. Fill in this method and run `ant TermCounter` again. The result should be `4798`.

Run `ant test1` to confirm that this part of the exercise is complete and correct.

12.6 Finishing off Index

For the second part of the exercise, we'll present our implementation of an `Index` object and you will fill in a missing method. Here's the beginning of the class definition:

```
public class Index {

    private Map<String, Set<TermCounter>> index = new HashMap<String, Set<TermCounter>>();

    public void add(String term, TermCounter tc) {
        Set<TermCounter> set = get(term);

        // if we're seeing a term for the first time, make a new Set
        if (set == null) {
            set = new HashSet<TermCounter>();
        }
    }
}
```

```

        index.put(term, set);
    }
    // otherwise we can modify an existing Set
    set.add(tc);
}

public Set<TermCounter> get(String term) {
    return index.get(term);
}

```

The instance variable, `index`, is a map from each search term to a set of `TermCounter` objects. Each `TermCounter` represents a page where the search term appears.

The `add` method adds a new `TermCounter` to the set associated with a term. When we index a term that has not appeared before, we have to create a new set. Otherwise we can just add a new element to an existing set. In that case, `set.add` modifies a set that lives inside `index`, but doesn't modify `index` itself. The only time we modify `index` is when we add a new term.

Finally, the `get` method takes a search term and returns the corresponding set of `TermCounter` objects.

This data structure is moderately complicated. To review, an `Index` object contains a map from each search term to a set of `TermCounter` objects, and each `TermCounter` is a map from search terms to counts. The method `printIndex` shows how to unpack this data structure:

```

public void printIndex() {
    // loop through the search terms
    for (String term: keySet()) {
        System.out.println(term);

        // for each term, print the pages where it appears and the count
        Set<TermCounter> tcs = get(term);
        for (TermCounter tc: tcs) {
            Integer count = tc.get(term);
            System.out.println("    " + tc.getLabel() + " " + count);
        }
    }
}

```

The outer loop iterates the search terms. The inner loop iterates the `TermCounter` objects that represent the pages where the search term appears.

- Run `ant build` to make sure your source code is compiled, and then run `ant Index`. It downloads two Wikipedia pages, indexes them, and prints the results; but when you run it you won't see any output because we've left one of the methods empty.

Your job is to fill in `indexPage`, which takes a URL (as a `String`) and an `Elements` object, and updates the index. The comments below sketch what it should do:

```
public void indexPage(String url, Elements paragraphs) {  
    // make a TermCounter and count the terms in the paragraphs  
  
    // for each term in the TermCounter, add the TermCounter to the index  
}
```

- When it's working, run `ant Index` again, and you should see output like this:

```
...  
configurations  
  https://en.wikipedia.org/wiki/Programming_language 1  
  https://en.wikipedia.org/wiki/Java_(programming_language) 1  
claimed  
  https://en.wikipedia.org/wiki/Java_(programming_language) 1  
servletresponse  
  https://en.wikipedia.org/wiki/Java_(programming_language) 2  
occur  
  https://en.wikipedia.org/wiki/Java_(programming_language) 2
```

The order of the search terms might be different when you run it.

Also, run `ant test2` to confirm that this part of the exercise is complete.

This exercise involved more reading than some of the others, and not as much programming. But we hope you learned a lot. Congratulations on getting to the end!

DRAFT: Not ready for distribution!

Chapter 13

The Map interface

13.1 Learning goals

1. Write a simple implementation of the `Map` interface.

13.2 Overview

In the next few exercises, we present several implementations of the `Map` interface. One of them, like the `HashMap` provided by Java, is based on a **hash table**, which is arguably the most magical data structure ever invented. Another, which is similar to `TreeMap`, is not quite as magical, but it has the added capability that it can iterate the element in order.

You will have a chance to implement these data structures, and then we will analyze their performance.

But before we can explain hash tables, we'll start with a simple implementation of a `Map` using a `List` of key-value pairs.

13.3 Implementing `MyLinearMap`

As usual, we provide starter code and you will fill in the missing methods. Here's the beginning of the `MyLinearMap` class definition:

```
public class MyLinearMap<K, V> implements Map<K, V> {  
  
    private List<Entry> entries = new ArrayList<Entry>();
```

This class uses two type parameters, *K*, which is the type of the keys, and *V*, which is the type of the values. `MyLinearMap` implements `Map`, which means it has to provide the methods in the `Map` interface.

A `MyLinearMap` object has a single instance variable, `entries`, which is an `ArrayList` of `Entry` objects. Each `Entry` contains a key-value pair. Here is the definition:

```
public class Entry implements Map.Entry<K, V> {  
    private K key;  
    private V value;  
  
    public Entry(K key, V value) {  
        this.key = key;  
        this.value = value;  
    }  
  
    @Override  
    public K getKey() {  
        return key;  
    }  
  
    @Override  
    public V getValue() {  
        return value;  
    }  
}
```

There's not much to it; an `Entry` is just a container for a key and a value. This definition is nested inside `MyLinearList`, so it uses the same type parameters, *K* and *V*.

That's all you need to do the exercise, so let's get started.

13.4 Instructions

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory with the code for this exercise, `javacs-lab07`.

In the subdirectory `javacs-lab07/src/com/flatiron/school/javacs` you'll find the source files for this exercise:

- * `'MyLinearMap.java'` contains starter code for the first part of the exercise.
- * `'MyLinearMapTest.java'` contains the unit tests for `'MyLinearMap'`.

And in `javacs-lab07`, you'll find the Ant build file `build.xml`.

- In `javacs-lab07`, run `ant build` to compile the source files. Then run `ant test`, which runs `MyLinearMapTest`. Several tests should fail, because you have some work to do!
- First, fill in the body of `findEntry`. This is a helper function that is not part of the `Map` interface, but once you get it working, you can use it for several methods. Given a target key, it should search through the entries and return the entry that contains the target (as a key, not a value) or `null` if it's not there. Notice that we have provided an `equals` method that compares two keys and handles `null` correctly.

You can run `ant test` again, but even if your `findEntry` is correct, the tests won't pass because `put` is not complete.

- Fill in `put`. You should read the documentation of [https://docs.oracle.com/javase/7/docs/api/java/util/Map.html#put\(K,%20V\)](https://docs.oracle.com/javase/7/docs/api/java/util/Map.html#put(K,%20V)) so you know what it is supposed to do. You might want to start with a version of `put` that always adds a new entry and does not modify an existing entry; that way you can test the simple case first. Or if you feel more confident, you can write the whole thing at once.

Once you've got `put` working, the test for `containsKey` should pass.

- Read the documentation of [https://docs.oracle.com/javase/7/docs/api/java/util/Map.html#get\(java.lang.Object\)](https://docs.oracle.com/javase/7/docs/api/java/util/Map.html#get(java.lang.Object)) and then fill in the method. Run the tests again.

- Finally, read the documentation of [https://docs.oracle.com/javase/7/docs/api/java/util/Map.html#remove\(java.lang.Object\)](https://docs.oracle.com/javase/7/docs/api/java/util/Map.html#remove(java.lang.Object)) `Map.remove` and then fill in the method.

At this point, all tests should pass. Congratulations!

In the next lesson, we'll present our solutions, analyze the performance of the core `Map` methods, and introduce a more efficient implementation.

DRAFT: Not ready for distribution!

Chapter 14

Hashing

14.1 Learning goals

1. Analyze the performance of a simple `Map` implementation.
2. Implement the `Map` interface using multiple lists.
3. Implement a hash table.
4. Write a hash function for String-like objects.

14.2 Overview

In this readme, we'll look at solutions to the previous exercise and analyze the performance of `MyLinearMap`. Then we'll define `MyBetterMap`, an implementation of the `Map` interface that uses many instances of `MyLinearMap`; and we'll introduce **hashing**, which makes `MyBetterMap` more efficient. Finally, we'll demonstrate a function that implements hashing for a String-like object.

14.3 Analyzing `MyLinearMap`

We'll start with our solutions to the previous exercise and then analyze their performance. Here are `findEntry` and `equals`:

```
private Entry findEntry(Object target) {
    for (Entry entry: entries) {
        if (equals(target, entry.getKey())) {
            return entry;
        }
    }
    return null;
}

private boolean equals(Object target, Object obj) {
    if (target == null) {
        return obj == null;
    }
    return target.equals(obj);
}
```

The runtime of `equals` might depend on the size of the `target` and the `keys`, but does not generally depend on the number of entries, `n`. So `equals` is constant time.

In `findEntry`, we might get lucky and find the key we're looking for at the beginning, but we can't count on it. In general, the number of entries we have to search is proportional to `n`, so `findEntry` is linear.

Most of the core methods in `MyLinearMap` use `findEntry`, including `put`, `get`, and `remove`. Here's what they look like:

```
public V put(K key, V value) {
    Entry entry = findEntry(key);
    if (entry == null) {
        entries.add(new Entry(key, value));
        return null;
    } else {
        V oldValue = entry.getValue();
        entry.setValue(value);
        return oldValue;
    }
}
```

```
public V get(Object key) {
    Entry entry = findEntry(key);
    if (entry == null) {
        return null;
    }
    return entry.getValue();
}

public V remove(Object key) {
    Entry entry = findEntry(key);
    if (entry == null) {
        return null;
    } else {
        V value = entry.getValue();
        entries.remove(entry);
        return value;
    }
}
```

After `put` calls `findEntry`, everything else is constant time. Remember that `entries` is an `ArrayList`, so adding an element *at the end* is constant time, on average. If the key is already in the map, we don't have to add an entry, but we have to call `entry.getValue` and `entry.setValue`, and those are both constant time. Adding it all up, `put` is linear.

By the same reasoning, `get` is also linear.

`remove` is slightly more complicated because `entries.remove` might have to remove an element from the beginning or middle of the `ArrayList`, and that takes linear time. But that's ok: two linear operations are still linear.

In summary, the core methods are all linear, which is why we called this implementation `MyLinearMap` (ta-da!).

If we know that the number of entries will be small, this implementation might be good enough, but we can do better. In fact, there is an implementation of `Map` where all of the core methods are constant time. When you first hear that, it might not seem possible. What we are saying, in effect, is that you can find a needle in a haystack in constant time, regardless of how big the haystack is. It's magic.

We'll explain how it works in two steps:

1. Instead of storing entries in one big `List`, we'll break them up into lots of short lists. For each key, we'll use a **hash code** (explained in the next section) to determine which list to use.
2. Using lots of short lists is faster than using just one, but as we'll explain, it doesn't change the order of growth; the core operations are still linear. But there is one more trick: if we increase the number of lists to limit the number of entries per list, the result is a constant-time map. You'll see the details in the next exercise, but first: hashing!

14.4 Hashing

To improve the performance of `MyLinearMap`, we'll define a new class, called `MyBetterMap`, that contains a collection of `MyLinearMap` objects. It divides the keys among the embedded maps, so the number of entries in each map is smaller, which speeds up `findEntry` and the methods that depend on it.

Here's the beginning of the class definition:

```
public class MyBetterMap<K, V> implements Map<K, V> {  
  
    protected List<MyLinearMap<K, V>> maps;  
  
    public MyBetterMap(int k) {  
        makeMaps(k);  
    }  
  
    protected void makeMaps(int k) {  
        maps = new ArrayList<MyLinearMap<K, V>>(k);  
        for (int i=0; i<k; i++) {  
            maps.add(new MyLinearMap<K, V>());  
        }  
    }  
}
```

The instance variable, `maps` is a collection of `MyLinearMap` objects. The constructor takes a parameter, `k`, that determines how many maps to use, at least initially. Then `makeMaps` creates the embedded maps and stores them in an `ArrayList`.

Now, the key to making this work is that we need some way to look at a key and decide which of the embedded maps it should go into. When we **put** a new key, we choose one of the maps; when we **get** the same key, we have to remember where we put it.

One possibility is to choose one of the sub-maps at random and keep track of where we put each key. But how should we keep track? It might seem like we could use a `Map` to look up the key and find the right sub-map, but the whole point of this exercise is to write an efficient implementation of a `Map`. We can't assume we already have one.

A better approach is to use a **hash function**, which takes an `Object`, any `Object`, and returns an integer called a **hash code**. Importantly, if it sees the same `Object` more than once, it always returns the same hash code. That way, if we use the hash code to store a key, we'll get the same hash code when we look it up.

In Java, every `Object` provides a method called `hashCode` that computes a hash function. The implementation of this method is different for different objects; we'll see an example soon.

Here's the helper function we wrote to choose the right sub-map for a given key:

```
protected MyLinearMap<K, V> chooseMap(Object key) {
    int index = key==null ? 0 : key.hashCode() % maps.size();
    return maps.get(index);
}
```

If `key` is `null`, we choose the sub-map with index 0, arbitrarily. Otherwise we use `hashCode` to get an integer and then apply the modulus operator, `%`, which guarantees that the result is between 0 and `maps.size()-1`. So `index` is always a valid index into `maps`. Then `chooseMap` returns a reference to the map it chose.

We use `chooseMap` in both `put` and `get`, so when we look up a key, we get the same map we chose when we added the key. At least, we should — we'll explain a little later why this might not work.

Here's our implementation of `put` and `get`:

```
public V put(K key, V value) {  
    MyLinearMap<K, V> map = chooseMap(key);  
    return map.put(key, value);  
}  
  
public V get(Object key) {  
    MyLinearMap<K, V> map = chooseMap(key);  
    return map.get(key);  
}
```

Pretty simple, right? In both methods, we use `chooseMap` to find the right sub-map and then invoke a method on the sub-map. In the next lab, you'll have a chance to finish off a few other methods. But first, let's think about performance.

If there are n entries split up among k sub-maps, there will be n/k entries per map, on average. When we look up a key, we have to compute its hash code, which takes some time, then we search the corresponding sub-map. Because the entry lists in `MyBetterMap` are k times shorter than the entry list in `MyLinearMap`, we expect the search to be k times faster. But the run time is still proportional to n , so `MyBetterMap` is still linear. In the next exercise, you'll see how we can fix that.

But first, a little more about hashing.

14.5 How does hashing work?

The fundamental requirement for a hash function is that the same object should produce the same hash code every time. For immutable objects, that's relatively easy. For objects with mutable state, we have to think harder.

As an example of an immutable object, we'll define a class called `SillyString` that encapsulates a `String`:


```
public class SillyString {
    private final String innerString;

    public SillyString(String innerString) {
        this.innerString = innerString;
    }

    public String toString() {
        return innerString;
    }
}
```

This class is not very useful, which is why it's called `SillyString`, but we'll use it to show how a class can define its own hash function:

```
@Override
public boolean equals(Object other) {
    return this.toString().equals(other.toString());
}

@Override
public int hashCode() {
    int total = 0;
    for (int i=0; i<innerString.length(); i++) {
        total += innerString.charAt(i);
    }
    System.out.println(total);
    return total;
}
```

Notice that `SillyString` overrides both `equals` and `hashCode`. This is important. In order to work properly, `equals` has to be consistent with `hashCode`, which means that if two objects are considered equal — that is, `equals` returns `true` — they should have the same hash code. But this requirement only works one way; if two objects have the same hash code, they don't necessarily have to be equal.

`equals` works by invoking `toString`, which returns `innerString`. So two `SillyString` objects are equal if their `innerString` attributes are equal.

`hashCode` works by iterating through the characters in the `String` and adding them up. When you add a character to an `int`, Java converts the character

to an integer using its Unicode code point. You don't need to know anything about Unicode to understand this example, but if you are curious, https://en.wikipedia.org/wiki/Code_point you can read more here.

This hash function satisfies the requirement: if two `SillyString` objects contain embedded strings that are equal, they will get the same hash code.

This works correctly, but the performance might not be very good, because it returns the same hash code for many different strings. If two strings contain the same letters in any order, they will have the same hash code. And even if they don't contain the same letters, they might yield the same total, like `"ac"` and `"bb"`.

If many objects have the same hash code, they end up in the same sub-map, and some sub-maps have more entries than others. In that case, the speedup when we have `k` maps might be much less than `k`. So one of the goals of a hash function is to be uniform; that is, it should be equally likely to produce any value in the range. You can https://en.wikipedia.org/wiki/Hash_function read more about designing good hash functions here.

14.6 Hashing and mutation

Strings are immutable and `SillyString` is also immutable because `innerString` is declared to be `final`. Once you create a `SillyString`, you can't make `innerString` refer to a different `String`, and you can change the `String` it refers to. Therefore, it will always have the same hash code.

But let's see what happens with a mutable object. Here's a definition for `SillyArray`, which is identical to `SillyString`, except that it uses an array of characters instead of a `String`:

```
public class SillyArray {
    private final char[] array;

    public SillyArray(char[] array) {
        this.array = array;
    }
}
```

```
public String toString() {
    return Arrays.toString(array);
}

@Override
public boolean equals(Object other) {
    return this.toString().equals(other.toString());
}

@Override
public int hashCode() {
    int total = 0;
    for (int i=0; i<array.length; i++) {
        total += array[i];
    }
    System.out.println(total);
    return total;
}
```

`SillyArray` also provides `setChar`, which makes it possible to modify the characters in the array:

```
public void setChar(int i, char c) {
    this.array[i] = c;
}
```

Now suppose we create a `SillyArray` and add it to a map:

```
SillyArray array1 = new SillyArray("Word1".toCharArray());
map.put(array1, 1);
```

The hash code for this array is 461. Now if we modify the contents of the array and they try to look it up:

```
array1.setChar(0, 'C');
Integer value = map.get(array1);
```

The hash code after the mutation is 441. With a different hash code, there's a good chance we'll go looking in the wrong sub-map. In that case, we won't find the key, even though it is in the map. And that's bad.

In general, it is dangerous to use mutable objects as keys in data structures that use hashing, which includes `MyBetterMap` and `HashMap`. If you can guarantee that the keys won't be modified while they are in the map, or that any changes won't affect the hash code, it might be okay. But it is probably a good idea to avoid it.

14.7 Resources

https://en.wikipedia.org/wiki/Code_pointUnicode code point: Wikipedia.

https://en.wikipedia.org/wiki/Hash_functionHash function: Wikipedia.

DRAFT: Not ready for distribution!

Chapter 15

Hash table implementation

15.1 Learning goals

1. Implement the `Map` interface using a hash table.
2. Analyze the performance of the hash table implementation.
3. Identify and fix performance bugs.

15.2 Overview

In this exercise, you will finish off the implementation of `MyBetterMap`. In the previous README, we showed that the core methods of `MyBetterMap` are faster than `MyLinearMap`, but they are still linear. We'll see how to solve that problem by allowing the hash table to grow, which you will implement in `MyHashMap`. Then we will profile the performance of `MyHashMap` and the `HashMap` provided by Java.

15.3 Finishing `MyBetterMap`

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory

contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory with the code for this exercise, `javacs-lab08`.

In the subdirectory `javacs-lab08/src/com/flatiron/school/javacs` you'll find the source files for this exercise:

- `MyLinearMap.java` contains our solution to the previous exercise, which we will build on in this exercise.
- `MyBetterMap.java` contains the code from the previous README with some methods you will fill in.
- `MyHashMap.java` contains the outline of a hash table that grows when needed, which you will complete.
- `MyLinearMapTest.java` contains the unit tests for `MyLinearMap`.
- `MyBetterMapTest.java` contains the unit tests for `MyBetterMap`.
- `MyHashMapTest.java` contains the unit tests for `MyHashMap`.
- `Profiler.java` contains code for measuring and plotting runtime versus problem size.
- `ProfileMapPut.java` contains code that profiles the `Map.put` method.

Also, in `javacs-lab08`, you'll find the Ant build file `build.xml`.

- In `javacs-lab08`, run `ant build` to compile the source files. Then run `ant test1`, which runs `MyBetterMapTest`. Several tests should fail, because you have some work to do!
- Review the implementation of `put` and `get` from the previous README. Then fill in the body of `containsKey`. Hint: use `chooseMap`. Run `ant test1` again and confirm that `testContainsKey` passes.
- Fill in the body of `containsValue`. Hint: *don't* use `chooseMap`. Run `ant test1` again and confirm that `testContainsValue` passes. Notice that we have to do more work to find a value than to find a key.

Like `put` and `get`, this implementation of `containsKey` is linear, because it has to search one of the embedded sub-maps. If there are n entries and k sub-maps, the size of the sub-maps is n/k on average, which is still proportional to n .

But if we increase k along with n , we can limit the size of n/k . For example, suppose we double k every time n exceeds k ; in that case the number of entries per map would be less than 1 on average, and pretty much always less than 10, as long as the hash function spreads out the keys reasonably well.

If the number of entries per sub-map is constant, we can search a single sub-map in constant time. And computing the hash function is generally constant time (it might depend on the size of the key, but does not depend on the number of keys). So we might be able to make `put` and `get` constant time.

We'll see how that works in the next section.

15.4 Implementing MyHashMap

In `MyHashMap.java`, we provide the outline of a hash table that grows when needed. Here's the beginning of the definition:

```
public class MyHashMap<K, V> extends MyBetterMap<K, V> implements Map<K, V> {

    // average number of entries per sub-map before we rehash
    private static final double FACTOR = 1.0;

    @Override
    public V put(K key, V value) {
        V oldValue = super.put(key, value);

        // check if the number of elements per sub-map exceeds the threshold
        if (size() > maps.size() * FACTOR) {
            rehash();
        }
        return oldValue;
    }
}
```

`MyHashMap` extends `MyBetterMap`, so it inherits the methods defined there. The only method it overrides is `put` which calls `put` in the superclass — that is, it calls the version of `put` in `MyBetterMap` — and then it checks whether it has to rehash. Calling `size` returns the total number of entries, `n`. Calling `maps.size` returns the number of embedded maps, `k`.

The constant `FACTOR`, which is called the **load factor**, determines the maximum number of entries per sub-map, on average. If $n > k * \text{FACTOR}$, that means $n/k > \text{FACTOR}$, which means the number of entries per sub-map exceeds the threshold, so we call `rehash`.

- In `javacs-lab08`, run `ant build` to compile the source files. Then run `ant test2`, which runs `MyHashMapTest`. It should fail because our implementation of `rehash` throws an exception. Your job is to fill it in.

Fill in the body of `rehash` to collect the entries in the table, resize the table, and then put the entries back in. We provide two methods that might come in handy: `MyBetterMap.makeMaps` and `MyLinearMap.getEntries`. Your solution should double the number of maps, `k`, each time it is called.

15.5 Analyzing MyHashMap

If the number of entries per sub-map, n/k , is constant, several of the core `MyBetterMap` methods become constant time:

```
public boolean containsKey(Object target) {
    MyLinearMap<K, V> map = chooseMap(target);
    return map.containsKey(target);
}

public V get(Object key) {
    MyLinearMap<K, V> map = chooseMap(key);
    return map.get(key);
}

public V remove(Object key) {
    MyLinearMap<K, V> map = chooseMap(key);
    return map.remove(key);
}
```


Each method hashes a key, which is constant time, and then invokes a method on a sub-map, which is constant time.

So far, so good. But the other core method, `put`, is a little harder to analyze. When we don't have to rehash, it is constant time, but when we do, it's linear. In that way, it's similar to `ArrayList.add`, and for the same reason, it turns out to be constant time if we average over a series of `put` operations. Again, the argument is based on https://en.wikipedia.org/wiki/Amortized_analysis amortized analysis.

Suppose the initial number of sub-maps, `k`, is 2, and the load factor is 1. Now let's see how much work it takes to `put` a series of keys. As the basic "unit of work", we'll count the number of times we have to hash a key and add it to a sub-map.

The first time we call `put` it takes 1 unit of work. The second time also takes 1 unit. The third time we have to rehash, so it takes 2 units to rehash the existing keys and 1 unit to hash the new key.

Now the size of the hash table is 4, so the next time we call `put`, it takes 1 unit of work. But the next time we have to rehash, which takes 4 units to rehash the existing keys and 1 unit to hash the new key.

The following figure shows the pattern, with the normal work of hashing a new key shown across the bottom and extra work of rehashing shown as a tower.

As the arrows suggest, if we knock down the towers, each one fills the space before the next tower. The result is a uniform height of 2 units, which shows that the average work per `put` is about 2 units. And that means that `put` is constant time on average.

This diagram also shows why it is important to double the number of sub-maps, `k`, when we rehash. If we only add to `k` instead of multiplying, the towers would be too close together and (eventually) they would start piling up. And that would not be constant time.

15.6 The tradeoffs

We've shown that `containsKey`, `get`, and `remove` are constant time, and `put` is constant time on average. We should take a minute to appreciate how

remarkable that is. The performance of these operations is pretty much the same no matter how big the hash table is. Well, sort of.

Remember that our analysis is based on a simple model of computation where each “unit of work” takes the same amount of time. Real computers are more complicated than that. In particular, they are usually fastest when working with data structures small enough to fit in cache; somewhat slower if the structure doesn’t fit in cache but still fits in memory; and *much* slower if the structure doesn’t fit in memory.

Another limitation of this implementation is that hashing doesn’t help if we want to look up a value: `containsValue` is still linear because it has to search all of the sub-maps. And there is no particularly efficient way to look up a value and find the corresponding key (or possibly keys).

And there’s one more limitation: some of the methods that were constant time in `MyLinearMap` have become linear. For example;

```
public void clear() {  
    for (int i=0; i<maps.size(); i++) {  
        maps.get(i).clear();  
    }  
}
```

`clear` has to clear all of the sub-maps, and the number of sub-maps is proportional to `n`, so it’s linear. Fortunately, this operation is not used very often, so for most applications this tradeoff is acceptable.

15.7 Profiling MyHashMap

Before we go on, we should check whether `MyHashMap` is really constant time.

- In `javacs-lab08`, run `ant build` to compile the source files. Then run `ant ProfileMapPut`. It measures the runtime of `HashMap.put` (provided by Java) with a range of problem sizes, and plots runtime versus problem size on a log-log scale. If this operation is constant time, the total time for `n` operations should be linear, so the result should be a straight line with slope 1. When we ran this code, the estimated slope was close to 1, which is consistent with our analysis. You should get something similar.

- Modify `ProfileMapPut.java` so it profiles your implementation, `MyHashMap`, instead of Java's `HashMap`. Run the profiler again and see if the slope is near 1. You might have to adjust `startN` and `endMillis` to find a range of problem sizes where the runtimes are more than a few milliseconds, but not more than a few thousand.
- When we ran this code, we got a surprise: the slope was about 1.7, which suggests that our implementation is not constant time after all. It contains a “performance bug”. As the last exercise for this exercise, you should track down the error, fix it, and confirm that `put` is constant time, as expected.

15.8 Resources

- https://en.wikipedia.org/wiki/Amortized_analysis Amortized analysis: Wikipedia.

DRAFT: Not ready for distribution!

Chapter 16

Fixing the hash tbale

16.1 Learning goals

1. Fix a performance bug.
2. Use UML class diagrams to represent relationships between classes.

16.2 Overview

This README presents our solution to the previous exercise: a map implementation called `MyFixedHashMap` because it fixes the performance bug in `MyHashMap`. We'll also introduce the UML class diagram, which is a graphical representation of the relationships between classes.

16.3 Fixing MyHashMap

The problem with `MyHashMap` is in `size`, which is inherited from `MyBetterMap`:

```
public int size() {  
    int total = 0;  
    for (MyLinearMap<K, V> map: maps) {  
        total += map.size();  
    }  
}
```

```
    }  
    return total;  
}
```

To add up the total size it has to iterate the sub-maps. Since we increase the number of sub-maps, **k**, as the number of entries, **n**, increases, **k** is proportional to **n**, so **size** is linear.

And that makes **put** linear, too, because it uses **size**:

```
public V put(K key, V value) {  
    V oldValue = super.put(key, value);  
  
    if (size() > maps.size() * FACTOR) {  
        rehash();  
    }  
    return oldValue;  
}
```

Everything we did to make **put** constant time is wasted if **size** is linear!

Fortunately, there is a simple solution, and we have seen it before: we have to keep the number of entries in an instance variable and update it whenever we call a method that changes it.

If you check out the **solution** branch of the previous exercise, **javacs-lab08**, you'll find our solution in **MyFixedHashMap.java**.

Here's the beginning of the class definition:

```
public class MyFixedHashMap<K, V> extends MyHashMap<K, V> implements Map<K, V> {  
  
    private int size = 0;  
  
    public void clear() {  
        super.clear();  
        size = 0;  
    }  
}
```

Rather than modify `MyHashMap`, we define a new class that extends it. It adds a new instance variable, `size`, which is initialized to zero.

Updating `clear` is straightforward; we invoke `clear` in the superclass (which clears the sub-maps), and then update `size`.

Updating `remove` and `put` is a little more difficult because when we invoke the method on the superclass, we can't tell whether the size of the sub-map changed. Here's how we worked around that:

```
public V remove(Object key) {
    MyLinearMap<K, V> map = chooseMap(key);
    size -= map.size();
    V oldValue = map.remove(key);
    size += map.size();
    return oldValue;
}
```

`remove` uses `chooseMap` to find the right sub-map, then subtracts away the size of the sub-map. It invokes `remove` on the sub-map, which may or may not change the size of the sub-map, depending on whether it finds the key. But either way, we add the new size of the sub-map back to `size`, so the final value of `size` is correct.

The rewritten version of `put` is similar:

```
public V put(K key, V value) {
    MyLinearMap<K, V> map = chooseMap(key);
    size -= map.size();
    V oldValue = map.put(key, value);
    size += map.size();

    if (size() > maps.size() * FACTOR) {
        size = 0;
        rehash();
    }
    return oldValue;
}
```

We have the same problem here: when we invoke `put` on the sub-map, we don't know whether it added a new entry. So we use the same solution, subtracting off the old size and then adding in the new size.

Now the implementation of the `size` method is simple:

```
public int size() {  
    return size;  
}
```

And that's pretty clearly constant time.

When we profiled this solution, we found that the total time for putting `n` keys is proportional to `n`, which means that each `put` is constant time, as it's supposed to be.

16.4 UML class diagrams

One challenge of working with the code for this exercise is that we have several classes that depend on each other. Here are some of the relationships between the classes:

- `MyLinearMap` contains a `LinkedList` and implements `Map`.
- `MyBetterMap` contains many `MyLinearMap` objects and implements `Map`.
- `MyHashMap` extends `MyBetterMap`, so it also contains `MyLinearMap` objects, and it implements `Map`.
- `MyFixedHashMap` also extends `MyHashMap`, and it implements `Map`.

To help keep track of relationships like these, software engineers often use “UML class diagrams”. UML stands for https://en.wikipedia.org/wiki/Unified_Modeling_Language Unified Modeling Language; a class diagram is one of several graphical standards defined by UML.

In a class diagram, each class is represented by a box, and relationships between classes are represented by arrows. Here is a UML class diagram for the classes from the previous exercise:

Different relationships are represented by different arrows:

- Arrows with a solid head indicate HAS-A relationships. For example, each instance of `MyBetterMap` contains multiple instances of `MyLinearMap`, so they are connected by a solid arrow.
- Arrows with a hollow head and a solid line indicate IS-A relationships. For example, `MyHashMap` extends `MyBetterMap`, so they are connected by an IS-A arrow.
- Arrows with a hollow head and a dashed line indicate that one class implements another; in this diagram, every class implements `Map` (except `Map` itself).

UML class diagrams provide a concise way to represent a lot of information about a collection of classes. They are used during design phases to communicate about alternative designs, during implementation phases to maintain a shared mental map of the project, and during deployment to document the design.

16.5 Resources

<http://yuml.me/> yUML is the online tool we used to draw the diagram in this README.

DRAFT: Not ready for distribution!

Chapter 17

TreeMap

17.1 Learning goals

1. Implement the `Map` interface using a binary search tree.
2. Analyze the performance of a tree-backed map.

17.2 Overview

At this point you should be familiar with the `Map` interface and the `HashMap` implementation provided by Java. And by making your own `Map` using a hash table, you should understand how `HashMap` works and why we expect its core methods to be constant time.

Because of this performance, `HashMap` is widely used, but it is not the only implementation of `Map`. There are two reasons you might want another implementation:

1. Hashing can be slow, so even though `HashMap` operations are constant time, the “constant” might be big.
2. The keys in a hash table are not stored in any particular order; in fact, the order might change when the table grows and the keys are rehashed. For some applications, it is necessary, or at least useful, to keep the keys in order.

It turns out to be hard to solve both of these problems at the same time, but Java provides an implementation called `TreeMap` that comes close:

1. The order of growth is not quite as good. Instead of constant time, a `TreeMap` takes time proportional to $\log n$.
2. Inside the `TreeMap`, the keys are not stored in order; they are stored in a **binary search tree**, which makes it possible to traverse the keys, in order, in linear time.

In this exercise, we'll explain how binary search trees work and then you will use one to implement a `Map`. Along the way, we'll analyze the performance of the core map methods when implemented using a tree.

17.3 Binary search tree

A binary search tree (BST) is a tree where each node contains a key, and every node has the “BST property”:

1. If `node` has a left child, all keys in the left subtree must be less than the key in `node`.
2. If `node` has a right child, all keys in the right subtree must be greater than the key in `node`.

The following diagram shows a tree of integers that has this property.

This figure is from the https://en.wikipedia.org/wiki/Binary_search_tree Wikipedia page on binary search trees, which you might find useful while you work on this exercise.

The key in the root is 8, and you can confirm that all keys to the left of the root are less than 8, and all keys to the right are greater. You can also check that the other nodes have this property.

Looking up a key in a binary search tree is fast because we don't have to search the entire tree. Starting at the root, we can use the following algorithm:

1. Compare the key you are looking for, `target`, to the key in the current node. If they are equal, you are done.

2. If **target** is less than the current key, search the left tree. If there isn't one, **target** is not in the tree.
3. If **target** is greater than the current key, search the right tree. If there isn't one, **target** is not in the tree.

At each level of the tree, you only have to search one child. For example, if you look for **target** = 4 in the previous diagram, you start at the root, which contains the key 8. Because **target** is less than 8, you go left. Because **target** is greater than 3 you go right. Because **target** is less than 6, you go left. And then you find the key you are looking for.

In this example, it takes 4 comparisons to find the target, even though the tree contains 9 keys. In general, the number of comparisons is proportional to the height of the tree, not the number of keys in the tree.

So what can we say about the relationship between the height of the tree, h , and the number of nodes, n ? Starting small and working up:

- If $h=1$, the tree only contains one node, so $n=1$.
- If $h=2$, we can add two more nodes, for a total of $n=3$.
- If $h=3$, we can add up to 4 more nodes, for a total of $n=7$.
- If $h=4$, we can add up to 8 more nodes, for a total of $n=15$.

By now you might see the pattern. If we number the levels tree from 1 to h , the level with index i can have up to 2^{i-1} nodes. And the total number of nodes in h levels is $2^h - 1$. If we have

$$n = 2^h - 1$$

We can take the log base 2 of both sides:

$$\log_2 n = h$$

Which means that the height of the tree is proportional to $\log n$, if the tree is full; that is, if each level contains the maximum number of nodes.

So we expect that we can look up a key in a binary search tree in time proportional to $\log n$. This is true if the tree is full, and even if the tree is only partially full. But it is not always true, as we will see.

An algorithm that takes time proportional to $\log n$ is called “logarithmic” or “log time”, and it belongs to the order of growth $O(\log n)$.

17.4 Implementing a tree-backed Map

In this section we'll describe the starter code for this exercise; in the next section we'll give you instructions for filling in the missing methods.

Here's the beginning of our implementation, called `MyTreeMap`:

```
public class MyTreeMap<K, V> implements Map<K, V> {
```

```
    private int size = 0;
    private Node root = null;
```

The instance variables are `size`, which keeps track of the number of keys, and `root`, which is a reference to the root node in the tree. When the tree is empty, `root` is `null` and `size` is 0.

Here's the definition of `Node`, which is defined inside `MyTreeMap`:

```
protected class Node {
    public K key;
    public V value;
    public Node left = null;
    public Node right = null;

    public Node(K key, V value) {
        this.key = key;
        this.value = value;
    }
}
```

Each node contains a key-value pair and references to two child nodes, `left` and `right`. Either or both of the child nodes can be `null`.

Some of the `Map` methods are easy to implement, like `size` and `clear`:

```
public int size() {
    return size;
}

public void clear() {
    size = 0;
    root = null;
}
```

`size` is clearly constant time.

`clear` appears to be constant time, but consider this: when `root` is set to `null`, the garbage collector reclaims the nodes in the tree, which takes linear time. Should work done by the garbage collector count? We think so.

In the next section, you'll fill in some of the other methods, including the most important ones, `get` and `put`.

17.5 Instructions

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory with the code for this exercise, `javacs-lab09`.

In the subdirectory `javacs-lab09/src/com/flatiron-school/javacs` you'll find these source files:

- `MyTreeMap.java` contains the code from the previous section with outlines for the missing methods.
- `MyTreeMapTest.java` contains the unit tests for `MyTreeMap`.

Also, in `javacs-lab09`, you'll find the Ant build file `build.xml`.

- In `javacs-lab09`, run `ant build` to compile the source files. Then run `ant test`, which runs `MyTreeMapTest`. Several tests should fail, because you have some work to do!
- We've provided outlines for `get` and `containsKey`. Both of them use `findNode`, which is a private method we defined; it is not part of the `Map` interface. Here's how it starts:

```
private Node findNode(Object target) {
    if (target == null) {
        throw new NullPointerException();
    }
}
```

```
    @SuppressWarnings("unchecked")
    Comparable<? super K> k = (Comparable<? super K>) target;
    int cmp = k.compareTo(p.key);
    // TODO: Fill this in.
    return null;
}
```

The parameter `target` is the key we’re looking for. If `target` is `null`, `findNode` throws an exception. Some implementations of `Map` can handle `null` as a key, but in a binary search tree, we need to be able to compare keys, so dealing with `null` is problematic. To keep things simple, we decided that our implementation should not allow `null` as a key.

The next lines show how we can compare `target` to a key in the tree. From the signature of `get` and `containsKey`, the compiler considers `target` to be an `Object`. But we need to be able to compare keys, so we typecast `target` to `Comparable<? super K>`, which means that it is comparable to an instance of type `K`, or any superclass of `K`.

If you are not familiar with this use of “type wildcards”, <http://docs.oracle.com/javase/tutorial/extra/generics/morefun.html> you can read more here.

Fortunately, dealing with Java’s type system is not the point of this exercise. Your job is to fill in the rest of `findNode`. If it finds a node that contains `target` as a key, it should return the node. Otherwise it should return `null`. When you get this working, the tests for `get` and `containsKey` should pass.

Note that your solution should only search one path through the tree, so it should take time proportional to the height of the tree. You should not search the whole tree!

- Your next task is to fill in `containsValue`. To get you started, we’ve provided a helper method, `equals`, that compares `target` and a given key. Note that the values in the tree (as opposed to the keys) are not necessarily comparable, so we can’t use `compareTo`; we have to invoke `equals` on `target`.

Unlike your previous solution for `findNode`, your solution for `containsValue` *does* have to search the whole tree, so its runtime will be proportional to the number of keys, `n`, not the height of the tree, `h`.

- The next method you should fill in is `put`. We've provided starter code that handles the simple cases:

```
public V put(K key, V value) {
    if (key == null) {
        throw new NullPointerException();
    }
    if (root == null) {
        root = new Node(key, value);
        size++;
        return null;
    }
    return putHelper(root, key, value);
}

private V putHelper(Node node, K key, V value) {
    // TODO: Fill this in.
}
```

If you try to put `null` as a key, `put` throws an exception.

If the tree is empty, `put` creates a new node and initializes the instance variable `root`.

Otherwise, it calls `putHelper`, which is a private method we defined; it is not part of the `Map` interface.

Fill in `putHelper` so it searches the tree and:

1. If `key` is already in the tree, it replaces the old value with the new, and returns the old value.
2. If `key` is not in the tree, it creates a new node, finds the right place to add it, and returns `null`.

Your implementation of `put` should take time proportional to `h`, not `n`. Ideally you should search the tree only once, but if you find it easier to search twice, you can do that; it will be slower, but it doesn't change the order of growth.

- Finally, you should fill in the body of `keySet`. [https://docs.oracle.com/javase/7/docs/api/java/util/TreeMap.html#keySet\(\)](https://docs.oracle.com/javase/7/docs/api/java/util/TreeMap.html#keySet()) According to the documentation, this method should return some kind of `Set` that iterates the keys in order; that is, in increasing order according to the `compareTo` method. There are several ways to do that, but the simplest is to return an implementation of `Set` that maintains the order of the elements. The most commonly-used implementations, like `HashSet`, don't have this property, but Java provides one that does: `LinkedHashSet`.

We provide an outline of `keySet` that creates and returns a `LinkedHashSet`:

```
public Set<K> keySet() {  
    Set<K> set = new LinkedHashSet<K>();  
    return set;  
}
```

You should finish off this method so it adds the keys from the tree to `set` in ascending order. Hint: you might want to write a helper method, you might want to make it recursive, and you might want to https://en.wikipedia.org/wiki/Tree_traversal#In-order read about in-order tree traversal.

When you are done, all tests should pass. In the next README, we'll go over our solutions and test the performance of the core methods.

17.6 Resources

https://en.wikipedia.org/wiki/Binary_search_tree Binary search tree: Wikipedia.

<http://docs.oracle.com/javase/tutorial/extra/generics/morefun.html> Type wildcards: Java tutorial.

https://en.wikipedia.org/wiki/Tree_traversal Tree traversal: Wikipedia.

Chapter 18

Binary search tree

18.1 Learning goals

1. Review a `Map` implementation using a binary search tree (BST).
2. Analyze the performance of BST methods.
3. Explain the problem of unbalanced trees and how self-balancing trees solve this problem.

18.2 Overview

This README presents solutions to the previous exercise, then tests the performance of our tree-backed map. We present a problem with our implementation and explain how Java's `TreeMap` solves this problem.

18.3 Our version of `MyTreeMap`

In the previous exercise we gave you the outline of `MyTreeMap` and asked you to fill in the missing methods. Now we'll present our solutions, starting with `findNode`:

```
private Node findNode(Object target) {
    // some implementations can handle null as a key, but not this one
    if (target == null) {
        throw new NullPointerException();
    }

    // something to make the compiler happy
    @SuppressWarnings("unchecked")
    Comparable<? super K> k = (Comparable<? super K>) target;

    // the actual search
    Node node = root;
    while (node != null) {
        int cmp = k.compareTo(node.key);
        if (cmp < 0)
            node = node.left;
        else if (cmp > 0)
            node = node.right;
        else
            return node;
    }
    return null;
}
```

`findNode` is a private method used by `containsKey` and `get`; it is not part of the `Map` interface. The parameter `target` is the key we're looking for. We explained the first part of this function in the previous exercise:

- In this implementation, `null` is not a legal value for a key.
- Before we can invoke `compareTo` on `target`, we have to typecast it to some kind of `Comparable`. The “type wildcard” used here is as permissive as possible; that is, it works with any type that implements `Comparable` and whose `compareTo` method accepts `K` or any supertype of `K`.

After all that, the actual search is relatively simple. We initialize a loop variable `node` so it refers to the root node. Each time through the loop, we compare the target to `node.key`. If the target is less than the current key, we

move to the left child. If it's greater, we move to the right child. And if it's equal, we return the current node.

If we get to the bottom of the tree without finding the target, we conclude that it is not in the tree and return `null`.

18.4 Searching for values

As we explained in the previous exercise, the runtime of `findNode` is proportional to the height of the tree, not the number of nodes, because we don't have to search the whole tree. But for `containsValue`, we have to search the values, not the keys; the BST property doesn't apply to the values, so we have to search the whole tree.

Our solution is recursive:

```
public boolean containsValue(Object target) {
    return containsValueHelper(root, target);
}

private boolean containsValueHelper(Node node, Object target) {
    if (node == null) {
        return false;
    }
    if (equals(target, node.value)) {
        return true;
    }
    if (containsValueHelper(node.left, target)) {
        return true;
    }
    if (containsValueHelper(node.right, target)) {
        return true;
    }
    return false;
}
```

`containsValue` takes the target value as a parameter and immediately invokes `containsValueHelper`, passing the root of the tree as an additional parameter.

Here's how `containsValueHelper` works

- The first `if` statement checks the base case of the recursion. If `node` is `null`, that means we have recursed to the bottom of the tree without finding the `target`, so we should return `false`. Note that this only means that the target did not appear on one path through the tree; it is still possible that it will be found on another.
- The second case checks whether we've found what we're looking for. If so, we return `true`. Otherwise, we have to keep going.
- The third case makes a recursive call to search for `target` in the left subtree. If we find it, we can return `true` immediately, without searching the right subtree. Otherwise, we keep going.
- The fourth case searches the right subtree. Again, if we find what we are looking for, we return `true`. Otherwise, having searched the whole tree, we return `false`.

This method “visits” every node in the tree, so it takes time proportional to the number of nodes.

18.5 Implementing `put`

The `put` method is a little more complicated than `get` because it has to deal with two cases: (1) If the given key is already in the tree, it replaces and returns the old value; (2) otherwise it has to add a new node to the tree, in the right place.

In the previous exercise, we provided this starter code:

```
public V put(K key, V value) {
    if (key == null) {
        throw new NullPointerException();
    }
    if (root == null) {
        root = new Node(key, value);
        size++;
        return null;
    }
    return putHelper(root, key, value);
}
```

And we asked you to fill in `putHelper`. Here's our solution:

```
private V putHelper(Node node, K key, V value) {
    Comparable<? super K> k = (Comparable<? super K>) key;
    int cmp = k.compareTo(node.key);

    if (cmp < 0) {
        if (node.left == null) {
            node.left = new Node(key, value);
            size++;
            return null;
        } else {
            return putHelper(node.left, key, value);
        }
    }
    if (cmp > 0) {
        if (node.right == null) {
            node.right = new Node(key, value);
            size++;
            return null;
        } else {
            return putHelper(node.right, key, value);
        }
    }
    V oldValue = node.value;
    node.value = value;
    return oldValue;
}
```

The first parameter, `node`, is initially the root of the tree, but each time we make a recursive call, it refers to a different subtree.

As in `get`, we use the `compareTo` method to figure out which path to follow through the tree.

If `cmp < 0`, the key we're adding is less than `node.key`, so we want to look in the left subtree. There are two cases:

- If the left subtree is empty; that is, if `node.left` is `null`, we have reached the bottom of the tree without finding `key`. At this point, we know that

`key` isn't in the tree, and we know where it should go. So we create a new node and add it as the left child of `node`.

- Otherwise we make a recursive call to search the left subtree.

If `cmp > 0`, the key we're adding is greater than `node.key`, so we want to look in the right subtree. And we handle the same two cases as in the previous branch.

Finally, if `cmp == 0`, we found the key in the tree, so we replace and return the old value.

We wrote this method recursively to make it more readable, but it would be straightforward to re-write it iteratively, which you might want to do as an exercise.

18.6 In-order traversal

The last method we asked you to write is `keySet`, which returns a `Set` that contains the keys from the tree in ascending order. In other implementations of `Map`, the keys returned by `keySet` are in no particular order, but one of the capabilities of the tree implementation is that it is simple and efficient to sort the keys. So we should take advantage of that.

Here's our solution:

```
public Set<K> keySet() {
    Set<K> set = new LinkedHashSet<K>();
    addInOrder(root, set);
    return set;
}

private void addInOrder(Node node, Set<K> set) {
    if (node == null) return;
    addInOrder(node.left, set);
    set.add(node.key);
    addInOrder(node.right, set);
}
```


In `keySet`, we create a `LinkedHashSet`, which is a `Set` implementation that keeps the elements in order (unlike most other `Set` implementations). Then we call `addInOrder` to traverse the tree.

The first parameter, `node`, is initially the root of the tree, but as you should expect by now, we use it to traverse the tree recursively. `addInOrder` performs a classic “in-order traversal” of the tree.

If `node` is `null`, that means the subtree is empty, so we return without adding anything to `set`. Otherwise we:

1. Traverse the left subtree.
2. Add `node.key`.
3. Traverse the right subtree.

Remember that the BST property guarantees that all nodes in the left subtree are less than `node.key`, and all nodes in the right subtree are greater. So we know that `node.key` has been added in the correct order.

Applying the same argument recursively, we know that the elements from the left subtree are in order, as well as the elements from the right subtree.

And the base case is correct: if the subtree is empty, no keys are added.

So we can conclude that this method adds all keys in the correct order.

Because this method visits every node in the tree, like `containsValue`, it takes time proportional to n .

18.7 The logarithmic methods

In `MyTreeMap`, the methods `get` and `put` take time proportional to the height of the tree, h . In the previous lab, we showed that if the tree is full — if every level of the tree contains the maximum number of nodes — the height of the tree is proportional to $\log n$.

And we claimed that `get` and `put` are logarithmic; that is, they take time proportional to $\log n$. But for most applications, there’s no guarantee that

the tree is full. In general, the shape of the tree depends on the keys and what order they are added.

To see how this works out in practice, we'll test our implementation with two sample datasets: a list of random strings and a list of timestamps in increasing order.

Here's the code that generates random strings:

```
Map<String, Integer> map = new MyTreeMap<String, Integer>();

for (int i=0; i<n; i++) {
    String uuid = UUID.randomUUID().toString();
    map.put(uuid, 0);
}
```

UUID is a class in the `java.util` package that can generate a random “universally unique identifier”. UUIDs are useful for a variety of applications, but in this example we're taking advantage of an easy way to generate random strings.

We ran this code with `n=16384` and measured the runtime and the height of the final tree. Here's the output:

```
Time in milliseconds = 151
Final size of MyTreeMap = 16384
log base 2 of size of MyTreeMap = 14.0
Final height of MyTreeMap = 33
```

We computed “log base 2 of size of MyTreeMap” to see how tall the tree would be if it were full. The result indicates that a full tree with height 14 would contain 16,384 nodes.

The actual tree of random strings has height 33, which is substantially more than the theoretical minimum, but not too bad. To find one key in a collection of 16,384, we only have to make 33 comparisons. Compared to a linear search, that's almost 500 times faster.

This performance is typical with random strings or other keys that are added in no particular order. The final height of the tree might be 2-3 times the theoretical minimum, but it is still proportional to $\log n$, which is much less

than n . In fact, $\log n$ grows so slowly as n increases, it can be difficult to distinguish logarithmic time from constant time in practice.

However, binary search trees don't always behave so well. Let's see what happens when we add keys in increasing order. Here's an example that measures timestamps in nanoseconds and uses them as keys:

```
MyTreeMap<String, Integer> map = new MyTreeMap<String, Integer>();

for (int i=0; i<n; i++) {
    String timestamp = Long.toString(System.nanoTime());
    map.put(timestamp, 0);
}
```

`System.nanoTime` returns an integer with type `long` that indicates elapsed time in nanoseconds. Each time we call it, we get a somewhat bigger number. When we convert these timestamps to strings, they appear in increasing alphabetical order.

And let's see what happens when we run it:

```
Time in milliseconds = 1158
Final size of MyTreeMap = 16384
log base 2 of size of MyTreeMap = 14.0
Final height of MyTreeMap = 16384
```

The runtime is longer than in the previous case, more than 7 times longer. If you wonder why, take a look at the final height of the tree: 16384!

If you think about how `put` works, you can figure out what's going on. Every time we add a new key, it's larger than all of the keys in the tree, so we always choose the right subtree, and always add the new node as the right child of the rightmost node. The result is an "unbalanced" tree that only contains right children.

The height of this tree is proportional to n , not $\log n$, so the performance of `get` and `put` is linear, not logarithmic.

18.8 Self-balancing trees

There are two possible solutions to this problem:

- You could avoid adding keys to the `Map` in order. But this is not always possible.
- You could make a tree that does a better job of handling keys if they happen to be in order.

The second solution is better, and there are several ways to do it. The most common is to modify `put` so that it detects when the tree is starting to become unbalanced and, if so, rearranges the nodes. Trees with this capability are called “self-balancing”. Common self-balancing trees include the AVL tree (“AVL” are the initials of the inventors), and the red-black tree, which is what the Java `TreeMap` uses.

In our example code, if we replace `MyTreeMap` with the Java `TreeMap`, the runtimes are about the same for the random strings and the timestamps. In fact, the timestamps run faster, even though they are in order, probably because they take less time to hash.

In summary, a binary search tree can implement `get` and `put` in logarithmic time, but only if the keys are added in an order that keeps the tree sufficiently balanced. Self-balancing trees avoid this problem by doing some additional work each time a new key is added.

18.9 One more exercise

In the previous exercise we decided not to implement `remove`, but you might want to try it. If you remove a node from the middle of the tree, you have to rearrange the remaining nodes to restore the BST property. You can probably figure out how to do that on your own, or you can https://en.wikipedia.org/wiki/Binary_search_tree#Deletion read the explanation here.

Removing a node and rebalancing a tree are similar operations: if you do this exercise, you will have a better idea of how self-balancing trees work.

18.10 Resources

- <https://docs.oracle.com/javase/7/docs/api/java/util/TreeMap.html> The documentation of `TreeMap` provides more information about the implementation.
- You can https://en.wikipedia.org/wiki/Red%E2%80%93black_tree read more about red-black trees here.

DRAFT: Not ready for distribution!

DRAFT: Not ready for distribution!

Chapter 19

Redis

19.1 Learning goals

1. Read and write data to a Redis database.
2. Translate Java data structures to Redis data structures.
3. Implement a Redis-backed Web search index.

19.2 Overview

In the next few exercises we will get back to the application we started in the first unit: building a web search engine. To review, the components of a search engine are:

- Crawling: We'll need a program that can download a web page, parse it, and extract the text and any links to other pages.
- Indexing: We'll need an index that makes it possible to look up a search term and find the pages that contain it.
- Retrieval: And we'll need a way to collect results from the Index and identify pages that are most relevant to the search terms.

If you did the exercise on indexing, you implemented an index using Java maps. In this exercise, we'll revisit the Indexer and make a new version that stores the results in a database.

If you did the exercise about the “Getting to Philosophy” conjecture, you built a crawler that follows the first link it finds. In the next exercise, we'll make a more general version that stores every link it finds in a queue and explores them in order.

And then, finally, you will work on the retrieval problem.

In these exercises, we will provide less starter code, and you will make more design decisions. These exercises are also more open-ended. We will suggest some minimal goals you should try to reach, but there are many ways you can go farther if you want to challenge yourself.

Now, let's get started on a new version of the Indexer.

19.3 Persistence

The previous version of the Indexer stores the index in two data structures: a **TermCounter** that maps from a search term to the number of times it appears on a web page, and an **Index** that maps from a search term to the set of pages where it appears.

These data structures are stored in the memory of a running Java program, which means that when the program stops running, the index is lost. Data stored only in the memory of a running program is called “volatile”, because it vaporizes when the program ends.

Data that persists after the program that created it ends is called “persistent”. In general, files stored in a file system are persistent, as well as data stored in databases.

A simple way to make data persistent is to store it in a file. Before the program ends, it could translate its data structures into a format like <https://en.wikipedia.org/wiki/JSON> and then write them into a file. When it starts again, it could read the file and rebuild the data structures.

But there are several problems with this solution:

1. Reading and writing large data structures (like a Web index) would be slow.
2. The entire data structure might not fit into the memory of a single running program.
3. If a program ends unexpectedly (for example, due to a power outage), any changes made since the program last started would be lost.

A better alternative is a database that provides persistent storage and the ability to read and write parts of the database without reading and writing the whole thing.

There are many kinds of database management systems (DBMS) that provide different capabilities. You can read an <https://en.wikipedia.org/wiki/Databaseoverview> on this Wikipedia page.

The database we recommend for this exercise is Redis, which provides persistent data structures that are similar to Java data structures. Specifically, it provides:

- Lists of strings, similar to Java Lists.
- Hashes, similar to Java Maps.
- Sets of strings, similar Java Sets.

Redis is a “key-value database”, which means that the data structures it contains (the values) are identified by unique strings (the keys). A key in Redis plays the same role as a reference in Java: it identifies an object. We’ll see some examples soon.

19.4 Redis clients and servers

Redis is usually run as a remote service; in fact, the name stands for “REmote DIctionary Server”. To use Redis, you have to run the Redis server somewhere and then connect to it using a Redis client. There are many ways to set up a server and many clients you could use. For this lab, we recommend:

1. Rather than install and run the server yourself, consider using a service like <https://redistogo.com/> RedisToGo, which runs Redis in the cloud. They offer a free plan with enough resources for this exercise.

2. For the client we recommend Jedis, which is a Java library that provides classes and methods for working with Redis.

Here are more detailed instructions to help you get started:

- <https://redistogo.com/signup> Create an account on RedisToGo and select the plan you want (probably the free plan to get started).
- Create an “instance”, which is a virtual machine running the Redis server. If you click on the “Instances” tab, you should see your new instance, identified by a host name and a port number. For example, we have an instance named “dory-10534”.
- Click on the instance name to get the configuration page. Make a note of the URL near the top of the page, which looks like this:
`redis://redistogo:1234567890feedfacebeefa1e1234567@dory.redistogo.com:10534`

This URL contains the server’s host name, `dory.redistogo.com`, the port number, `10534`, and the password you will need to connect to the server, which is the long string of letters and numbers in the middle. You will need this information for the next step.

19.5 Hello, Jedis

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory with the code for this exercise, `javacs-lab10`.

In the subdirectory `javacs-lab10/src/com/flatiron/school/javacs` you’ll find the source files for this exercise:

- `JedisMaker.java` contains example code for connecting to a Redis server and running a few Jedis methods.
- `JedisIndex.java` contains starter code for this exercise.
- `JedisIndexTest.java` contains test code for `JedisIndex`.

- `WikiFetcher.java` contains the code we saw in previous exercises to read web pages and parse them using JSoup.

You'll also find these files, which are part of our solution to previous labs.

- `Index.java` implements an index using Java data structures.
- `TermCounter.java` represents a map from terms to their frequencies.
- `WikiNodeIterable.java` iterates through the nodes in a DOM tree produced by JSoup.

Also, in `javacs-lab10`, you'll find the Ant build file `build.xml`.

The first step is to use Jedis to connect to your Redis server. `RedisMaker.java` shows how to do this. It reads information about your Redis server from a file, connects to it and logs in using your password, then returns a `Jedis` object you can use to perform Redis operations.

If you open `JedisMaker.java`, you should see something like this (but with more error checking):

```
public class JedisMaker {

    public static Jedis make() {
        // assemble the directory name
        String slash = File.separator;
        String filename = System.getProperty("user.dir") + slash +
            "src" + slash + "resources" + slash + "redis_url.txt";

        // read the contents of the file into a string
        StringBuilder sb = new StringBuilder();
        BufferedReader br = new BufferedReader(new FileReader(filename));
        while (true) {
            String line = br.readLine();
            if (line == null) break;
            sb.append(line);
        }
        br.close();
    }
}
```

```
// extract the components from the URI
URI uri = new URI(sb.toString());
String host = uri.getHost();
int port = uri.getPort();
String[] array = uri.getAuthority().split("[:@]");
String auth = array[1];

// connect to the server and authenticate
Jedis jedis = new Jedis(host, port);
jedis.auth(auth);

return jedis;
}
```

JedisMaker is a helper class that provides one static method, **make**, which creates a **Jedis** object. Once this object is authenticated (by invoking **auth**), you can use it to communicate with your Redis database.

JedisMaker reads information about your Redis server from a file named **redis_url.txt**, which you should put in the directory **javacs-lab10/src/resources**:

- Use a text editor to create and edit **javacs-lab10/src/resources/redis_url.txt**.
- Paste in the URL of your server. If you are using RedisToGo, the URL will look like this:

```
redis://redistogo:1234567890feedfacebeefa1e1234567@dory.redistogo.com:10534
```

Because this file contains the password for your Redis server, you should not put this file in a public repository. To help you avoid doing that by accident, we have provided a **.gitignore** file with this lab that should make it harder (but not impossible) to put this file in your repo.

Now in **javacs-lab10**, run **ant build** to compile the source files and **ant JedisMaker** to run the example code in **main**:

```
public static void main(String[] args) {

    Jedis jedis = make();
```

```
// String
jedis.set("mykey", "myvalue");
String value = jedis.get("mykey");
System.out.println("Got value: " + value);

// Set
jedis.sadd("myset", "element1", "element2", "element3");
System.out.println("element2 is member: " + jedis.sismember("myset", "element2"));

// List
jedis.rpush("mylist", "element1", "element2", "element3");
System.out.println("element at index 1: " + jedis.lindex("mylist", 1));

// Hash
jedis.hset("myhash", "word1", Integer.toString(2));
jedis.hincrBy("myhash", "word2", 1);
System.out.println("frequency of word1: " + jedis.hget("myhash", "word1"));
System.out.println("frequency of word2: " + jedis.hget("myhash", "word2"));

jedis.close();
}
```

This example demonstrates the data types and methods you are most likely to use for this exercise. When you run it, the output should be:

```
Got value: myvalue
element2 is member: true
element at index 1: element2
frequency of word1: 2
frequency of word2: 1
```

In the next section, we'll explain how the code works.

19.6 Redis data types

Redis is basically a map from keys, which are Strings, to values, which can be one of several data types. The most basic Redis data type is a String. To add a String to the database, use `jedis.set`, which is similar to `Map.put`; the

parameters are the new key and the corresponding value. To look up a key and get its value, use `jedis.get`:

```
jedis.set("mykey", "myvalue");  
String value = jedis.get("mykey");
```

In this example, the key is **"mykey"** and the value is **"myvalue"**.

A Redis **Set** is similar to a Java `Set<String>`. To add elements to a set, you have to choose a key to identify the set and then use `jedis.sadd`:

```
jedis.sadd("myset", "element1", "element2", "element3");  
boolean flag = jedis.sismember("myset", "element2");
```

You don't have to create the Set as a separate step. If it doesn't exist, Redis creates it. In this case, it creates a Set named `myset` that contains three elements.

The method `jedis.sismember` checks whether an element is in a Set. Adding elements and checking membership are constant time operations.

A Redis **List** is similar to a Java `List<String>`. The method `jedis.rpush` adds elements to the end (right side) of a List:

```
jedis.rpush("mylist", "element1", "element2", "element3");  
String element = jedis.lindex("mylist", 1);
```

Again, you don't have to create the data structure before you start adding elements. This example creates a List named `"mylist"` that contains three elements.

The method `jedis.lindex` takes an integer index and returns the indicated element. Adding and accessing elements are constant time operations.

Finally, a Redis **Hash** is similar to a Java `Map<String, String>`. The method `jedis.hset` adds a new entry to the hash:

```
jedis.hset("myhash", "word1", Integer.toString(2));  
String value = jedis.hget("myhash", "word1");
```

This example creates a Hash named `myhash` that contains one entry.

The value it stores is an `Integer`, so we have to convert it to a `String`. When we look up the value using `jedis.hget`, the result is a `String`.

Working with Redis hashes can be confusing, because we use a key to identify which hash we want, and then another key to identify a value in the hash. In the context of Redis, the second key is called a “field”, which might help keep things straight. So a “key” like `myhash` identifies a particular hash, and then a “field” like `word1` identifies a value in the hash.

For many applications, the values in a Redis hash are integers, so Redis provides a few special methods, like `hincrby`, that treat the values as integers:

```
jedis.hincrBy("myhash", "word2", 1);
```

This method accesses `myhash`, gets the current value associated with `word2` (or 0 if it doesn’t already exist), increments it by 1, and writes the result back to the hash.

Setting, getting, and incrementing entries in a hash are constant time operations.

You can <http://redis.io/topics/data-types> read more about Redis data types here.

19.7 Making a Jedis-backed index

At this point you have the information you need to make a Web search index that stores results in a Redis database.

In `javacs-lab10`, run `ant build` to compile the source files and `ant test` to run `JedisIndexTest`. It should fail, because you have some work to do!

`JedisIndexTest` tests these methods:

- `JedisIndex`, which is the constructor that takes a `Jedis` object as a parameter.

- `indexPage`, which adds a Web page to the index; it takes a `String` URL and a `JSoup Elements` object that contains the elements of the page that should be indexed.
- `getCounts`, which takes a search term and returns a `Map<String, Integer>` that maps from each URL that contains the search term to the number of times it appears on that page.

Here's an example of how these methods are used:

```
WikiFetcher wf = new WikiFetcher();
String url1 = "https://en.wikipedia.org/wiki/Java_(programming_language)";
Elements paragraphs = wf.readWikipedia(url1);

Jedis jedis = JedisMaker.make();
JedisIndex index = new JedisIndex(jedis);
index.indexPage(url1, paragraphs);
Map<String, Integer> map = index.getCounts("the");
```

If we look up `url1` in the result, `map`, we should get 339, which is the number of times the word “the” appears on [https://en.wikipedia.org/wiki/Java_\(programming_language\)](https://en.wikipedia.org/wiki/Java_(programming_language)) the Java Wikipedia page (that is, the version we saved).

If we index the same page again, the new results should replace the old ones.

You might want to start with a copy of `Index.java`, which contains our solution to the previous exercise where we built an index using Java data structures.

One suggestion for translating data structures from Java to Redis: remember that each object in a Redis database is identified by a unique key, which is a string. If you have two kinds of objects in the same database, you might want to add a prefix to the keys to distinguish between them. For example, in our solution, we have two kinds of objects:

- We define a `URLSet` to be a Redis `Set` that contains the URLs that contain a given search term. The key for each `URLSet` starts with `"URLSet:"`, so to get the URLs that contain the word “the”, we access the `Set` with the key `"URLSet:the"`.

- We define a `TermCounter` to be a Redis `Hash` that maps from each term that appears on a page to the number of times it appears. The key for each `TermCounter` starts with `"TermCounter:"` and ends with the URL of the page we're looking up.

In our implementation, we have one `URLSet` for each term and one `TermCounter` for each indexed page. We provide two helper methods, `urlSetKey` and `termCounterKey`, to assemble these keys.

19.8 More suggestions if you want them

At this point you have all the information you need to do the exercise, so you can get started if you are ready. But we have a few suggestions you might want to read first:

- For this exercise we provide less guidance than we did in previous exercises. You will have to make some design decisions; in particular, you will have to figure out how to divide the problem into pieces that you can test one at a time, and then assemble the pieces into a complete solution. If you try to write the whole thing at once, without testing smaller pieces, it might take a very long time to debug.
- One of the challenges of working with persistent data is that it is persistent. The structures stored in the database might change every time you run the program. If you mess something up in the database, you will have to fix it or start over before you can proceed. To help you keep things under control, we've provided methods called `deleteURLSets`, `deleteTermCounters`, and `deleteAllKeys`, which you can use to clean out the database and start fresh. You can also use `printIndex` to print the contents of the index.
- Each time you invoke a `Jedis` method, your client sends a message to the server, then the server performs the action you requested and sends back a message. If you perform many small operations, it will probably take a long time. You can improve performance by grouping a series of operations into a `Transaction`.

For example, here's a simple version of `deleteAllKeys`:

```
public void deleteAllKeys() {
    Set<String> keys = jedis.keys("*");
    for (String key: keys) {
        jedis.del(key);
    }
}
```

Each time you invoke `del` requires a round-trip from the client to the server and back. If the index contains more than a few pages, this method would take a long time to run. We can speed it up with a `Transaction` object:

```
public void deleteAllKeys() {
    Set<String> keys = jedis.keys("*");
    Transaction t = jedis.multi();
    for (String key: keys) {
        t.del(key);
    }
    t.exec();
}
```

`jedis.multi` returns a `Transaction` object, which provides all the methods of a `Jedis` object. But when you invoke a method on a `Transaction`, it doesn't run the operation immediately, and it doesn't communicate with the server. It saves up a batch of operations until you invoke `exec`. Then it sends all of the saved operations to the server at the same time, which is usually much faster.

19.9 A few design hints

Now you *really* have all the information you need; you should start working on the exercise. But if you get stuck, or if you really don't know how to get started, you can come back for a few more hints.

Don't read the following until you have run the test code, tried out some basic Redis commands, and written a few methods in `JedisIndex.java`.

Ok, if you are really stuck, here are some methods you might want to work on:

```
/**
 * Adds a URL to the set associated with 'term'.
 */
public void add(String term, TermCounter tc) {}

/**
 * Looks up a search term and returns a set of URLs.
 */
public Set<String> getURLs(String term) {}

/**
 * Returns the number of times the given term appears at the given URL.
 */
public Integer getCount(String url, String term) {}

/**
 * Pushes the contents of the TermCounter to Redis.
 */
public List<Object> pushTermCounterToRedis(TermCounter tc) {}
```

These are the methods we used in our solutions, but they are certainly not the only way to divide things up. So please take these suggestions if they help, but ignore them if they don't.

For each method, consider writing the tests first. When you figure out how to test a method, you often get ideas about how to write it.

Good luck!

19.10 Resources

- <https://en.wikipedia.org/wiki/JSON>JSON: Wikipedia.
- <https://en.wikipedia.org/wiki/Database>Database management systems: Wikipedia.
- <http://www.computerhope.com/issues/ch000549.htm>Environment variables in Windows: tutorial.
- <http://redis.io/topics/data-types>Redis data types: documentation.

DRAFT: Not ready for distribution!

Chapter 20

Crawling Wikipedia

20.1 Learning goals

1. Analyze the performance of Web indexing algorithms.
2. Implement a Web crawler.

20.2 Overview

In this exercise, we present our solution to the previous exercise and analyze the performance of Web indexing algorithms. Then we build a simple Web crawler.

20.3 Our Redis-backed indexer

In our solution, we store two kinds of structures in Redis:

- For each search term, we have a **URLSet**, which is a Redis Set of URLs that contain the search term.
- For each URL, we have a **TermCounter**, which is a Redis Hash that maps each search term to the number of times it appears.

We discussed these data types in the previous exercise. You can also <http://redis.io/topics/data-types> read about Redis Sets and Hashes here.

In `JedisIndex`, we provide a function that takes a search term and returns the Redis key of its `URLSet`:

```
private String urlSetKey(String term) {  
    return "URLSet:" + term;  
}
```

And a function that takes a URL and returns the Redis key of its `TermCounter`:

```
private String termCounterKey(String url) {  
    return "TermCounter:" + url;  
}
```

Here's our implementation of `indexPage`, which takes a URL and a `JSoup Elements` object that contains the DOM tree of the paragraphs we want to index:

```
public void indexPage(String url, Elements paragraphs) {  
    System.out.println("Indexing " + url);  
  
    // make a TermCounter and count the terms in the paragraphs  
    TermCounter tc = new TermCounter(url);  
    tc.processElements(paragraphs);  
  
    // push the contents of the TermCounter to Redis  
    pushTermCounterToRedis(tc);  
}
```

To index a page, we

1. Make a Java `TermCounter` for the contents of the page, using code from a previous exercise.
2. Push the contents of the `TermCounter` to Redis.

Here's the new code that pushes a `TermCounter` to Redis:

```
public List<Object> pushTermCounterToRedis(TermCounter tc) {
    Transaction t = jedis.multi();

    String url = tc.getLabel();
    String hashname = termCounterKey(url);

    // if this page has already been indexed; delete the old hash
    t.del(hashname);

    // for each term, add an entry in the termcounter and a new
    // member of the index
    for (String term: tc.keySet()) {
        Integer count = tc.get(term);
        t.hset(hashname, term, count.toString());
        t.sadd(urlSetKey(term), url);
    }
    List<Object> res = t.exec();
    return res;
}
```

This method uses a `Transaction` to collect the operations and send them to the server all at once, which is much faster than sending a series of small operations.

It loops through the terms in the `TermCounter`. For each one it

1. Finds or creates a `TermCounter` on Redis, then adds a field for the new term.
2. Finds or creates a `URLSet` on Redis, then adds the current URL.

If the page has already been indexed, we delete its old `TermCounter` before pushing the new contents.

That's it for indexing new pages.

The second part of the exercise asked you to write `getCounts`, which takes a search term and returns a map from each URL where the term appears to the number of times it appears there. Here is our solution:

```
public Map<String, Integer> getCounts(String term) {
    Map<String, Integer> map = new HashMap<String, Integer>();
    Set<String> urls = getURLs(term);
    for (String url: urls) {
        Integer count = getCount(url, term);
        map.put(url, count);
    }
    return map;
}
```

This method uses two helpers:

- `getURLs` takes a search term and returns the Set of URLs where the term appears.
- `getCount` takes a URL and a term and returns the number of times the term appears at the given URL.

Here are the implementations:

```
public Set<String> getURLs(String term) {
    Set<String> set = jedis.smembers(urlSetKey(term));
    return set;
}

public Integer getCount(String url, String term) {
    String redisKey = termCounterKey(url);
    String count = jedis.hget(redisKey, term);
    return new Integer(count);
}
```

Because of the way we designed the index, these methods are simple and efficient.

20.4 Analysis of lookup

Suppose we have indexed N pages and discovered M unique search terms. How long will it take to look up a search term? Think about your answer before you continue.

To look up a search term, we run `getCounts`, which

1. Creates a map.
2. Runs `getURLs` to get a Set of URLs.
3. For each URL in the Set, it runs `getCount` and adds an entry to a HashMap.

`getURLs` takes time proportional to the number of URLs that contain the search term. For rare terms, that might be a small number, but for common terms it might be as large as N .

Inside the loop, we run `getCount`, which finds a `TermCounter` on Redis, looks up a term, and adds an entry to a HashMap. Those are all constant time operations, so the overall complexity of `getCounts` is $O(N)$ in the worst case. However, in practice the runtime is proportional to the number of pages that contain the term, which is normally much less than N .

This algorithm is about as efficient as it can be, in terms of algorithmic complexity, but it is very slow because it sends many small operations to Redis. You can make it much faster using a `Transaction`. You might want to do that as an exercise, or you can see our solution in `RedisIndex.java`.

DRAFT: Not ready for distribution!

Chapter 21

Analysis of indexing

Using the data structures we designed, how long will it take to index a page? Again, think about your answer before you continue.

To index a page, we traverse its DOM tree, find all the `TextNode` objects, and split up the strings into search terms. That all takes time proportional to the number of words on the page.

For each term, we increment a counter in a `HashMap`, which is a constant time operation. So making the `TermCounter` takes time proportional to the number of words on the page.

Pushing the `TermCounter` to Redis requires deleting a `TermCounter`, which is linear in the number of unique terms. Then for each term we have to

1. Add an element to a `URLSet`, and
2. Add an element to a Redis `TermCounter`.

Both of these are constant time operations, so the total time to push the `TermCounter` is linear in the number of unique search terms.

In summary, making the `TermCounter` is proportional to the number of words on the page. Pushing the `TermCounter` to Redis is proportional to the number of unique terms.

Since the number of words on the page usually exceeds the number of unique search terms, the overall complexity is proportional to the number of words on the page. In theory a page might contain all search terms in the index, so the worst case performance is $O(M)$, but we don't expect to see the worse case in practice.

This analysis suggests a way to improve performance: we should probably avoid indexing very common words. First of all, they take up a lot of time and space, because they appear in almost every `URLSet` and `TermCounter`. Furthermore, they are not very useful because they don't help identify relevant pages.

Most search engines avoid indexing common words, which are known in this context as https://en.wikipedia.org/wiki/Stop_words stop words.

21.1 Graph traversal

If you did the “Getting to Philosophy” exercise, you already have a program that reads a Wikipedia page, finds the first link, uses the link to load the next page, and repeats. This program is a specialized kind of crawler, but when people say “Web crawler” they usually mean a program that:

- Loads a starting page and indexes the contents,
- Finds all the links on the page and adds the linked URLs to a queue, and
- Works its way through the queue, loading pages, indexing them, and adding new URLs to the queue.
- If it finds a URL in the queue that has already been indexed, it skips it.

You can think of the Web as a [https://en.wikipedia.org/wiki/Graph_\(discrete_mathematics\)](https://en.wikipedia.org/wiki/Graph_(discrete_mathematics)) graph where each page is a node and each link is a directed edge from one node to another. Starting from a source node, a crawler traverses this graph, visiting each reachable node once.

The behavior of the queue determines what kind of traversal the crawler performs:

- If the queue is first-in-first-out (FIFO), the crawler performs a breadth-first traversal.
- If the queue is last-in-first-out (LIFO), the crawler performs a depth-first traversal.
- More generally, the items in the queue might be prioritized. For example, we might want to give higher priority to pages that have not been indexed for a long time.

You can https://en.wikipedia.org/wiki/Graph_traversal read more about graph traversal here

21.2 Making a crawler

Now it's time to write the crawler.

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory with the code for this exercise, `javacs-lab11`.

In the subdirectory `javacs-lab11/src/com/flatiron-school/javacs` you'll find the source files for this exercise:

- * `'WikiCrawler.java'`, which contains starter code for your crawler.
- * `'WikiCrawlerTest.java'`, which contains test code for `'WikiCrawler'`.
- * `'JedisIndex.java'`, which is our solution to the previous exercise.

You'll also find some of the helper classes we've used in previous labs

- * `'JedisMaker.java'`
- * `'WikiFetcher.java'`
- * `'TermCounter.java'`
- * `'WikiNodeIterable.java'`

And as usual, in `javacs-lab11`, you'll find the Ant build file `build.xml`.

Before you run `JedisMaker`, you have to provide a file with information about your Redis server. If you did this in the previous lab, you can just copy it over. Otherwise you can find instructions in the previous exercise.

Run `ant build` to compile the source files, then run `ant JedisMaker` to make sure it is configured to connect to your Redis server.

Now run `ant test` to run `WikiCrawlerTest`. It should fail, because you have work to do!

Here's the beginning of the `WikiCrawler` class we provided:

```
public class WikiCrawler {

    public final String source;
    private JedisIndex index;
    private Queue<String> queue = new LinkedList<String>();
    final static WikiFetcher wf = new WikiFetcher();

    public WikiCrawler(String source, JedisIndex index) {
        this.source = source;
        this.index = index;
        queue.offer(source);
    }

    public int queueSize() {
        return queue.size();
    }
}
```

The instance variables are

- `source` is the URL where we start crawling.
- `index` is the `JedisIndex` where the results should go.
- `queue` is a `LinkedList` where we keep track of URLs that have been discovered but not yet indexed.
- `wf` is the `WikiFetcher` we'll use to read and parse Web pages.

Your job is to fill in `crawl`. Here's the prototype.

```
public String crawl(boolean testing) throws IOException {}
```

The parameter `testing` will be `true` when this method is called from `WikiCrawlerTest` and should be `false` otherwise.

When testing is `true`, the `crawl` method should:

- Choose and remove a URL from the queue in FIFO order.
- Read the contents of the page using `WikiFetcher.readWikipedia`, which reads cached copies of pages we have included in this repository for testing purposes (to avoid problems if the Wikipedia version changes).
- It should index pages regardless of whether they are already indexed.
- It should find all the internal links on the page and add them to the queue in the order they appear. “Internal links” are links to other Wikipedia pages.
- And it should return the URL of the page it indexed.

When testing is `false`, this method should:

- Choose and remove a URL from the queue in FIFO order.
- If the URL is already indexed, it should not index it again, and should return `null`.
- Otherwise it should read the contents of the page using `WikiFetcher.fetchWikipedia`, which reads current content from the Web.
- Then it should index the page, add links to the queue, and return the URL of the page it indexed.

`WikiCrawlerTest` loads the queue with about 200 links and then invokes `crawl` three times. After each invocation, it checks the return value and the new length of the queue.

When your crawler is working as specified, this test should pass. Good luck!

21.3 Resources

- https://en.wikipedia.org/wiki/Stop_words Stop words
- [https://en.wikipedia.org/wiki/Graph_\(discrete_mathematics\)](https://en.wikipedia.org/wiki/Graph_(discrete_mathematics)) Graphs
- https://en.wikipedia.org/wiki/Graph_traversal Graph traversal #
cs-application-retrieval-lab

DRAFT: Not ready for distribution!

Chapter 22

Boolean search

22.1 Learning goals

1. Analyze the performance of Web indexing algorithms.
2. Use boolean operators to generate search results for multiple search terms.
3. Score search results by relevance and sort them.

22.2 Overview

In this exercise, we present our solution to the previous exercise. Then you will write code to combine multiple search results and sort them by their relevance to the search terms.

22.3 Crawler solution

First, let's go over our solution to the previous exercise. We provided an outline of `WikiCrawler`; your job was to fill in `crawl`. As a reminder, here are the fields in the `WikiCrawler` class:

```
public class WikiCrawler {
    // keeps track of where we started
    private final String source;

    // the index where the results go
    private JedisIndex index;

    // queue of URLs to be indexed
    private Queue<String> queue = new LinkedList<String>();

    // fetcher used to get pages from Wikipedia
    final static WikiFetcher wf = new WikiFetcher();
}
```

When we create a `WikiCrawler`, we provide `source` and `index`. Initially, `queue` contains only one element, `source`.

Notice that the implementation of `queue` is a `LinkedList`, so we can add elements at the end – and remove them from the beginning – in constant time. By assigning a `LinkedList` object to a `Queue` variable, we limit ourselves to using methods in the `Queue` interface; specifically, we'll use `offer` to add elements and `poll` to remove them.

Here's our implementation of `crawl`:

```
public String crawl(boolean testing) throws IOException {
    if (queue.isEmpty()) {
        return null;
    }
    String url = queue.poll();
    System.out.println("Crawling " + url);

    if (testing==false && index.isIndexed(url)) {
        System.out.println("Already indexed.");
        return null;
    }

    Elements paragraphs;
    if (testing) {
```

```
        paragraphs = wf.readWikipedia(url);
    } else {
        paragraphs = wf.fetchWikipedia(url);
    }
    index.indexPage(url, paragraphs);
    queueInternalLinks(paragraphs);
    return url;
}
```

Most of the complexity in this method is there to make it easier to test. Here's the logic:

- If the queue is empty, it returns `null` to indicate that it did not index a page.
- Otherwise it removes and stores the next URL from the queue.
- If the URL has already been indexed, `crawl` doesn't index it again, unless it's in testing mode.
- Next it reads the contents of the page: if it's in testing mode, it reads from a file; otherwise it reads from the Web.
- It indexes the page.
- It parses the page and adds internal links to the queue.
- Finally, it returns the URL of the page it indexed.

We presented our implementation of `index.indexPage` in the previous exercise. So the only new function is `queueInternalLinks`.

We wrote two versions of this function with different parameters: one takes an `Elements` object containing one DOM trees per paragraph; the other takes an `Element` object that contains a single paragraph.

The first version just loops through the paragraphs. The second version does the real work.

```
void queueInternalLinks(Elements paragraphs) {
    for (Element paragraph: paragraphs) {
        queueInternalLinks(paragraph);
    }
}

private void queueInternalLinks(Element paragraph) {
    Elements elts = paragraph.select("a[href]");
    for (Element elt: elts) {
        String relURL = elt.attr("href");

        if (relURL.startsWith("/wiki/")) {
            String absURL = elt.attr("abs:href");
            queue.offer(absURL);
        }
    }
}
```

To determine whether a link is “internal,” we check whether the URL starts with “/wiki/”. This might include some pages we don’t want to index, like meta-pages about Wikipedia. And it might exclude some pages we want, like links to pages in non-English languages. But we kept it simple.

That’s all there is to it. This exercise didn’t have a lot of new material; it was mostly a chance to bring the pieces together.

22.4 Information retrieval

The next phase of this project is to implement a search tool. The pieces we’ll need include:

1. An interface where users can provide search terms and view results.
2. A lookup mechanism that takes each search term and returns the pages that contain it.
3. Mechanisms for combining search results from multiple search terms.

4. Algorithms for ranking and sorting search results.

The general term for processes like this is “information retrieval”, https://en.wikipedia.org/wiki/Information_retrieval which you can read more about here.

In this exercise, we’ll focus on steps 3 and 4. We’ve already build a simple version of 2. If you are interested in building Web applications, you’ll have to option to work on step 1.

22.5 Boolean search

Most search engines can perform “boolean searches”, which means you can combine the results from multiple search terms using boolean logic. For example:

- The search “java AND programming” might return only pages that contain both search terms: “java” and “programming”.
- “java OR programming” might return pages that contain either term but not necessarily both.
- “java -indonesia” might return pages that contain “java” and do not contain “indonesia”.

Expressions like these that contain search terms and operators are called “queries”.

When applied to search results, the boolean operators AND, OR, and - correspond to the set operations **intersection**, **union**, and **difference**. For example, suppose

- **s1** is the set of pages containing “java”,
- **s2** is the set of pages containing “programming”, and
- **s3** is the set of pages containing “indonesia”.

In that case:

- The intersection of `s1` and `s2` is the set of pages containing “java” AND “programming”.
- The union of `s1` and `s2` is the set of pages containing “java” OR “programming”.
- The difference of `s1` and `s2` is the set of pages containing “java” and not “indonesia”.

In the next section you will write method to implement these operations.

22.6 Relevance scores

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory with the code for this exercise, `javacs-lab12`.

In the subdirectory `javacs-lab12/src/com/flatiron/school/javacs` you’ll find the source files for this exercise:

- `WikiSearch.java`, which defines an object that contains search results and performs operations on them.
- `WikiSearchTest.java`, which contains test code for `WikiSearch`.
- `Card.java`, which demonstrates how to use the `sort` method in `java.util.Collections`.

You will also find some of the helper classes we’ve used in previous labs.

Here’s the beginning of the `WikiSearch` class definition:

```
public class WikiSearch {  
  
    // map from URLs that contain the term(s) to relevance score  
    private Map<String, Integer> map;  
  
    public WikiSearch(Map<String, Integer> map) {  
        this.map = map;  
    }  
}
```

```
}

public Integer getRelevance(String url) {
    Integer relevance = map.get(url);
    return relevance==null ? 0: relevance;
}
```

A `WikiSearch` object contains a map from URLs to their relevance score. In the context of information retrieval, a “relevance score” is a number intended to indicate how well a page meets the needs of the user as inferred from the query. There are many ways to construct a relevance score, but most of them are based on “term frequency”, which is the number of times the search terms appear on the page. A common relevance score is called TF-IDF, which stands for “term frequency – inverse document frequency”. <https://en.wikipedia.org/wiki/Tf%E2%80%93idf> You can read more about it here.

You’ll have the option to implement TF-IDF later, but we’ll start with something even simpler, TF:

- If a query contains a single search term, the relevance of a page is its term frequency; that is, the number of time the term appears on the page.
- For queries with multiple terms, the revelance of a page is the sum of the term frequencies; that is, the total number of times any of the search terms appear.

Now you’re ready to start the exercise.

22.7 Implementing boolean operators

Assuming you have already checked out the repository for this exercise, you should have a directory named `javacs-lab12` that contains the Ant build file `build.xml`.

Run `ant build` to compile the source files, then run `ant test` to run `WikiSearchTest`. As usual, it should fail, because you have work to do.

In `WikiSearch.java`, fill in the bodies of `and`, `or`, and `minus` so that the relevant tests pass. You don't have to worry about `testSort` yet.

You can run `WikiSearchTest` without using Jedis because it doesn't depend on the index in your Redis database. But if you want to run a query against your index, you have to provide a file with information about your Redis server. If you did this in the previous lab, you can just copy it over. Otherwise you can find instructions in [ADD%20THIS%20LINK%20AFTER%20DEPLOYMENT](#)cs-application-backing-with-redis-lab.

Run `ant JedisMaker` to make sure it is configured to connect to your Redis server. Then run `WikiSearch`, which prints results from three queries:

- “java”
- “programming”
- “java AND programming”

Initially the results will be in no particular order, because `WikiSearch.sort` is incomplete.

Fill in the body of `sort` so the results are returned in increasing order of relevance. We suggest you use the `sort` method provided by <https://docs.oracle.com/javase/7/docs/api/java/util/Collections.html> `java.util.Collections`, which sorts any kind of `List`. There are two version of `sort`:

- The one-parameter version takes a list and sorts the elements using the `compareTo` method, so the elements have to be `Comparable`.
- The two-parameter version takes a list of any object type and a `Comparator`, which is an object that provides a `compare` method that compares elements.

If you are not familiar with the `Comparable` and `Comparator` interfaces, we explain them in the next section.

22.8 Comparable and Comparator

This exercise includes `Card.java`, which demonstrates two ways to sort a list of `Card` objects. Here's the beginning of the class definition:

```
public class Card implements Comparable<Card> {

    private final int rank;
    private final int suit;

    public Card(int rank, int suit) {
        this.rank = rank;
        this.suit = suit;
    }
```

A `Card` object has two integer fields, `rank` and `suit`. `Card` implements `Comparable<Card>`, which means that it provides `compareTo`:

```
    public int compareTo(Card that) {
        if (this.suit < that.suit) {
            return -1;
        }
        if (this.suit > that.suit) {
            return 1;
        }
        if (this.rank < that.rank) {
            return -1;
        }
        if (this.rank > that.rank) {
            return 1;
        }
        return 0;
    }
```

The specification of `compareTo` indicates that it should return a negative number if **this** is considered less than **that**, a positive number if it is considered greater, and 0 if they are considered equal.

If you use the one-parameter version of `Collections.sort`, it uses the `compareTo` method provided by the elements to sort them. To demonstrate, we can make a list of 52 cards like this:

```
public static List<Card> makeDeck() {
    List<Card> cards = new ArrayList<Card>();
    for (int suit = 0; suit <= 3; suit++) {
        for (int rank = 1; rank <= 13; rank++) {
            Card card = new Card(rank, suit);
            cards.add(card);
        }
    }
    return cards;
}
```

And sort them like this:

```
Collections.sort(cards);
```

This version of `sort` puts the elements in what's called their "natural order" because it's determined by the objects themselves.

But it is possible to impose a different ordering by providing a `Comparator` object. For example, the natural order of `Card` objects treats Aces as the lowest rank, but in some card games they have the highest rank. We can define a `Comparator` that considers "Aces high", like this:

```
Comparator<Card> comparator = new Comparator<Card>() {
    @Override
    public int compare(Card card1, Card card2) {
        if (card1.getSuit() < card2.getSuit()) {
            return -1;
        }
        if (card1.getSuit() > card2.getSuit()) {
            return 1;
        }
        int rank1 = getRankAceHigh(card1);
        int rank2 = getRankAceHigh(card2);

        if (rank1 < rank2) {
            return -1;
        }
        if (rank1 > rank2) {
            return 1;
        }
    }
}
```

```
        }
        return 0;
    }

    private int getRankAceHigh(Card card) {
        int rank = card.getRank();
        if (rank == 1) {
            return 14;
        } else {
            return rank;
        }
    }
};
```

This code defines an anonymous class that implements `compare`, as required. Then it creates an instance of the class. If you are not familiar with anonymous classes in Java, <https://docs.oracle.com/javase/tutorial/java/java00/anonymousclasses.html> you can read about them here.

Using this `Comparator`, we can invoke `sort` like this:

```
Collections.sort(cards, comparator);
```

In this ordering, the Ace of Spaces is considered the highest class in the deck; the two of clubs is the lowest.

The code in this section is in `Card.java` if you want to experiment with it. As an exercise, you might want to write a comparator that sorts by `rank` first and then by `suit`, so all the Aces should be together, and all the twos, etc.

22.9 Extensions

If you get a basic version of this exercise working, you might want to work on these optional exercises:

- <https://en.wikipedia.org/wiki/Tf%E2%80%93idf> Read about TF-IDF and implement it. You might have to modify `JavaIndex` to compute document frequencies; that is, the total number of times each term appears on all pages in the index.

- For queries with more than one search term, the total relevance for each page is currently the sum of the relevance for each term. Think about when this simple version might not work well, and try out some alternatives.
- Build a user interface that allows users to enter queries with boolean operators. Parse the queries, generate the results, then sort them by relevance and display the highest-scoring URLs. Consider generating “snippets” that show where the search terms appeared on the page. If you want to make a Web application for your user interface, we recommend <https://devcenter.heroku.com/articles/getting-started-with-java> Heroku as simple options for developing and deploying Web applications using Java.

22.10 Resources

- https://en.wikipedia.org/wiki/Information_retrieval Information retrieval.
- <https://en.wikipedia.org/wiki/Tf%E2%80%93idf> TF-IDF
- <https://docs.oracle.com/javase/tutorial/java/java00/anonymousclasses.html> Anonymous classes # cs-application-sorting-lab

Chapter 23

Sorting

23.1 Learning goals

1. Implement and analyze sort algorithms.
2. Use a heap to implement a `PriorityQueue`.
3. Use a `PriorityQueue` to find the top elements in a large dataset.
4. Analyze the space needs of algorithms.

23.2 Overview

Computer science departments have an unhealthy obsession with sort algorithms. Based on the amount of time CS students spend on the topic, you would think that choosing sort algorithms is the cornerstone of modern software engineering. Of course, the reality is that software developers can go years, or entire careers, without thinking about how sorting works. For almost all applications, they use whatever general-purpose algorithm is provided by the language or libraries they use. And usually that's just fine.

So if you skip this exercise and learn nothing about sort algorithms, you can still be an excellent developer. But there are a few reasons you might want to do it anyway:

1. Although there are general-purpose algorithms that work well for the vast majority of applications, there are two special-purpose algorithms you might need to know about: radix sort and bounded heap sort.
2. One sort algorithm, merge sort, makes an excellent teaching example because it demonstrates an important and useful strategy for algorithm design, called “divide-conquer-glue”. Also, when we analyze its performance, you will learn about an order of growth we have not seen before. Finally, some of the most widely-used algorithms are hybrids that include elements of merge sort.
3. One other reason to learn about sort algorithms is that technical interviewers love to ask about them. If you want to get hired, it helps if you can demonstrate CS cultural literacy.

So, in this exercise we’ll analyze insertion sort, you will implement merge sort, we’ll tell you about radix sort, and you will write a simple version of a bounded heap sort.

23.3 Insertion sort

We’ll start with insertion sort, mostly because it is simple to describe and implement. It is not very efficient, but it has some redeeming qualities, as we’ll see.

Rather than explain the algorithm here, we suggest you read the https://en.wikipedia.org/wiki/Insertion_sort insertion sort Wikipedia page, which includes pseudocode and animated examples. Come back when you get the general idea.

Here’s our implementation of insertion sort in Java:

```
public class ListSorter<T> {  
  
    public void insertionSort(List<T> list, Comparator<T> comparator) {  
  
        for (int i=1; i < list.size(); i++) {  
            T elt_i = list.get(i);
```

```

        int j = i;
        while (j > 0) {
            T elt_j = list.get(j-1);
            if (comparator.compare(elt_i, elt_j) >= 0) {
                break;
            }
            list.set(j, elt_j);
            j--;
        }
        list.set(j, elt_i);
    }
}

```

We define a class, `ListSorter`, as a container for sort algorithms. By using the type parameter, `T`, we can write methods that work on lists containing any object type.

`insertionSort` takes two parameters, a `List` of any kind and a `Comparator` that knows how to compare type `T` objects. It sorts the list “in place”, which means it modifies the existing list and does not have to allocate any new space.

This example shows how to call this function with a `List` of `Integer`:

```

List<Integer> list = new ArrayList<Integer>(Arrays.asList(3, 5, 1, 4, 2));

Comparator<Integer> comparator = new Comparator<Integer>() {
    @Override
    public int compare(Integer n, Integer m) {
        return n.compareTo(m);
    }
};

ListSorter<Integer> sorter = new ListSorter<Integer>();
sorter.insertionSort(list, comparator);
System.out.println(list);

```

`insertionSort`, has two nested loops, so you might guess that its runtime is quadratic. In this case, that turns out to be correct, but before you jump to

that conclusion, you have to check that the number of times each loop runs is proportional to n , the size of the array.

The outer loop iterates from 1 to n , so it is clearly linear. The inner loop iterates from i to 0, so it is also linear, and the total number of times the inner loop runs is quadratic.

If you are not sure about that, here's the argument:

- The first time through, $i=1$ and the inner loop runs at most once.
- The second time, $i=2$ and the inner loop runs at most twice.
- The last time, $i=n-1$ and the inner loop runs at most $n-1$ times.

So the total number of times the inner loop runs is the sum of the series 1, 2, ... $n-1$, which is $n(n-1) / 2$. And the leading term of that expression (the one with the highest exponent) is n^2 .

In the worst case, insertion sort is quadratic. However:

1. If the elements are already sorted, or nearly so, insertion sort is linear. Specifically, if each element is no more than k locations away from where it should be, the inner loop never runs more than k times, and the total runtime is $O(kn)$.
2. Because the implementation is simple, the overhead is low; that is, even if the runtime is $a n^2$, the coefficient of the leading term, a , is probably small.

So if we know that the array is nearly sorted, or is not very big, insertion sort might not be a bad choice. But for large arrays, we can do better. In fact, much better.

23.4 Merge sort

Merge sort is one of several algorithms whose runtime is better than quadratic (we're still not saying how much better). Again, rather than explaining the algorithm here, we suggest you https://en.wikipedia.org/wiki/Merge_

[sort](#) read about it on Wikipedia. Once you get the idea, come back and you can test your understanding by writing your own implementation.

When you check out the repository for this exercise, you should find a file structure similar to what you saw in previous exercises. The top level directory contains `CONTRIBUTING.md`, `LICENSE.md`, `README.md`, and the directory with the code for this exercise, `javacs-lab13`.

In the subdirectory `javacs-lab12/src/com/flatiron-school/javacs` you'll find the source files for this exercise:

- `ListSorter.java`,
- `ListSorterTest.java`,

Run `ant build` to compile the source files, then run `ant test` to run `ListSorterTest`. As usual, it should fail, because you have work to do.

In `ListSorter.java`, we've provided an outline of two methods, `mergeSortInPlace` and `mergeSort`:

```
public void mergeSortInPlace(List<T> list, Comparator<T> comparator) {
    List<T> sorted = mergeSortHelper(list, comparator);
    list.clear();
    list.addAll(sorted);
}

private List<T> mergeSort(List<T> list, Comparator<T> comparator) {
    // TODO: fill this in!
    return null;
}
```

These two methods do the same thing but provide different interfaces. `mergeSort` takes a list and returns a new list with the same elements sorted in ascending order. `mergeSortInPlace` is a `void` method that modifies an existing list.

Your job is to fill in `mergeSort`. Before you write a fully recursive version of merge sort, start with something like this:

1. Split the list in half.

2. Sort the halves using `Collections.sort` or `insertionSort`.
3. Merge the sorted halves into a complete sorted list.

This will give you a chance to debug the merge code without dealing with the complexity of a recursive method.

Next, add a https://en.wikipedia.org/wiki/Recursion#base_case base case. If you are given a list with only one element, you could return it immediately, since it is already sorted, sort of. Or if the length of the list is below some threshold, you could sort it using `Collections.sort` or `insertionSort`. Test the base case before you proceed.

Finally, modify your solution so it makes two recursive calls to sort the halves of the array. When you get it working, `testMergeSort` and `testMergeSortInPlace` should pass.

23.5 Analysis of merge sort

To classify the runtime of merge sort, it helps to think about the algorithm in terms of levels of recursion and how much work is done on each level. Suppose we start with a list that contains n elements. Here are the steps of the algorithm:

1. Make two new arrays and copy half of the elements into each.
2. Sort the two halves.
3. Merge the halves.

The following figure shows these steps.

The first step copies each of the elements once, so it is linear. The third step also copies each element once, so it is also linear. Now we need to figure out the complexity of step 2. To do that, it helps to look at a different picture of the computation, which shows the levels of recursion:

At the top level, we have 1 list with n elements. At the next level there are 2 lists with $n/2$ elements. Then 4 lists with $n/4$ elements, and so on until we get to n lists with 1 element.

On every level we have a total of n elements. On the way down, we have to split the arrays in half, which takes time proportional to n on every level. On the way back up, we have to merge a total of n elements, which is also linear.

If the number of levels is h , the total amount of work for the algorithm is $O(nh)$. So how many levels are there? There are two ways to think about that:

1. How many times do we have to cut n in half to get to 1. Or,
2. How many times do we have to double 1 before we get to n .

Another way to ask the second question is “what power of 2 is n ”?

$$2^h = n$$

Taking the \log of both sides yields

$$h = \log_2 n$$

So the total time is $O(n \log n)$. I didn’t bother to write the base of the logarithm because logarithms with different bases differ by a constant factor, so all logarithms are in the same order of growth.

Algorithms in $O(n \log n)$ are technically called “linearithmic”, but most people just say “ $n \log n$ ”.

It turns out that $O(n \log n)$ is the theoretical lower bound for sort algorithms that work by comparing elements to each other. That means there is no https://en.wikipedia.org/wiki/Comparison_sort comparison sort whose order of growth is better than $n \log n$. Nevertheless, there are non-comparison sorts that take linear time!

23.6 Radix sort

During the 2008 United States Presidential Campaign, candidate Barack Obama was asked to perform an impromptu algorithm analysis when he visited Google. Chief executive Eric Schmidt jokingly asked him for “the most efficient way to sort a million 32-bit integers.” Obama had apparently been tipped off,

because he quickly replied, “I think the bubble sort would be the wrong way to go.” You can http://www.youtube.com/watch?v=k4RRi_ntQc8 watch the video here.

Obama was right: https://en.wikipedia.org/wiki/Bubble_sort bubble sort is conceptually simple but its runtime is quadratic; and even among quadratic sort algorithms, its performance is not very good.

The answer Schmidt was probably looking for is “radix sort”, which is a *non-comparison* sort algorithm that works if the size of the elements is bounded, like a 32-bit integer or a 20-character string.

To see how this works, imagine you have a deck of cards where each card contains a 3-letter word. Here’s how you could sort the cards:

1. Make one pass through the cards and divide them into buckets based on the first letter. So words starting with **a** should be together, followed by words starting with **b**, and so on.
2. Divide each bucket again based on the second letter. So words starting with **aa** should be together, followed by words starting with **ab**, and so on. Of course, not all buckets will be full, but that’s ok.
3. Divide each bucket again based on the third letter.

At this point each bucket contains one element, and the buckets are sorted in ascending order. The following diagram shows an example with 3-letter words:

The top row shows the unsorted words. The second row shows what the buckets look like after the first pass. The words in each bucket begin with the same letter.

After the second pass, the words in each bucket begin with the same two letters. After the third pass, there can be only one word in each bucket, and the buckets are in order.

During each pass, we iterate through the elements and add them to buckets. As long as the buckets allow addition in constant time, each pass is linear.

The number of passes, which we’ll call w , depends on the length of the words, but it doesn’t depend on the number of words, n . So the order of growth is $O(wn)$, which is linear in n .

There are many variations on radix sort, and many ways to implement each one. You can http://en.wikipedia.org/wiki/Radix_sort read more about them here. As an optional exercise, consider writing a version of radix sort.

23.7 Heap sort

In addition to radix sort, which applies when the things you want to sort are bounded in size, there is one other special-purpose sorting algorithm you might encounter: bounded heap sort. Bounded heap sort is useful if you are working with a very large dataset and you want to report the “Top 10”, or “Top k ” for some value of k much smaller than n .

For example, suppose you are monitoring a Web service that handles a billion transactions per day. At the end of each day, you want to report the k biggest transactions (or slowest, or any other superlative). One option is to store all transactions, sort them at the end of the day, and select the top k . That would take time proportional to $n \log n$, and it would be very slow because we probably can’t fit a billion transactions in the memory of a single program. We would have to use an “out of core” sort algorithm. You can https://en.wikipedia.org/wiki/External_sorting read about “external sorting” here.

Using a bounded heap, we can do much better! We’ll explain how in three steps:

1. We’ll explain (unbounded) heap sort.
2. You’ll implement it.
3. We’ll explain bounded heap sort and analyze it.

To understand heap sort, you have to understand a heap, which is a data structure similar to a binary search tree (BST). Here’s the difference:

- In a BST, every node, x , has the “BST property”: all nodes in the left subtree of x are less than x and all nodes in the right subtree are greater than x .
- In a heap, every node, x , has the “heap property”: all nodes in both subtrees of x are greater than x .

The smallest element in a heap is always at the root, so we can find it in constant time. Adding and removing elements from a heap takes time proportional to the height of the tree h . And it is easy to keep a heap balanced, so h is proportional to $\log n$. You can [https://en.wikipedia.org/wiki/Heap_\(data_structure\)](https://en.wikipedia.org/wiki/Heap_(data_structure)) read more about heaps here.

The Java `PriorityQueue` is implemented using a heap. `PriorityQueue` provides the methods specified in the `Queue` interface, including `offer` and `poll`:

- **offer**: Adds an element to the queue, updating the heap so that every node has the “heap property”. Takes $\log n$ time.
- **poll**: Removes the smallest element in the queue from the root and updates the heap. Takes $\log n$ time.

Given a `PriorityQueue`, you can easily sort of a collection of n elements like this:

1. Add all elements of the collection to a `PriorityQueue` using `offer`.
2. Remove the elements from the queue using `poll` and add them to a `List`.

Because `poll` returns the smallest element remaining in the queue, the elements are added to the `List` in ascending order. This way of sorting is called <https://en.wikipedia.org/wiki/Heapsort> heap sort.

Adding n elements to the queue takes $n \log n$ time. So does removing n elements. So the run time for heap sort is $O(n \log n)$.

- In `ListSorter.java` you’ll find the outline of a method called `heapSort`. Fill it in and then run `ant test` to confirm that it works.

23.8 Bounded heap

A bounded heap is a heap that is limited to contain at most k elements. If you have n elements, you can keep track of the k largest elements like this:

For each element, x :

- Branch 1: If the heap is not full, add x to the heap.
- Branch 2: If the heap is full, compare x to the *smallest* element in the heap. If x is smaller, it cannot be one of the largest k elements, so you can discard it.
- Branch 3: If the heap is full and x is greater than the smallest element in the heap, remove the smallest element from the heap and add x .

Using a heap with the smallest element at the top, we can keep track of the largest k elements. Let's analyze the performance of this algorithm. For each element, we perform one of:

- Branch 1: Adding an element to the heap is $O(\log k)$.
- Branch 2: Finding the smallest element in the heap is $O(1)$.
- Branch 3: Removing the smallest element is $O(\log k)$. Adding x is also $O(\log k)$.

In the worst case, if the elements appear in ascending order, we always run Branch 3. In that case, the total time to process n elements is $O(n \log k)$, which is linear in n .

- In `ListSorter.java` you'll find the outline of a method called `topK` that takes a `List`, a `Comparator`, and an integer k . It should return the k largest elements in the `List` in ascending order. Fill it in and then run `ant test` to confirm that it works.

23.9 Space complexity

Until now we have talked a lot about runtime analysis, but for many algorithms we are also concerned about space. For example, one of the drawbacks of merge sort is that it makes copies of the data. In our implementation, the total amount of space it allocates is $O(n \log n)$. With a more clever implementation, you can get the space requirement down to $O(n)$.

In contrast, insertion sort doesn't copy the data because it sorts the elements in place. It uses temporary variables to compare two elements at a time, and it uses a few other local variables. But its space use doesn't depend on n .

Our implementation of heap sort creates a new `PriorityQueue` to store the elements, so the space is $O(n)$; but if you are allowed to sort the list in place, you can run heap sort with $O(1)$ space.

One of the benefits of the bounded heap algorithm you just implemented is that it only needs space proportional to k (the number of elements we want to keep), and k is often much smaller than n .

Software developers tend to pay more attention to runtime than space, and for many applications, that's appropriate. But for large datasets, space can be just as important or more so. For example,

1. If a dataset doesn't fit into the memory of one program, the runtime often increases dramatically, or it might not run at all. If you choose an algorithm that needs less space, and that makes it possible to fit the computation into memory, it might run much faster. In the same vein, a program that uses less space might make better use of https://en.wikipedia.org/wiki/CPU_cache CPU caches and run faster.
2. On a server that runs many programs at the same time, if you can reduce the space needed for each program, you might be able to run more programs on the same server, which reduces hardware and energy costs.

So those are some reasons you should know at least a little bit about the space needs of algorithms.

23.10 Resources

- https://en.wikipedia.org/wiki/Comparison_sort Comparison sort
- https://en.wikipedia.org/wiki/Merge_sort Merge sort
- https://en.wikipedia.org/wiki/Insertion_sort Insertion sort

- http://en.wikipedia.org/wiki/Radix_sortRadix sort
- https://en.wikipedia.org/wiki/External_sortingExternal sorting
- [https://en.wikipedia.org/wiki/Heap_\(data_structure\)](https://en.wikipedia.org/wiki/Heap_(data_structure))Heap data structure
- <https://en.wikipedia.org/wiki/Heapsort>Heap sort
- https://en.wikipedia.org/wiki/CPU_cacheCPU cache

DRAFT: Not ready for distribution!

DRAFT: Not ready for distribution!

Index

GitHub, xiv

repository, xiv

DRAFT: Not ready for distribution!