

# Principles of Computer Science

Michael C Horsch

Course readings for  
CMPT 145

Copyright © 2017–2020 Michael C. Horsch

PRODUCED BY THE AUTHOR FOR STUDENTS IN CMPT 145.

CS.USASK.CA

LaTeX style files used under the Creative Commons Attribution-NonCommercial 3.0 Unported License, Mathias Legrand ([legrand.mathias@gmail.com](mailto:legrand.mathias@gmail.com)) downloaded from [www.LaTeXTemplates.com](http://www.LaTeXTemplates.com).

Cover and chapter heading images are in the public domain downloaded from <http://wallpaperspal.com>.

This document is licensed under the Creative Commons Attribution-NonCommercial 3.0 Unported License (the “License”). You may not use this file except in compliance with the License. You may obtain a copy of the License at <http://creativecommons.org/licenses/by-nc/3.0>. Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

*Sixth Edition, July, 2020*

*Fifth Edition, January, 2019*

*Fourth Edition, May, 2018*

*Third Edition, January, 2018*

*Second Edition, July, 2017*

*First Edition, January, 2017*

The background of the top section is a dark red field with a complex network of glowing red lines and dots, resembling a molecular structure or a data network. These lines and dots are interconnected, creating a web-like pattern that fills the upper half of the page.

## Preface

Computer science is the study of computationally effective compositions. You may not think of programming as composition, but it is. It's actually a fairly challenging form of composition, because, unlike essays, poems, plays, songs, and novels, which have communicative meaning, a program has (or is intended to have) an operational meaning. Programs have to *work*, and even a tiny defect can make the difference between a program that works correctly, and one that works incorrectly, or doesn't work at all.

To begin the study of computer science, we need to start with a basic literacy in the various programming constructs. We need to understand that computing is about working with data to answer a question. We need to master basic computational devices, including data types, atomic and compound variables, functions, conditional execution, and repetitive execution. These can be learned in various courses, using almost any programming language; all computer programming languages, at this point in history, employ essentially the very same computational devices, despite superficial syntactic differences. We assume that the reader has at least one course behind them, and that they arrive here with basic competence with these concepts. We also assume that the reader knows at least a little Python; this is not absolutely necessary, as Python is a friendly language, and fairly easy to learn if you already know another programming language.

To make the jump from learning basic programming to computer science, we study algorithms, and the data structures that make algorithms effective. All of this could be done without any computers at all, but computer science is a lot more fun because we have computers to demonstrate the concepts that we might otherwise experience only in our imaginations. The fact that we have computers makes computer science pre-eminently practical, and that's an endless source of motivation.

Composing a program is the pinnacle of problem solving. The practical application of computer science requires a mindful discipline to manage the complexity of problem solving, and a scientific attitude towards the result of this work. There is no single right answer to any interesting problem, but many wrong answers, and many stumbling blocks on any number of paths. The number of possible combinations of computational concepts is literally infinite, but computer scientists do not have infinite time, and so have to find some way to navigate the infinite to arrive at a practical



solution. In response to a decade of teaching students faced with this challenging and exciting occupation, the author has put together this course material.

This book is not a traditional textbook. It's a collection of topical readings, arranged into a single document. The chapters are deliberately kept short, so that readers are not discouraged by long reading assignments. An author is always tempted to add insights and content to every page; I hope I have kept these impulses in check sufficiently.

This book is not intended to teach Python. We will use Python, and we will learn a few things about this important language, but that's not our primary focus. We can't and won't cover every detail that a deeply knowledgeable expert would know, because there are other books and future courses for that kind of thing. Anything we teach about the Python language is intended to clarify an important principle.

This book has a number of very strong influences. Primarily, the course material developed for the predecessor course, CMPT115, developed by Mark Eramian and Ian McQuillan. Several chapters were inspired by the work of Robert Sedgewick and Kevin Wayne; their introductory text, *Introduction to Programming in Python*, and their *Algorithms* text. I borrowed ideas shamelessly from both. My teaching career started as a teaching assistant for a brand new course and work-in-progress that turned into *The Schematics of Computation* by Vincent Manis and James J. Little. Less directly, *Structure and Interpretation of Computer Programs*, by Harold Abelson and Gerald Sussman, was also a source of inspiration. This book is only a shadow of their accomplishments. Chapter 14 has drawn a lot of inspiration from the classic *The Elements of Programming Style*, by Kernighan and Plauger.

I'd like to acknowledge the following people who have contributed to this effort. First and foremost, Mark Eramian and Jeff Long, for long discussions about how to change the way Computer Science is taught in first year at uSask, and for their inspired and dedicated work on those other two new courses in the program. Brittany Chan worked tirelessly to help bring into reality the ideas that the three of us came up with. I'd also like to thank Mark for contributing a ton of advice and hard work with L<sup>A</sup>T<sub>E</sub>X, and Python. For lecture material from previous courses: Mark Eramian, Ian McQuillan, Lingling Jin, and Dmytro Dyachuk. Cyril Coupal provided helpful feedback about the course as a co-instructor.

Undoubtedly there are omissions, and errors in the text. I would be happy to acknowledge corrections made by readers, so feel free to contact me when you find any. Students who have helped improve the text with comments and corrections: Yunfei Xiao, Paul Larson, David Graf, Jian Su, Tristen MacPherson, Allan McPherson, Killian Stacey, Kun Niu, David Sese, Peggy Anderson, Duane Bergstrom, Aurora Dubnyk, Tyler Akins, Melanie Gibbons. The responsibility for any and all remaining errors in this book lies with its author.



# Contents

<b>1</b>	<b>Computers and Software</b>	<b>15</b>
1.1	Computers	15
1.2	Languages for Computation	17
1.3	Software	22
1.4	Summary	24
<b>2</b>	<b>Computing with Python</b>	<b>25</b>
2.1	Two example problems, with solutions	25
2.1.1	Sorting a list	25
2.1.2	Counting prime numbers	26
2.2	A grab-bag of other listy stuff	28
2.2.1	Perspective on lists	29
2.3	Dictionaries	30
2.3.1	An example program	31
2.4	Modules	32
2.5	Three Example Problems	33
2.5.1	The Gambler's Ruin Problem	33
2.5.2	The Coupon Collector Problem	33
2.5.3	Self-Avoiding Random Walks	34
<b>3</b>	<b>References</b>	<b>35</b>
3.1	Introduction	35

3.2	Values, objects, and addresses	36
3.3	Expressions	36
3.4	Variables and frames	37
3.5	Equality	39
3.6	More about frames and functions	39
3.7	Frames and environments	41
3.8	Lists and Dictionaries	42
3.9	Python caches values for efficiency	43
3.9.1	Mutable vs Immutable Data	43
3.10	Casual language	44
3.11	Summary	44
<b>4</b>	<b>Software Design Goals</b>	<b>45</b>
4.1	Design Goals	45
4.2	Implementation Goals	46
4.3	Example	47
4.4	How to Use the Design and Implementation Goals Effectively	49
4.5	Pedagogical Goals	49
<b>5</b>	<b>Software Development Processes</b>	<b>51</b>
5.1	Software development processes	51
5.1.1	Waterfall model	52
5.1.2	Iterative and incremental model	53
5.1.3	Requirements and Specifications	53
5.2	Software Design Strategies	54
5.3	Version Control	55
5.4	Prototypes	55
5.5	Summary	56
<b>6</b>	<b>Procedural Abstraction</b>	<b>57</b>
6.1	Procedural Abstraction	57
6.1.1	The benefits of procedural abstraction	58
6.1.2	Multiple levels of procedural abstraction	58
6.2	Describing a Function: Interface Documentation	58
6.2.1	Documenting the interface after defining the function	60

6.3	As part of the implementation process	61
6.4	Using Procedural Abstraction in the Design Process	64
<b>7</b>	<b>Software Testing</b> .....	<b>65</b>
7.1	Basic Terminology	66
7.2	The importance of test scripts	68
7.3	Testing as error discovery	68
7.4	Unit test case design	69
7.5	No, really, how much testing do I have to do?	71
7.6	Debugging	71
<b>8</b>	<b>Abstract Data Types</b> .....	<b>75</b>
8.1	Data Types	76
8.2	Data Structures	76
8.3	Abstract Data Types (ADTs)	77
8.4	Examples of Abstract Data Types	78
8.4.1	Registry .....	78
8.4.2	Statistics ADT .....	79
8.4.3	Application example .....	80
8.4.4	Discussion .....	81
8.5	Summary	82
<b>9</b>	<b>Objects and Classes</b> .....	<b>83</b>
9.1	Introduction	83
9.2	Defining classes and instantiating objects	84
9.3	Simple examples	85
9.4	Encapsulation and access control	87
9.5	Using classes to define ADTs	88
9.5.1	Perspective .....	89
<b>10</b>	<b>Stacks and Queues</b> .....	<b>91</b>
10.1	Queues: First-In First-Out (FIFO)	91
10.2	Queue Application: Buffering communications	93
10.3	Stacks: Last-In First-Out (LIFO)	94
10.4	Stack Application: Backtracking	95

<b>11</b>	<b>Applications of Stacks and Queues</b>	<b>97</b>
11.1	Introduction	97
11.2	Queueing simulation	97
11.2.1	Discussion	101
11.3	Bracket checking	101
11.4	Post-fix Arithmetic evaluator	103
<b>12</b>	<b>List-based Stacks and Queues</b>	<b>107</b>
12.1	Why use Queues and Stacks when we have lists?	107
12.2	Adapting Python Lists to Implement Queues	108
12.2.1	Operation: Initialization	108
12.2.2	Operations: size() and is_empty()	109
12.2.3	Operation: enqueue()	109
12.2.4	Operation: dequeue()	110
12.2.5	Operation: peek()	110
12.3	Adapting Python Lists to Implement Stacks	110
12.3.1	Operation: create()	111
12.3.2	Operations: size() and is_empty()	111
12.3.3	Operation: push()	112
12.3.4	Operation: pop()	112
12.3.5	Operation: peek()	112
12.3.6	Summary	113
<b>13</b>	<b>The Node ADT</b>	<b>115</b>
13.1	Nodes	115
13.2	Node chains	117
<b>14</b>	<b>Programming practice and style</b>	<b>121</b>
14.1	Introduction: Style counts	121
14.2	Style Guidelines	122
14.3	Improving your code	125
14.4	Making your program faster	130
14.5	How to improve your programming skills	130
<b>15</b>	<b>Defensive Programming</b>	<b>133</b>
15.1	What could possibly go wrong?	133
15.2	Defensive Programming	133



<b>16</b>	<b>Node-based Stacks and Queues</b>	<b>141</b>
<b>16.1</b>	<b>Why use nodes to implement Queues and Stacks?</b>	<b>141</b>
<b>16.2</b>	<b>Using Nodes to to Implement Stacks</b>	<b>141</b>
16.2.1	Operation: Stack()	142
16.2.2	Operations: size() and is_empty()	142
16.2.3	Operation: peek()	143
16.2.4	Operation: pop()	144
16.2.5	Operation: push()	145
<b>16.3</b>	<b>Using Nodes to Implement Queues</b>	<b>145</b>
16.3.1	Operation: Queue()	146
16.3.2	Operations: size() and is_empty()	147
16.3.3	Operation: peek()	148
16.3.4	Operation: dequeue()	148
16.3.5	Operation: enqueue()	150
<b>16.4</b>	<b>Summary</b>	<b>151</b>
<b>17</b>	<b>Linked Lists</b>	<b>153</b>
<b>17.1</b>	<b>Linked List ADT</b>	<b>153</b>
<b>17.2</b>	<b>Linked List Record</b>	<b>154</b>
<b>17.3</b>	<b>Linked List Operations</b>	<b>154</b>
17.3.1	Operation: create()	154
17.3.2	Operations: size() and is_empty()	155
17.3.3	Summary of Linked List Operations	155
17.3.4	Implementation details	156
17.3.5	Perspective	156
<b>18</b>	<b>Algorithm Analysis</b>	<b>159</b>
<b>18.1</b>	<b>The importance of understanding costs</b>	<b>159</b>
18.1.1	Abstraction #1: Time	160
18.1.2	Abstraction #2: Input size	161
18.1.3	Abstraction #3: Categories	161
18.1.4	Abstraction #4: Worst case behaviour	162
<b>18.2</b>	<b>Asymptotic analysis</b>	<b>163</b>
18.2.1	Basics of Working with Big-O	163
18.2.2	Counting Steps in Algorithms	164
<b>18.3</b>	<b>Examples</b>	<b>166</b>
18.3.1	Step analysis	166
18.3.2	Method 2: Simplified Steps	166
18.3.3	Logarithmic loops	167
18.3.4	Nested loops with independent limits	167
18.3.5	Nested loops with dependent limits	168

<b>18.4</b>	<b>Warnings and Advice</b>	<b>168</b>
<b>19</b>	<b>Recursion</b>	<b>169</b>
<b>19.1</b>	<b>Introduction</b>	<b>169</b>
19.1.1	Terminology	170
19.1.2	The Delegation Metaphor	171
19.1.3	Avoid the rabbit hole	171
<b>19.2</b>	<b>How to Design a Recursive Function</b>	<b>171</b>
19.2.1	The recursive case	172
<b>19.3</b>	<b>Examples</b>	<b>173</b>
<b>19.4</b>	<b>Functions, recursion, and the call stack</b>	<b>174</b>
<b>19.5</b>	<b>Algorithm analysis for recursive functions</b>	<b>176</b>
<b>20</b>	<b>Trees and Binary Trees</b>	<b>179</b>
<b>20.1</b>	<b>Trees organize data hierarchically</b>	<b>179</b>
<b>20.2</b>	<b>Informal Terminology</b>	<b>180</b>
<b>20.3</b>	<b>Technical definitions</b>	<b>181</b>
<b>20.4</b>	<b>Binary trees</b>	<b>182</b>
<b>20.5</b>	<b>Binary Treenode ADT</b>	<b>183</b>
<b>20.6</b>	<b>Trees from treenodes</b>	<b>184</b>
<b>21</b>	<b>Algorithms on Binary Trees</b>	<b>187</b>
<b>21.1</b>	<b>Sequences and Tree Traversals</b>	<b>187</b>
21.1.1	Breadth-order sequence	187
21.1.2	Pre-order sequence	188
21.1.3	In-order sequence	188
21.1.4	Post-order sequence	188
<b>21.2</b>	<b>Tree Traversals</b>	<b>189</b>
21.2.1	Pre-order traversal	189
21.2.2	In-order traversal	190
21.2.3	Post-order traversal	190
21.2.4	Breadth-order traversal	191
<b>21.3</b>	<b>Algorithms based on Tree Traversals</b>	<b>192</b>
21.3.1	Calculating the height of a binary tree	192
21.3.2	Counting the number of nodes in a tree	193
21.3.3	Searching for a value in the tree	194
21.3.4	Generalizations about functions on trees	195

<b>22</b>	<b>Binary Search Trees</b>	<b>197</b>
22.1	Introduction	197
22.2	Binary Search Trees	198
22.3	Basic operations	199
22.3.1	Checking if the value is in the BST	199
22.3.2	Adding a value to the BST	200
22.3.3	Removing a value from the BST	202
22.4	The importance of Binary Search Trees	203
<b>23</b>	<b>The Table ADT</b>	<b>205</b>
23.1	Introduction	205
23.2	Table ADT implementation	206
23.3	The efficiency of the Table ADT implementation	207
<b>24</b>	<b>Object Oriented Programming</b>	<b>209</b>
24.1	Mutable vs Immutable	209
24.2	Inheritance: Giving several classes some common properties	212
24.3	Over-riding a method	216
24.4	Summary	217
<b>25</b>	<b>Three Kinds of Tasks</b>	<b>219</b>
25.1	Introduction	219
25.2	Example Problems	220
25.2.1	Subset Sum	221
25.2.2	Maximum Slice	221
25.2.3	Making Change	221
25.2.4	Maximum Tree Path	221
25.2.5	Leap Line	222
<b>26</b>	<b>Algorithms</b>	<b>223</b>
26.1	Introduction	223
26.2	Algorithm styles	224
26.2.1	Brute Force Algorithms	224
26.2.2	Backtracking Algorithms	224
26.2.3	Divide and Conquer Algorithms	224
26.2.4	Greedy Algorithms	225
26.2.5	Dynamic Programming Algorithms	225
26.3	How to solve it	226

<b>27</b>	<b>Algorithm Examples</b>	<b>227</b>
27.1	Introduction	227
27.2	Subset Sum	227
27.2.1	Brute Force	227
27.2.2	Divide-and-conquer	229
27.3	Maximum Slice	230
27.3.1	Brute Force	230
27.3.2	Divide and conquer	232
27.4	Make Change	234
27.4.1	Greedy	235
27.5	Maximum Tree Path	236
27.5.1	Brute force	236
27.5.2	Greedy	236
27.6	Leap Line	238
27.6.1	Brute Force	238
27.6.2	Dynamic Programming	239
27.6.3	Greedy	240
<b>28</b>	<b>Case Study: Huffman codes</b>	<b>243</b>
28.1	Introduction	243
28.2	Using frequency information to design a code	245
28.3	Encoding and decoding	247
28.4	Efficiency considerations	248
<b>A</b>	<b>Abstract Data Types</b>	<b>253</b>
A.1	The full ADT Registry implementation	253
A.2	The full ADT Statistics implementation	254
<b>B</b>	<b>Applications of Stacks and Queues</b>	<b>257</b>
B.1	Queueing simulation	257
B.2	Bracket checking	259
B.3	Post-fix Arithmetic evaluator	260
<b>C</b>	<b>List-based Stacks and Queues</b>	<b>263</b>
C.1	The full list-based Queue implementation	263
C.2	The full list-based Stack implementation	265

<b>D</b>	<b>Node-based Stacks and Queues .....</b>	<b>267</b>
D.1	The full list-based Queue implementation	267
D.2	The full list-based Stack implementation	270
<b>E</b>	<b>The treenode ADT .....</b>	<b>273</b>
E.1	The full Treenode ADT implementation	273
<b>F</b>	<b>Algorithms .....</b>	<b>275</b>
F.1	Maximum Slice	275
F.2	Subset Sum	279
F.3	Maximum Tree Path	285
F.4	Make Change	288
F.5	Leap Line	290





# 1 — Computers and Software

## Learning Objectives

After studying this chapter, a student should be able to:

- Describe the different components of a computer system.
- Explain how the components of a computer system interact.
- Describe the operational behaviour of the machine cycle.
- Explain the difference between machine language and programming language.
- Explain the difference between a compiled language and an interpreted language.
- Explain the difference between an Interactive Development Environment, a compiler, and an interpreter.
- Describe the purpose of a computer operating system.
- Describe the distinguishing features of interactive and non-interactive applications.

We'll spend a great deal of time concerned with software in this course, but it will benefit us to have some context. In this chapter, we'll cover some of this background, starting with a simple review of computer architecture, before proceeding to talk about software. There are places in this chapter where the facts have been simplified to aid in creating a consistent and general perspective on the material.

## 1.1 Computers

A modern computer is a complex electronic machine, consisting of electronic circuits and wires, with other components made from glass, metal, plastic, and ceramics. The term computer is quite flexible: we use the term for desktop computers, notebook computers, super-computers, smartphones, and tablets. Embedded computers can also be found in mechanical systems like cars and refrigerators and vacuum cleaners. Figure 1.1 shows a conceptual diagram of a modern computer. The story we tell in this section, like the figure, is highly simplified, but conceptually similar for all of these different

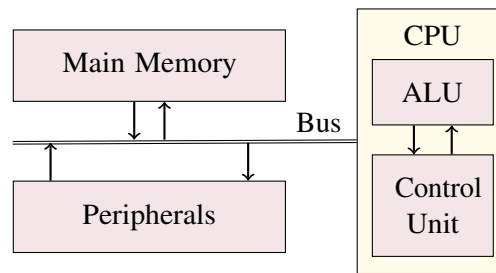


Figure 1.1: An abstract representation of a computer system.

kinds of devices. This story, even in simplified form, will help us understand the relationship between our main purpose, writing software, and the hardware that the software runs on.

### Main memory

The *main memory* is where applications and data are stored when the applications are running. The computer stores all data (including programs and applications) in a number of different binary codes. Fundamental to all computer systems is the fact that the basic memory unit can represent 2 values, either 0 or 1, and this unit is called a binary digit, or *bit*. Most modern computers organize bits into 8-bit groups called *bytes*; modern computers organize main memory as a very long sequence of bytes (see Figure 1.2). Each byte in the sequence has a numeric *address* that can be used by the computer to access the data stored in the byte very quickly. The address itself is not actually stored in memory; it is implicit in the design of the memory circuitry. Data gets transferred from place to place in a computer using a system called the *data bus*. In a modern computer, the data bus is 8 bytes wide (64 bits), meaning that 8 bytes of data can be transferred in a single action; in older computers, and in some embedded applications, the bus might be 1, 2 or 4 bytes wide. The bus connects the main memory, the peripheral devices, and the CPU.

### The Central Processing Unit

The *central processing unit* (CPU), is where most of the computation happens. It's made up of an *arithmetic logic unit* (ALU) and a *control unit*. The ALU knows how to do simple arithmetic operations on integers and floating point numbers, how to make comparisons to see which of two numbers is larger, and how to perform boolean operations like AND, OR, and NOT. The control unit contains the circuitry for coordinating the execution of a computer program; it uses the bus to fetch data and program instructions from main memory as needed, feeds data to the ALU to perform computation, and obtains the results of such computations, possibly also storing them back into main memory.

A very small amount of data can be stored right inside the CPU in extremely high-speed memories called *registers*. Registers can be used as temporary storage for small amounts of data that will be used in the very immediate future to avoid the relatively high time costs of communication with main memory. Different CPU designs have different numbers and sizes of registers, but there are some special-purpose registers that are common to all CPUs (though they may go by different names). One of these common registers is the *instruction pointer* (IP); also called the *program counter* (PC). Another is the *instruction register* (IR).

The CPU is commonly called the brain of a computer, but it might be better to call it the engine: all the components of a computer operate like electronic clockwork, and the CPU drives it all. The CPU performs a very simple algorithm called the *machine cycle* consisting of three steps:

1. The CPU *fetches* an instruction, by sending a signal to main memory asking for the instruction stored at the memory address contained in the IP register. This instruction is sent along the data bus from memory and stored in the CPU's instruction register (IR). During the time that the data was travelling over the bus, the CPU updates the instruction pointer, to contain the address of the next instruction.
2. Once the instruction register contains the instruction, the CPU *decodes* it, which means that the CPU uses the bit-patterns in the instruction to activate the appropriate circuits in the ALU, or sometimes, in the control unit itself.
3. When the correct circuits are ready to go, the CPU calls for the *execution* of the circuit: electronic signals pass through the circuit.

The CPU's only work is to repeat these three steps, which it can perform at extremely high rates, because they are encoded by the circuitry of the control unit.

Since the invention of the microprocessor, CPUs have been getting smaller, faster, and cheaper. For most of the history of computers, a CPU had only one control unit, and one ALU. In the past 15 years or so, CPUs have been built with multiple so-called "cores." Each core has a control unit, and an ALU, some registers, and a larger amount of memory called a cache. In principle, each core could be running a different application independently; in practice, some coordination between the cores is necessary, because all the cores still share the same data bus, the main memory, and all the peripherals. A lot of sophisticated hardware and software design contributes to a computer system making effective use of multi-core CPUs.

### Peripherals

Disk drives, keyboards, monitors, and touch screens are probably essential to a modern understanding of computer use, but these are represented abstractly in Figure 1.1 as *peripherals*. A peripheral device, like a disk drive or a graphics card, typically has its own control unit, some amount of memory that can only be used by its control unit, and some hardware connector to a physical device. Data can enter the computer through a peripheral device, like a keyboard or a network connection, or data can be sent from the computer to a peripheral device, like a monitor or a network connection. This data gets transferred to or from the computer's main memory using the data bus. Peripheral devices are designed to be able to act somewhat independently of each other, and the CPU. For example, the network connection listens for incoming data, and sends outgoing data; the computer's CPU is not listening directly. Similarly, the graphics card processes any data that a program wants to display; the computer's CPU is not doing the graphics processing. This division of labour means that the CPU is not doing everything itself. But all peripherals take direction from the computer's CPU.

## 1.2 Languages for Computation

Computers are useful because they can be programmed to perform different computational tasks. This statement warrants reflection: a computer's behaviour can be changed by software, which is basically encoded information. Consider, in contrast, an internal-combustion engine: no amount of information can get it to behave like a printing press; this change in behaviour can only be achieved by making drastic physical changes to the machine. To communicate the necessary information that computers need to perform different tasks, we write software, and writing requires a language.

Address	Data
000	11011001
001	00100101
002	10010000
003	00000000
004	11000010
005	11001100
006	01010101
007	10101010
⋮	⋮

Figure 1.2: Organization of main memory. Memory units consisting of 8-bit bytes are individually addressable.

### Machine Language

A computer application (also called an *executable*) consists of a sequence of machine instructions, represented as binary data. Each instruction is typically very simple. For example:

- A *load* instruction retrieves data from a given address in main memory, and stores it in a given register.
- An *add* instruction adds the data in 2 given registers together, storing the result in a third register.
- A *store* instruction sends data from a given register out to a given address in main memory.

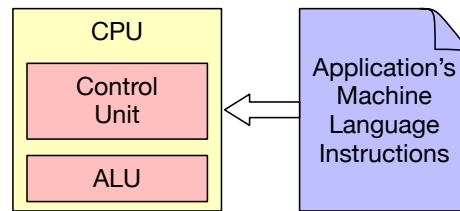
There are other kinds of instructions that affect the control unit itself. For example, a *jump* instruction tells the control unit to change the value stored in the instruction pointer; this instruction changes which instruction is loaded into the instruction register next.

One of the fundamental concepts in computer science is that *there is no real difference between program and data*. In fact, if you were to look at any sequence of bytes in main memory, you would have a very hard time figuring out whether you were looking at numeric data, text data, or machine instructions. A machine language program consists of machine language instructions, which are fetched, decoded, and executed by the CPU. In other words, the CPU's machine cycle algorithm treats the machine language instructions of an application as data. It is a useful abstraction to treat the machine language instructions as if they were controlling the CPU, effectively creating a specialized machine for a specific purpose. Consider Figure 1.3, showing that a machine language application interacts with a computer. The diagram uses an arrow, which suggests that the application is in control of the machine, but another way to interpret this diagram is to understand that the machine is taking machine language instructions as data to inform its behaviour. Both points of view are valid, and useful.

While it is possible for programmers to write applications in the language of machine instructions, it's almost never done. Machine language is, literally speaking, expressed in the zeroes and ones of a binary code. To make machine language slightly easier to work with, each machine language instruction code is given a short English label: *jump* or *load*, as mentioned above. This language is sometimes also called machine code, because there is a one-to-one translation between these English labels and the zeroes and ones of the machine code. Strictly speaking, a program expressed in terms



Figure 1.3: A computer changes its behaviour based on the instructions provided by the machine language application.



of English labels is called *assembly language*, and an application called an *assembler* is used to translate the assembly language code into machine code.

For computer engineers, and low level system software developers, it is still necessary and useful to write relatively small control programs or device drivers in assembly language. Using assembly language, these experts can exert the finest control over the hardware, allowing the computer to interact with peripherals in the most efficient way possible. However, working in assembly language requires the programmer to break their algorithms down into a relatively large number of machine instructions. Assembly language has none of the built-in conveniences for programmers that you may have learned in a previous programming course: no variables, no arrays, no lists, no dictionaries, no strings; it uses memory locations and addresses directly. It has no if-statement or while-loops; it uses jump instructions and other similar kinds of instructions directly. These conveniences can be programmed in assembly language, but doing so increases the programmer's burden substantially. The number of details that an assembly language programmer needs to keep track of is quite high, and it's very easy to make mistakes. Furthermore, applications written in assembly language can only be executed on computers that have the same CPU; this is a severe disadvantage, considering the number of different kinds of computers, and the rate at which computer hardware is updated.

### Programming Languages

Because of the difficulties of working in assembly language, computer scientists invented programming languages, which allow programmers to express their ideas in terms of larger steps, with more structure. Programming languages provide high-level abstractions like variables, loops, if-statements, and various compound data structures like lists, arrays, strings, and dictionaries. The rather large advantage is that the programmer can express algorithmic concepts more easily in a programming language than they could in assembly language or machine language. One minor disadvantage is that the programmer and the computer are no longer speaking the same language. A far more serious disadvantage is that a programmer might not have a firm understanding of what a computer is actually doing. This is a source of inefficiency and error!

There are hundreds of different programming languages being used by programmers. There are programming languages for different purposes, for different hardware systems, and for different users with different goals. The best programming language for a task depends on the task itself. A computer engineer might prefer to use C and assembly language. A data scientist might prefer to use Python. A statistician might prefer to use R. Engineering professionals might prefer to use MATLAB. Large commercial applications that require sophisticated software engineering might be programmed in Java, C++, or C#.

All programming languages have to communicate algorithms, to be executed on computer

hardware, so there are necessarily some common fundamental concepts, such as data, conditionals, repetition, etc. The key consequence of this observation is that learning a second programming language is not really as daunting as learning to program. It's common for novices to conflate the idea of learning a programming language with the idea of learning to program, because their first experience is learning both at the same time. On the contrary, once a novice knows what it means to express an algorithm in computational terms, picking up a second (or twentieth) language is a lot easier, because programming is the same; only the language details are different. In other words, each language you learn will be easier to learn, because you already know the deeper fundamentals of programming already.

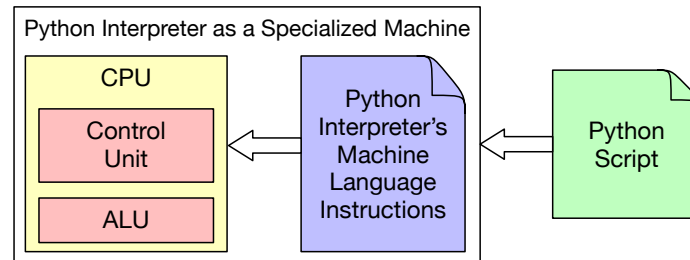
The wide variety of programming languages is a consequence of the active area of research in Computer Science into the design and implementation of programming languages. The main issue in programming language research is to provide tools and abstractions so that programmers can control the complexity that arises in designing and implementing large software applications. These abstractions make some programming tasks simpler, but consequently, they make other tasks more difficult. For example, assembly language provides very minimal abstraction, so even simple programming tasks can take a long time to write. On the other hand, Python provides a lot of abstractions (lists, dictionaries, for-loops, etc) which make simple programs a lot easier to write. The down-side for Python is that it is very difficult to access the computer hardware's abilities directly. A programming language like C/C++ gives a middle ground, allowing greater control of the computer more directly than Python, but using fewer abstractions. A C/C++ program is likely to be four times longer than an equivalent Python program, but might run ten times faster.

### Compilers and Interpreters

A *compiler* is a computer application whose job is to take a program (that is, a text document containing instructions written in a programming language) as input and translate it to machine language, producing an application (also called an executable file) as output. For example, a programmer can write a program using the C programming language, and the compiler would analyze the program, and translate it into machine code, saving the results in a document called an *executable file*. The compiler's input is the program (a text document), and its output is the executable; it's important to understand that the compiler does not actually cause the application to be executed. Once the program has been compiled, the resulting executable can be executed directly by the computer, as often as needed without compiling it again. In principle, a program written in the C programming language could be compiled and executed by any other computer system that has a compiler for the C language; in practice, porting a program from one system to another not quite as easy as it seems it should be. Compiler design, and compilation is so important to the field that it has become a major sub-field of computer science. Without the enormous effort invested into compiler design and the development of compilers, we would not have computers in our pockets today.

An *interpreter* is a computer application that executes a script (that is, a text document containing instructions) one line at a time. The interpreter is an executable application, running directly on the computer system. Its input is a script, but an interpreter does not create an executable like a compiler does. Instead, the interpreter analyzes each line of the script, and gets the computer to execute it immediately, before moving to the next line. For example, when we run a Python script, it may seem as if the computer is running it, but what's actually happening is that a Python interpreter executes our script. See Figure 1.4. The Python script cannot be executed directly by the computer, and if we want to run the script again, it has to be interpreted by the Python interpreter again. Some interpreters can be used interactively, allowing a user to type commands in the language, and see the

Figure 1.4: A Python interpreter is a machine language application that controls the computer system. This application can make the computer behave like a specialized machine for running Python scripts. A Python script informs the interpreter, and the interpreter informs the machine.



results, one instruction at a time.

A *debugger* is basically an interpreter that allows interactive exploration of the script's variables and data, as the script is being executed interactively. If your program has an error that you can't find, it's very useful to be able to watch the execution of your program step by step, and compare the actual effect of each instruction against the effect intended by the programmer. Interpreted languages are highly portable: applications written in an interpreted language (like Python) can be executed by any computer system that has the appropriate interpreter written for it. The disadvantage of interpreted applications is that they are usually noticeably slower than compiled applications using the machine's own language.

Some programming languages combine the compiled and interpreted approaches. The Java language has a compiler that translates Java programs to a form that looks like machine code; they call this Java byte-code. However, there is no CPU that runs the Java byte-code directly. Instead, there is an interpreter, called a Java Virtual Machine (JVM), for the compiled Java byte-code; this interpreter executes the compiled Java program by pretending to be a computer with a different machine language. This approach is useful because Java byte-code can be executed on any computer system that has a JVM written for it, making Java applications very portable. Furthermore, compiling to byte-code means that Java applications are faster than applications written in languages that are interpreted directly (but more slowly than applications written in the computer's own machine code). There are some Python interpreters that take the same approach as Java, using a virtual machine that can execute byte-code produced by the interpreter.

An *interactive development environment*, or IDE, is an application that coordinates the interaction between a number of other applications, including an application that allows programmers to edit programs, and other applications like compilers, interpreters, and debuggers. The IDE is something like the conductor of a symphony, making sure that there is coordination between the associated applications. It's important to understand that the IDE is not itself an interpreter or a compiler; the IDE only tells the computer to execute the interpreter, and can tell the interpreter which script to execute.

## 1.3 Software

### Operating Systems

An operating system is a collection of highly integrated and interdependent control software that provides a programmer convenient and useful access to the computer's hardware resources (CPU, graphics, network, mouse, keyboard, disk drives, etc). For the computer user, the operating system is necessary to support the basic use of the system, as well as any software that they might use.

The operating system allows applications to share the hardware resources, so for a software developer or a computer engineer, the appreciation for an operating system is much less superficial. Consider the following examples of the support provided by a modern operating system:

- The operating system allows multiple applications to run at the same time. This is called *multi-tasking*.
- The operating system is responsible for providing organizational tools for the storage of data files on peripherals like disk drives.
- When an application is started, the operating system allocates memory for the application to use, copies the application executable code to the allocated memory, and starts executing it.
- When multiple applications are sending and receiving information (from the Internet, or other communications), the operating system makes sure that each application's requests get sent out fairly, and that the requested data gets delivered to the application that requested it.
- When multiple applications have windows open on the screen, the operating system takes care of the way the windows work, including the window controls, and each application only has to be concerned with making the right information appear within the window.
- When multiple applications are reading data from or writing data to the disk, the operating system makes sure the data is transferred correctly and efficiently.

Without an operating system, every application program would take direct control of the whole computer for as long as the program was running. Furthermore, every programmer would have to figure out how to do common tasks on the computer, such as display information to a screen, read data from a file, or send data to another computer connected by a cable, or by wireless communications. Clearly this would be so much trouble and so prone to errors that computer systems would hardly be practical, and certainly, you wouldn't have one in your pocket if programmers had to contend with the hardware without an operating system. See Figure 1.5.

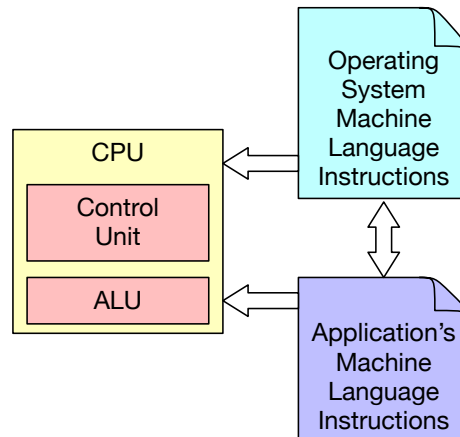
Most casual computer users identify the visual aspects of a computer system (windows, menus, touch screen, etc) as the essential aspect of an operating system, because that's what they see when they use the computer. This visual aspect is important for branding, but for well-trained software professionals, the branding is less important than the consequences of the operating system software each brand provides. The most popular computer operating systems include Microsoft Windows and Apple macOS.<sup>1</sup> Because the operating system software is different, applications that run on one of these systems usually cannot run on the other. So if a software application is released for both systems, either the developers have implemented two different versions, or the application is running in an interpreter or virtual machine available for both systems. Designing and maintaining two (or more) versions is very expensive.

A general operating system design called UNIX was developed in the 1970s for main-frame computers at a time when dozens of users had to share a single computer. The popularity of UNIX among software professionals, starting from the 1970s, resulted in the development of thousands of

---

<sup>1</sup>These names are trademarks owned by their respective owners.

Figure 1.5: A machine language application interacts with the machine for some of its work, but calls the operating system software to interact with shared resources, like disks, peripherals, etc. This balance makes compiled applications very fast, but also allows them to work cooperatively when many applications are executing at the same time.



applications written by programmers for programmers, as part of the UNIX toolbox. As desktop computers became more powerful, UNIX migrated to the desktop; currently Apple's macOS and various forms of Linux are the most well-known UNIX systems, but there are others. Part of our work in this course will be to begin to use some of these tools.

### Software applications

Modern software applications take a wide variety of forms. Most students have experience with interactive graphical software, including Microsoft Office applications, web-browsers, and computer games. These are called *graphical* because they communicate with users via visual tools, including icons, figures, images, and animation. Applications are *interactive* if they primarily operate by performing tasks based on the user's actions, including actions of the computer's keyboard, mouse, touch screen, and more recently, voice. Interactive applications that use graphical tools are commonly called Graphical User Interface (GUI) applications.

There are other forms for applications. Interactive applications can be built entirely in terms of text-based communications. Such applications are often called Console Applications. The console application can communicate with the user by displaying text in a console window, and obtaining user input from the keyboard. We will focus on Console Applications in this course. Building a GUI application requires concepts that will come up in later courses.

Interactive applications are good for many purposes, but not all; professionals who are well-trained in their computer use will often prefer limited interactivity, just for efficiency. In previous courses, you may have made every program or script an interactive application: politely greeting the user when the program starts, and asking for various data, before outputting some kind of answer. This is a good way to get novices to think about their homework, but if you are using a computer to get work done, sometimes it's very effective to minimize the interactivity. A lot of the time professionals will want to build a program that acts like a tool, not a game, based on input that comes from a data file, or from some initial condition provided by the user. The advantage is that we can



run the tool many times, and we will not be impeded by friendly banter at the console.

## 1.4 Summary

In this chapter, we have provided a simplified view of the organization of hardware in a computer. A computer is made up of a Central Processing Unit (CPU), which communicates with memory and peripheral devices using a bus. The CPU's hardware implements a very simple algorithm called the machine cycle, which fetches machine code instructions, decodes them, and then executes the instructions. Machine code instructions are very simple operations that can be performed very quickly. An application is basically a sequence of these simple instructions. This simplified description is not exact, but it is useful as a model to understand how software works. There are other courses you can take to learn more about how computer hardware is organized.

Applications are executable files containing machine language instructions, but they are usually written in a programming language, not machine code. Assembly language is basically a language that gives each machine instruction an English word, like *load*, *store*, or *add*. It's considered a low-level language because each part of the language corresponds closely to something that the CPU can do directly. High level programming languages allow programmers to express ideas in larger concepts, using variables, loops, if-statements, etc. Programs written in high-level languages express need to be translated to machine code, either by an interpreter or a compiler. An interpreter is an application that reads and executes a program file one line at a time; the interpreter is executed by the CPU, but the program is executed by the interpreter. A compiler is an application that reads the whole program file, and translates the program into an application without running the program file. Once the application is created, it can be executed directly by the CPU. An interactive development environment (IDE) is an application that coordinates the tools that a software developer uses in daily work, such as compilers, editors, and interpreters.

An operating system is a collection of highly integrated and interdependent control software that provides a programmer convenient and useful access to the computer's hardware resources (CPU, and peripherals). It's the operating system that allows applications to share these resources fairly and conveniently. When programs try to send output to the console, or read data from a file, and many other similar kinds of activities, they are actually interacting with the operating system, which provides these abilities. High-level languages provide abstractions that hide most of these details from the programmer.

A computer application can be interactive, meaning that the user participates in the application's work, either by using a keyboard, mouse, or touch screen to direct the application's activity. Sometimes it's better for an application to be non-interactive, and just get some work done.

### Two example problems, with solutions

- Sorting a list
- Counting prime numbers

### A grab-bag of other listy stuff

- Perspective on lists

### Dictionaries

- An example program

### Modules

### Three Example Problems

- The Gambler's Ruin Problem
- The Coupon Collector Problem
- Self-Avoiding Random Walks

## 2 — Computing with Python

### Learning Objectives

After studying this chapter, a student should be able to:

- Make use of lists to solve computational problems.
- Make use of dictionaries to keep data organized.
- Import Python modules and use functions defined in them.

In this first chapter, we'll do a bit of review of lists, dictionaries, objects, and modules. The purpose is to give students a sense of where they should be, as they are starting this course.

### 2.1 Two example problems, with solutions

Python lists have a very high degree of utility, and it would be rare to do much programming without using them. We'll spend a significant amount of time working with lists, so it's good that we start with a review.

In Python, a list is a compound data type used to organize data into linear sequences. To access a data value in a list, you need to know its *index*, i.e., its position in the list, represented by an integer. The first value or item in a list has index 0 (zero), and the last item's index is one less than the number of items in the list. For example, if a list contains 10 items, the index of the last item is 9.

#### 2.1.1 Sorting a list

To get things started, consider the problem of sorting a list. Specifically, we assume that we have a list named `unsorted`, which has numeric data in it, and we want to create a new list, named `sorted`, containing all of the values in `unsorted` but in increasing order. For example, if we start with the list `unsorted = [3, 1, 2]`, after sorting, the list `sorted` would have the values `[1, 2, 3]`. The short Python program below accomplishes this task.

```
1 unsorted = [3, 2, 5, 7, 6, 8, 0, 1, 2, 8, 2]
```

```
2 sorted = list()
3
4 while len(unsorted) > 0:
5     out = min(unsorted)
6     unsorted.remove(out)
7     sorted.append(out)
8
9 print(sorted)
```

Notice that Python lists can be created by using a literal expression, e.g. Line 1, or by using the function `list()`; on Line 2, we create an empty list, by giving the function no arguments. The assignment statements on Lines 1 and 2 also demonstrate that lists can have names; more formally, we say that the variables `unsorted` and `sorted` each refer to a specific list. On Lines 4 and 5, we see that Python has functions that can be applied to lists: `len()` returns the number of values in the list, and `min()` returns the smallest value in the list. On Lines 6 and 7, we see that a list is actually an object, having a method called `remove()` that removes a given value from the list, and `append()` which puts a given value after all the current values in the list.

The program itself is an algorithm that can be summed up as follows: *Keep on removing the smallest value from `unsorted` and adding it to end of `sorted` until there are no more values left in `unsorted`.* There are a number of questions you might think about while considering this program. Is it a good algorithm? By what standards can we consider an algorithm to be *good*? Does it actually work? How would we be certain?

There are a number of questions you might think about while considering the use of Python in this program. What does it mean that a variable *refers* to a value? What's the difference between an *object* and a *value*? Why is `len()` a function, but `append()` is a method? And why is one in blue text colour, but not the other?

You may be aware that Python has a built-in function whose name is `sorted()`, which takes a list argument and produces a new list with all the values in ascending order. While it is almost obviously easier to use a built-in function than to write your own code, is it always *better*? And what happens to the function `sorted()` when a program creates a variable of exactly the same name?

### 2.1.2 Counting prime numbers

The question we want to answer here is this: given a positive integer  $n$ , how many prime numbers are there that are less than or equal to  $n$ ? A *prime number* is an integer greater than 1, that cannot be divided evenly (i.e., without a remainder) by any other integer except 1 and itself. For example, 2, 3, 5, and 7 are prime, but 4, 6, 8, and 9 are not. So if we pick  $n = 11$ , the prime numbers are {2, 3, 5, 7, 11}, so in total, there are 5 primes less than or equal to 11.

There is no nice simple formula for this, so we have to calculate it by counting. One way to count primes is to step up from 2 to  $n$ , asking whether each number is prime or not. The Python script below accomplishes this more or less directly.

```
1 n = 100000 # end of range of numbers to check for primes
2
3 count = 0
4 for i in range(2, n + 1):
5     no_factors_found = True # assume prime until disproven
```

```
6     f = 2
7     # check if i is prime by checking remainders
8     while no_factors_found and f < i:
9         if i % f == 0:
10             no_factors_found = False
11             f += 1
12
13     if no_factors_found:
14         count += 1
15
16 print("# Prime numbers between 2 and " + str(n) + ":", count)
```

In this program we made use of the `range()` function, and the remainder operator `%`. Before going on, you might consider the same kinds of questions as before: Is it a good algorithm? Does it actually work?

Another way to count primes makes use of a list that records whether or not a number is prime. The list records Boolean values, one for each integer from 0 to  $n$  (though we won't be looking at 0 or 1, which are not prime by definition). The list will start with all entries `True`, i.e., the list will assert that until proven otherwise, all the numbers are possibly prime. We'll walk up through the list, starting at 2. We know that 2 is prime, so we'll keep the list entry as `True`. But all the multiples of 2, namely 4, 6, 8, etc., are all not prime, so we'll jump through the list recording the value `False` for all the multiples of 2. In this way, we've marked the numbers we've definitely proven to be not prime, and we're left with numbers that might still possibly be prime. This technique is implemented below:

```
1 n = 20 # end of range of numbers to check for primes
2
3 still_is_prime = (n+1)*[True] # assume prime until disproven
4
5 for i in range(2, n):
6     if still_is_prime[i]:
7         # mark multiples of i as not prime
8         j = 2*i
9         while j <= n:
10             still_is_prime[j] = False
11             j += i
12
13 # now, every possible prime is a definite prime
14 count = sum([1 for v in still_is_prime[2:] if v])
```

In this example, we created a list by multiplication; on Line 3, we created a list by joining together  $n+1$  copies of the list `[True]`. Lists can also be *concatenated* or joined together using `+`, e.g. `[1,3,5] + [2,4,6]`.

We made use of a very important aspect of lists, namely the use of indexing (Lines 6, 9).

Notice Line 12, which uses a slice to drop the first two elements: `still_is_prime[2:]`, as well as a comprehension, to form a list containing one value 1 for each `True` in the slice. The Python function `sum()` adds up all the values in a given list; we used it instead of a for-loop to count how many `True` values are in the list.

The program itself is an implementation algorithm named after an ancient Greek mathematician. This algorithm, known as “The Sieve of Eratosthenes,” can be summed up as follows:

Sift the Twos and Sift the Threes,  
The Sieve of Eratosthenes.  
When the multiples sublime,  
The numbers that remain are Prime.  
– Anonymous

The algorithm was known millennia before computers were invented.

Before you skip ahead to the next section, you might consider whether Eratosthenes’ algorithm is better than the previous algorithm, and how you might make that assessment one way or the other.

## 2.2 A grab-bag of other listy stuff

A list comprehension is a very useful way to build a list, and is reminiscent of mathematical set theory. Consider the set of all the even integers in the range from 1 to 9. In Python, you can create this list a couple of ways:

```
>>> [x for x in range(1,10) if x % 2 == 0]
[2, 4, 6, 8]
>>> [2*x for x in range(1,5)]
[2, 4, 6, 8]
```

If you stare closely enough at these Python expressions, you might see the resemblance to mathematical set notation, as follows:

$$\{x \mid x \in \{1, \dots, 9\}, x \bmod 2 = 0\}$$

$$\{2x \mid x \in \{1, \dots, 4\}\}$$

It’s important to say that lists are not, strictly speaking, identical to sets. A list can contain repeated values, but a set, strictly speaking, does not have repeated values.

Lists can contain anything at all, even though we’ve only used integers. A list can contain all of the same kind of thing, e.g., a list containing only integers, or different kinds of things, like a list containing numbers and strings.

Python lets you extract parts of a list using a technique called *slicing*, as we saw earlier. For example, the expression `L[2:5]` creates a new list consisting of the third through fifth items in a previous list `L`. Note that the first index in a slice is inclusive, but the second is not inclusive. The expression `L[0:len(L):2]` creates a new list consisting of all the items in the list `L` stored at even indices (in English, one way to say this is “every other item,” but this is slightly ambiguous). Slicing is so useful that it can be applied to any kind of sequence, including tuples, and strings.

The Python for-loop allows a program to step through a list in two ways. The simplest is to use each item directly, such as:

```
L = [1,3,5] + list(range(2,7,2))
total = 0
for item in L:
    total = total + item
```



Here, the variable `item` is bound to each item in `L` in sequence. A for-loop like this can be very fast in Python. But if you're summing a list, as in the above example, it's not as fast as the function `sum()`.

The word `in` can be used as a way to test if a given value is stored in a given list. For example:

```
L = [1,3,5] + list(range(2,7,2))
if 5 in L:
    print('Yes')
else:
    print('No')
```

The expression `5 in L` is a relational expression that is `True` if `L` contains the value 5, and `False` otherwise. You can combine `in` with `not`: The expression `5 not in L` is a relational expression that is `False` if `L` contains the value 5, and `True` otherwise.

A list is mutable, which means that we can change its contents, as follows:

```
L[2] = 33
```

We can add items to the list using methods `extend()` and `append()`, as in the following examples:

```
L = [1,3,5]
L.extend(list(range(2,7,2)))
L.append(77)
```

In this example, we are using the list `L` as an object, by giving the name of the object `L`, the dot, and then something that looks like a function call. The method `extend()` adds a sequence of items to the end of the list; the method `append()` adds a single item to the end of the list. Other important methods of this sort include `pop()`, `remove()`, and `index()`, among many others.

List objects have lots of methods, which may be useful in various ways. One of the most important things to learn about Python is that it is very common for programmers to look up the details about any particular function, object, or method. Some you will learn simply by re-use. Others you will need to look up quickly to remind yourself how to use them. There is a standard website for Python information: <https://docs.python.org/3/>, which should be one of the first places you check.

It's very common for a list to contain other lists. For example:

```
L = [[1], [2,3], [4,5,6], [7,8,9,10]]
print(L[2], L[3][1])
```

In the above example, `L` is a list containing 4 sub-lists (we say sub-list only to help distinguish which list we're talking about; sub-lists are completely capable lists in every respect). The first sub-list of `L` has exactly one item. The print statement displays the third sub-list, i.e., `[4,5,6]`, followed by the value 8. Indexing lists-of-lists requires two pairs of `[]`. In the example, `L[3]` refers to the whole sub-list `[7,8,9,10]`; a second set of `[]` allows us to index into the sublist, so `L[3][1]` is the way to access the second item in the fourth sublist of `L`.

### 2.2.1 Perspective on lists

Python lists are undoubtedly very useful, and you should definitely master the skills to use them effectively. However, novices tend to over-use lists, primarily because of the time invested in

mastering them. Because we will spend a fair amount of time investigating effective use of tools that are not Python lists, you should be prepared to master those tools without relying exclusively on lists. In this course, we will emphasize the use of lists when data is naturally *sequential*, or when data is effectively managed using *integer indices*. These two criteria are quite common, but not universal.

If you are familiar with the concept of an array (because you studied C, C++, or Java previously), you should be warned that lists are not arrays. Arrays can be thought of as very primitive lists. Or, in other words, a list is a highly advanced array. An array is like a bookshelf with a fixed size, and it cannot grow or shrink. A list is exactly as big as it needs to be to store the data it contains. There is no real-world analogy for a list; the closest would be a school bus that changes the number of seats every time a new pupil boards or exits. Clearly such things are technically difficult in the real world, but really quite plausible in a computer. One of the goals of this course is to reveal how these computational devices are designed.

## 2.3 Dictionaries

A Python dictionary is a collection of *key-value pairs*. A *key* must be an immutable value (e.g., an integer, string, or tuple), and a key must be unique, i.e. any particular key can appear at most once in the dictionary. A dictionary key is roughly analogous to a list index. To access a data value in a list, you need to know its *index*, i.e., its position in the list. To access a data value in a dictionary, you need to know its *key*. In a given dictionary, the keys have to be comparable to each other; usually this means the keys all have to be the same data type, e.g., all strings, or all integers.

A dictionary can be created in two ways: using a literal, or using the Python function `dict()`.

```
database = {} # an empty dictionary literal
db2 = dict() # creates an empty dictionary

# a dictionary literal with 2 key-value pairs
assignment1 = {'Grade' : 33, 'Out of' : 35}

# calling dict() with a list of tuples
assignment2 = dict([('Grade', 28), ('Out of', 40)])
```

Note that dictionary literals use the curly brackets {}, and the colon character `:` to separate literal key-value pairs. The `dict()` function is very versatile; it can create an empty dictionary, or if given a list of tuples, it will create a dictionary using the first value in each tuple as the key, and the second value of each tuple as the corresponding value.

We use dictionaries to store the data values; the key is simply the information we use to access the data values. For example:

```
print('You got', assignment1['Grade'], 'on A1')

assignment2['Grade'] += 3 # make a correction to the grade
```

The syntax for accessing the dictionary uses the square brackets [] just like lists do, using the key for a key-value pair, rather than an index.

You can add a new key-value pair to a dictionary using a simple assignment statement, as follows:

```
assignment2['Bonus'] = 3 # a new key-value pair is created
```

Take care to remember that you can't add a new value to a list this way!

You can check if a dictionary has a given key using the `in` operator:

```
if 'Bonus' in assignment1:
    print('You got', assignment1['Grade']+assignment1['Bonus'], 'on A1')
```

And you can iterate over the keys in a dictionary using a for-loop:

```
for k in assignment2:
    assignment2[k] += 3 # make a correction to each value
```

Lists store data values in an explicit sequence, but dictionaries do not. The sequential aspect of a list is crucial. A dictionary's crucial property is association, and there is no standard order for a dictionary's data. If you print a dictionary twice, you might see the data in different order. This is normal and expected, though at this point in your study, unexplained. Consider it Python's business to guarantee that you can always access the data quickly by its key, but it is not Python's business to ensure any particular order of the data when it is printed. This idea will become clearer in later chapters (esp. Chapter 23).

Dictionaries generally represent associations: key-value pairs. You could use dictionaries to represent data mappings, i.e., translating one data value to another. For example, you could represent a mapping between months and integers:

```
month_to_number = {'January':1, 'February':2, 'March':3, # and the rest
```

You can also use dictionaries to store records, a collection of organized data values associated with a single topic or purpose. For example, we can store information about an assignment as a record, as we saw above:

```
assignment1 = {'Grade' : 33, 'Out of' : 35, 'Completed':'8 January'}
```

If you had a list of records, you could call it a database.

Lists can contain dictionaries as values, and dictionaries can contain lists as values (but not as keys, because lists are mutable; a key must be immutable).

### 2.3.1 An example program

Suppose you were doing some analysis of data from some experiment you had designed. Your data is stored in a text file, and the file consists of 3719 rows (lines), each with 11 comma-separated columns of numeric data. You might write a script to read this data as follows:

```
1 # open the file
2 fh = open('datafile.txt')
3
4 # read the file
5 data = []
6 for line in fh:
7     # strip the junk off the end of the line,
8     # and split the line into pieces
9     row = [float(d) for d in line.rstrip().split(',')]
10    data.append(row)
11
```

```
12 # close the file
13 fh.close()
```

You might need to calculate some statistics on the data stored in any given column. One way to accomplish this task is to write a function that takes a list of data, and returns a record containing the maximum, minimum and average of the data in the list, as below:

```
1 def calc_stats(col):
2     '''Calculate min, ave, max of a column of data'''
3     stats = {}
4     stats['minimum'] = min(col)
5     stats['maximum'] = max(col)
6     stats['average'] = sum(col)/len(col)
7     return stats
```

Notice how the function creates and populates a dictionary, which serves as a record of the statistics for the column. This function makes good use of a dictionary, for several reasons. If you wanted to calculate additional statistics (e.g., median, standard deviation), you could simply add a new key-value pair to the dictionary. The use of a dictionary to return a bunch of values is better than using a list, because a string like `'average'` on line 6, is more intuitive than an index, like 2.

## 2.4 Modules

The basic Python language itself provides a lot of tools that are useful for programmers (lists, dictionaries, strings, etc). Additionally, many useful tools are provided by *modules*. Modules allow programmers to share functions (and other computational abstractions, which we will discuss later) that they have written so that other people can use them in their own programs.

Modules are stored in separate files, outside a programmer's scripts. Thus, in order to use the functions and objects defined by a module, we have to tell a script to look in those files and read those definitions. We do this using a Python keyword called `import`.

In general, the syntax for importing modules is: `import x as y`. Here, *x* must be the name of a module, and *y* must be a valid variable name (often an abbreviation for *x*). This technique creates an *object* called *y* that contains the functions defined in *x*. We'll say more about objects later in the course. There are other ways to use the `import` command, but we will use this one exclusively, just so that there is no unnecessary confusion.

As an example, we will use the random module:

```
import random as rand

rval = rand.randint(0,10)
rlist = [rand.random() for i in range(10)]
```

Here, `rval` is chosen randomly from the integers from 0 to 10; note that the range of `randint()` is inclusive on both sides, unlike the Python function `range()`, which is exclusive on one side! At the end of the script, the variable `rlist` refers to a list that contains 10 random floating point values between 0 and 1, which is also exclusive on one side. Notice that we used the functions in the random module using object dot-notation: we gave the name of the module (we called it `rand` using the `import` command), and then the function name, with a dot between them.

## 2.5 Three Example Problems

In this section, we'll describe three example problems that will set the stage for several later topics. The problems themselves are interesting and have many applications and extensions.

### 2.5.1 The Gambler's Ruin Problem

In science, and in other human endeavours, we occasionally run into phenomena that we cannot explain with a simple formula. When the behaviour of the phenomenon is too complicated, or unknown, we often turn to models we call "stochastic," which basically means that some amount of randomness is used in the model. The randomness is used as a replacement for a mechanism that is either too complex or too obscure to model directly.

One very simple example of a stochastic model is called "The Gambler's Ruin." In this scenario, we imagine a gambler playing some kind of betting game. The game itself is not important, as we are not presently interested in modelling the game or the gambler's strategies. The game gives the gambler a 50-50 chance of winning each time the gambler plays. If the gambler wins a game, he earns 1 unit of currency, and if he loses, he pays 1 unit of currency. The gambler starts with a given quantity of currency; we'll call this the *stake*; for sake of clarity, let's say the stake is 20 units. The gambler also has a *goal*: to earn a known quantity of the currency; let's imagine the goal is to reach 100 units. If the gambler reaches his goal, he stops gambling, as a success. If the gambler reaches zero units of currency, he also stops gambling (the house does not allow gambling on credit), as a failure.

One question we want to answer is this: what is the probability that the gambler will reach the goal? The second question is: success or failure, how many bets, on average, does the gambler place before stopping?

We will write a program to play the game a number of times, and use the data generated to estimate the answers to our questions. It's a form of simulation called "Monte Carlo Simulation" which is very common in applications of statistics to science. Without a computer, this kind of simulation would be very difficult to carry out.

### 2.5.2 The Coupon Collector Problem

Another scenario in which Monte Carlo simulation is very useful is called the "Coupon Collector". In this scenario, we imagine that someone is trying to collect coupons from some kind of advertising or customer loyalty program. There are 100 different kinds of coupons, and the collector gets any one of these coupons at random (presumably every time she visits the store, or the website, but it really doesn't matter). The chance of obtaining any of the coupons is equal, i.e., 1 in 100, and the collector will almost certainly collect more than one copy of some of the coupons.

The question we want to answer is: how many coupons will the collector have in total by the time she's gathered one of every type? It's at least 100, but because repeats are allowed, it is probably higher than that. Allowing repeats is sometimes called "sampling with replacement." Because it's a stochastic process, we'll want to run the simulation a number of times, and take an average as the estimate.

This problem may seem somewhat trivial, but this model can be applied in science to determine if a phenomenon is random or not.


### 2.5.3 Self-Avoiding Random Walks

The Gambler's Ruin problem can be interpreted as a one dimensional random walk, as follows. Imagine starting at a point somewhere between two ends of a line segment, and randomly choosing to walk one step forward or back. In this interpretation, success is when you reach the one end, and failure is when you reach the other. In this interpretation, the one dimensional random walker is allowed to go back and forth, and is allowed to re-occupy positions multiple times. The probability of success depends on where the walker starts.

The self-avoiding random walk problem is similar, but has a few important differences. The walker starts in the middle of a large room, and can take a step at random in any of the 4 cardinal directions (forward, back, left, right), subject to a constraint: the walker is not allowed to step anywhere he has already stepped. In other words, we have to keep track of where the walker has been. If a random step is chosen that would lead the walker into a position previously occupied, that step is discarded, and a new random step is chosen. If the walker reaches the "wall" of the room, it is called a success, and if the walker gets caught in a position where there are no legal steps that can be made, it is a failure.

The question we want to answer is this: what is the probability of failure?





Introduction  
Values, objects, and addresses  
Expressions  
Variables and frames  
Equality  
More about frames and functions  
Frames and environments  
Lists and Dictionaries  
Python caches values for efficiency  
    Mutable vs Immutable Data  
Casual language  
Summary

## 3 — References

### Learning Objectives

After studying this chapter, a student should be able to:

- Explain the difference between a value and an object.
- Explain what a reference is.
- Explain how Python uses frames and references to associate variables and values.
- Explain how Python evaluates expressions involving mutable and immutable values.
- Draw diagrams showing the frame(s), values, and objects, given a sequence of Python expressions or statements.
- Explain how frames are used for a function's local variables.
- Explain what happens when a local variable shadows a global variable.
- Explain the difference between `==` and `is`.
- Explain the difference between copying references and values.

### 3.1 Introduction

To prepare for later chapters, we need to revisit and clarify a concept that we briefly mentioned in the introductory course. To do this, we will review some familiar concepts in a very precise way, and show how Python depends on a concept called a *reference*. References are ubiquitous in Python; every data value calculated by a script has a reference to it, and every variable records a reference to a data value. We need references to build effective tools to organize data, which we will start in Chapter 13 and continue in Chapter 20. References are used in every programming language to some degree or other. For example, Java uses references for its objects. In C and C++, special variables called *pointers* allow programmers to store and manipulate references. These are all the same basic concept. Mastering the concepts of this chapter will deepen your understanding of Python for your own projects, and will be the foundation on which you will build skills in later courses.

## 3.2 Values, objects, and addresses

A data *value* is any piece of information that you want to use in a script or a program. Data values can be integers, Booleans, floating point numbers, or strings. We will also include compound data values such as lists, dictionaries, tuples, and any other kind of collection of information one could construct (as we will see in Chapter 9). Data values are meaningful to human beings, and we don't need computers to work with them. We can do arithmetic “in our heads,” or on paper; we can keep track of appointments and addresses in a little book.

In Python, all data values are stored inside the computer as *objects*. An object encodes, or *represents*, a data value inside a computer. Because objects represent data values, we can say that every object *has* a data value. If we know Python's internal rules, we should be able to look at an object and figure out what data value it represents. We won't need to learn those rules in this course, but they will come up in later courses.

To distinguish between data values and objects, we could say that a data value is information that is independent of the way in which a computer encodes it, but an object is completely dependent on the rules used to encode a value. In Python, values like numbers are represented by simple kinds of objects; lists and dictionaries require more complicated objects. Other programming languages may have slightly different rules for encoding data values.

Because objects are stored in computer memory, every object also has an *address*. An address is nothing more than an integer number, representing the location of an object in computer memory; in many ways, an address is like an index (also called a subscript, or an offset) for a very long list. For two objects to be different, it is only required that they have different addresses. But it is possible for two different objects to have the same data value. For example, if you make an exact copy of a list, you have two lists that represent exactly the same information, but they will have different addresses. If later on you change one of the copies, but not the other, the objects will have different values, as well as different addresses.

In the follow sections, we will use the term *value* as a short-hand for the more complicated phrase *object representing a value*. This is just a way to simplify the writing, even if it is not quite precise.

## 3.3 Expressions

In the previous course, we relied on your previous education and your intuitions when discussing expressions. We were pretty certain that you'd understand that `3.14` is a number. We hoped that you'd be willing to see the literal `True` also as a data value (as opposed to an abstract concept). We pretty much had to assert without justification that `'Hello, world'` should be considered a data value, and that `3.14` and `'3.14'` are different kinds of data values.

Indeed, we relied on your understanding that `2 + 1` should result in the data value 3. We relied on you to remember that relational expressions like `5 > 7` could be `True` or `False`, but we may have had to convince you that this kind of expression actually results in Boolean data values (in this case, the value `False`). We relied on your linguistic intuitions when we described the meaning of Boolean expressions like `3 < 0 or 7 > 5`. For expressions involving lists, e.g., `[1] + [2]` and strings, e.g., `'Hello!' * 3`, we relied on your ability to see logical patterns and apply them. At this point, we're ready to go into a bit more detail.

People write *expressions* when they want to calculate data values. More specifically, a Python programmer would write an expression in a Python script to cause the computer to calculate a data

value. We have *literal* expressions, like 0, 3.14, 'Hello, world.', [1, 2, 3], etc.; these are literals because they look literally like their data value. More generally, an *expression* can involve literals, operators (like + or \*), function calls, and variables. We have arithmetic expressions like 2 + 1, relational expressions like 5 > 7, Boolean expressions like 3 < 0 or 7 > 5, as well as expressions involving lists, e.g., [1] + [2] and strings, e.g., 'Hello! ' \* 3.

Python expressions are not data values, nor are they objects; they are nothing more than strings in a Python program. In order to become data values, Python has to *evaluate* those expressions. To *evaluate* any kind of expression, Python examines the expression to determine what computations are required, then performs those computations. These computations always result in the creation of a new object, which is automatically stored in computer memory when the evaluation is complete. Arithmetic expressions create number objects, string expressions create string objects, Boolean expressions result in Boolean objects, etc. It is appropriate to consider any expression involving an operator, like 2 + 1 or 'Hello! ' \* 3, as being a function call: the operator is the function's name, and the function creates the new object from its operands according to some code that we cannot see directly. They are very much like functions we would write ourselves in Python (especially after Chapter 24).

Once Python creates an object, Python never moves it. Python programs do not toss objects around like beach balls, even though we casually talk about passing arguments to functions, returning values from functions, and sending data to the console. Instead, Python conveys information about objects indirectly, using addresses. For this reason, we frequently use the term *reference* as a synonym for *address*. Python does everything with references. In Python, functions and operators do not receive data values directly; instead, arguments and operands are referred to indirectly by addresses (i.e., references to objects). If a Python function returns a value, the value itself is not returned, but a reference to it is returned instead.

For example, if we see the expression 1 + 2, we naturally think about the values 1 and 2, and combining in some way to produce the value 3. Python creates objects to represent the 1 and the 2, translates the expression involving values to a function call involving references. The Python addition operation takes the two references, creates a new object in computer memory, representing the value 3, and then returns a reference to the new object.

Python does all this without our help, and for a lot of the code we write, we didn't need to understand this aspect of Python to write useful Python scripts.

## 3.4 Variables and frames

In most of the introductory course, we needed only to concern ourselves with two main ideas about variables. We needed to understand that a variable has a name, and that it has a data value associated with it. The first thing we said about variables was that “[v]ariables are a way of giving names to data.” Later on we said that variable assignment “is used in Python to *assign* a variable name to a value.”

For our purposes, we need to realize that variables are not objects or data values. In Python, a variable is nothing but a name in a Python script. Using an assignment statement, we can give an object a name, which is a vast improvement over the alternative: describing an algorithm in terms of the location of the data values. A variable only knows one thing about the object it is assigned to: its address. In Python, a variable does not know or care what kind of data value the object represents, or what process created it, or anything else.

The association between variable and value is in fact accomplished by means of a table that

Python creates and manages. This table stores variable names used in the script, and references to objects created by the script. We call this table a *frame* (or name-space).

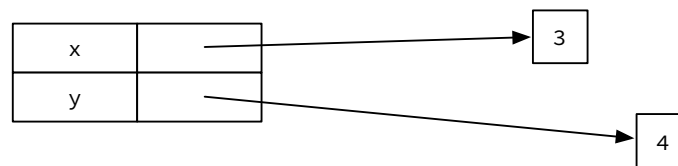
For example, consider the following simple script:

```
x = 3
y = x + 1
```

A frame showing the end result of the above script would look like this:

Name	Address
x	0x10397e7b8
y	0x10397e840

The numbers stored in the second column are machine addresses encoded in hexadecimal notation. If we looked in the computer memory at address 0x10397e7b8, we would find an object representing the integer 3. Because the actual address is a detail, we use arrows to represent a reference, as an abstraction. A simple picture using arrows would look like this:



The variable names *x* and *y* are in a frame on the left, and arrows refer to objects representing the values 3 and 4. In figures and diagrams, any box that stores a reference always shows an arrow leaving from inside the box to some object in computer memory. This reminds us that the arrow is a stored address, i.e., a reference. In this example, the arrows start in the second column of a frame, but there are other ways references are used in Python. We will talk about them later in this chapter.

Simple assignment statements, i.e., when we assign a variable to a value, change the references stored in a frame. When a variable is initialized, i.e., its first assignment statement, Python adds the variable to the frame, and also stores the address of its initial value. Assigning new values to the variable after the initialization causes the variable to refer to a different value. When Python executes a simple assignment statement, it first performs the calculations described by the expression on the right hand side, then stores the result, namely an address of a value, along side the variable in the frame.

When a variable is used in an expression, Python looks in the frame for that variable; if the variable is in the frame, Python uses the reference stored there to evaluate the expression. If an expression uses a variable that has not been initialized yet, Python will raise a `NameError`, which is basically Python's way of saying *that name is not in the frame yet*.

Let's walk through the script, laying out the full detail of how variables, frames, and objects are used:

```
x = 3
y = x + 1
```

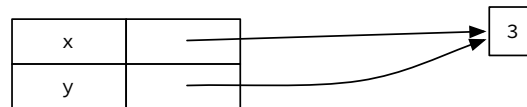
On the first line, Python adds the variable *x* to the frame, creates an object representing the value 3 in computer memory, and stores the address of 3 in the frame beside *x*. On the second line, a new

variable `y` is initialized, adding a new entry in the frame. The right hand side of the assignment is an expression involving the variable `x`. To evaluate the expression `x + 1`, Python uses the reference stored with the variable `x`, and does addition to produce a new object representing the value 4. The result of the expression is the address of the object created in the expression, and that reference is stored in the frame with the variable `y`.

It is important to understand that assignment statements do not copy values. Consider:

```
x = 3
y = x
```

One object is created, the object representing the value 3. On the first line, the variable `x` is added to the frame, and the address of the value 3 is stored with it. On the second line, a new variable `y` is added to the frame, but the reference stored with `y` is the address of the value 3. In other words, the right hand of the assignment does not make a copy; it's simply an address of something that `x` already refers to. So we have two variables referring to a single value.



This behaviour is not specific to numbers. It does not matter what `x` refers to, be it a list or a dictionary or anything else. In Python, writing `y = x` never makes a copy of `x`. It only makes `y` refer to the object that `x` refers to already.

### 3.5 Equality

Now we can explain the difference between `==` and `is`. The Python operator `==` compares the value of two objects, ignoring their addresses. The Python operator `is` compares two addresses.

These two operators are not interchangeable! The answer to `is` is stronger than the answer to `==`. If `x is y` happens to be true, then it must be true that `x == y`. This would be a case where two variables refer to the same object. However, if two objects have the same value, that is, if `x == y`, it is not necessary that they are at the same address. One could be a copy of the other.

### 3.6 More about frames and functions

In the preceding sections, we introduced the concept of a frame, but some details were left vague. In this section we'll make some of the details a bit clearer.

Python creates several frames to manage all the variables (and their references) that come up in a Python script. First, there is a *global frame* that contains all the variables you create when you write code outside of any functions. The global frame is created when Python starts running, and it contains all the names of Python's built-in functions, e.g., `len()`, etc. This global frame is where all your global variables are found.

Python also creates a new frame each time a function is called; this new frame contains all the parameters and variables created by the function. It's the function *call* that causes a new frame to be created; a function *definition* does not create any frame. Remember that a function's parameters are also variables, so they are also part of the frame that gets created when a function is called. When

the function is called, the function's parameters get their initial values by an implicit assignment using the arguments in the function call.

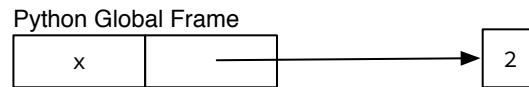
For example, consider the following:

```

1 def fun(a, b):
2     c = (a + b) * 2
3     return c
4
5 x = 2
6 y = fun(x, 4)

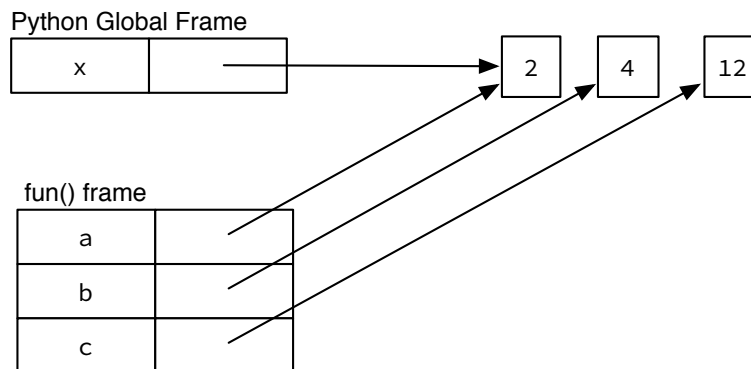
```

Line 5 causes the variable `x` to be added to the the global frame, with an address of the value 2. We could picture it as follows:



Line 6 is a function call; the arguments in the function call are given as initial values for the parameters of the function. So Python creates a frame with the variable `a` referring to the value 2 and variable `b` referring to the value 4. Note that `x` refers to the value 2, and it's the address of the value 2 that gets copied to the new frame. For the second argument, the literal 4, an address of the value 4 is entered into the new frame.

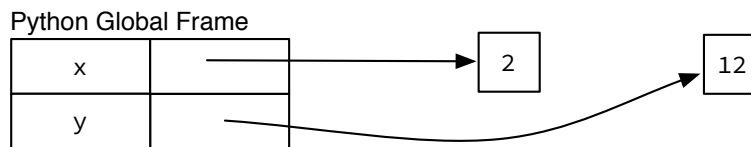
When line 2 (the first line inside the function) is executed, the variable `c` is added to the frame, with the address of the value 12 (which is the result of the expression  $(a + b) * 2$ ). The situation now looks like this:



The frame for the function `fun()` has three entries. It's very important to realize that the value referred to by `a` is not a copy of the value referred to by `x`. Arguments' values are never copied; when we use an argument, its reference is entered into the frame, so the variable `a` in the new frame stores a copy of the address of the value 2. In the global frame, `x` refers to it, and so does the variable `a` in the frame for `fun()`. Again, this is true no matter what kind of thing `x` refers to.

The return statement on line 3 returns an address of the value 12, and that reference is stored in the global frame beside the variable `y`. When the function is finished, i.e., after the return statement is executed, the frame created by the function call on line 5 is no longer needed, and can be destroyed. Only the global frame remains, with the new entry for `y`:





Notice that the frame for the function is no longer in the picture. The value 4 has disappeared too, because nothing in the global frame referred to it. But the value 12 is still there because *y* refers to it.

### 3.7 Frames and environments

Let's consider a simple Python script without any user-defined functions. We call such code *global*, because it is outside any function definition. This kind of script has access to all the Python built-in functions, and can freely create variables using assignment statements. All of these names are part of the global frame. In a very real sense, when the global code is running, we can say that the global frame is acting as its *environment*. An environment is essentially a collection of frames that indicate what names are accessible by a line of code. In global code, the environment is the global frame. If the global code uses a name that does not appear in the global frame, Python issues a `NameError`.

We have already said that a new frame is created when a function is called, and that the variables created in the function's code are entered into the function's frame. While the function is running, the function's new frame is part of the function's environment.

When code inside a function is executed, Python checks for names by first looking for the name in the function's frame; if the name is there, Python uses the reference stored there. However, if code in a function uses a name that does not appear in the function's frame, Python looks outside the function's frame to the global frame. If the name is not in the global frame either, we get that `NameError`. This establishes a precedence inside the function's environment: first the local frame, then the global frame.

Now we can explain what happens when a function creates a variable with the same name as a global variable. For example, consider the following slight modification of an earlier example:

```

1 x = 5
2
3 def fun(a, b):
4     x = a + b
5     return x * 2
6
7 y = fun(x, x + 1)
```

A quick scan of this code reveals that the script uses variables *x* and *y* in the global frame. The function has two parameters, *a* and *b*.

The assignment statement on line 4 uses a variable named *x*. Python's rule for variables causes this assignment statement to create a local variable called *x*. This rule says that the first time a variable is assigned a value in a function, a local variable is created. In this case, it happens to have the same name as the global variable, but it appears in a different frame. Because of this rule, and the precedence of local frames, the expression `x * 2` on line 5 uses the local variable, not the global variable. When we have a local variable and a global variable with the same name, we say that the local variable *shadows* the global one.

The rule about creating local variables can be defeated by using the keyword `global`. However, modifying global variables inside a function can create errors in your code that are very hard to find and fix.

Global code is permitted to create, make use of, and re-assign global variables. Code internal to a function is allowed to create, make use of, and re-assign local variables. Code inside a function can access the values of global variables, but by default, is not allowed to reassign global variables. Unless the keyword `global` is used, any attempt to reassign a global variable creates a local variable instead!

A full description of Python's treatment of frames is somewhat more complex than what we want to get into in this chapter. For the kinds of functions we have been writing, we can treat frames as being created when the function is called, and being destroyed when the function returns. We can treat functions' frames as being mutually exclusive, meaning that no code can look into the frame of another function, and use the variables there. We can treat name lookup as starting in the function's frame, and if necessary, looking through the global frame. These simplifications are not exactly true, but good enough for now. We will fill in the missing detail in later chapters.

### 3.8 Lists and Dictionaries

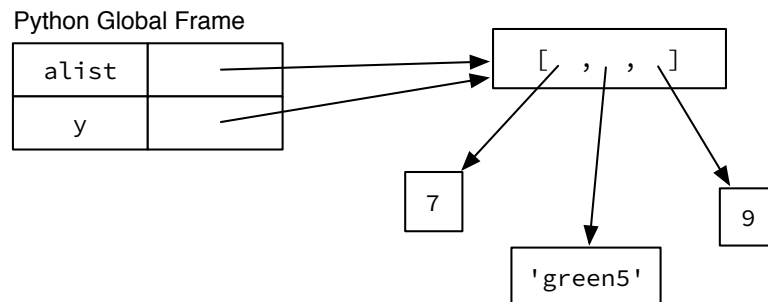
We have emphasized that variables do not store values; they have references to values. We have described that assignment statements change what a variable refers to: a simple assignment like `x = 3` changes the reference for `x` in a frame. But we can use lists and dictionaries in assignment statements as well. In this part of the chapter, we have to explain how lists and dictionaries contain other values.

Consider the following example, assuming it is global code:

```
alist = [7, 1, 9]
y = alist
alist[1] = 'green' + str(5)
```

The variable `alist` is stored in the frame, along with an address of a list (a list is a compound data value, but still a data value). Intuitively, we can understand that when the script is finished, the string `'green5'` will appear in the list at index 1.

Python lists do not store data values; they store references to values. For example, when the above script is finished, the data value `'green5'` was created and stored in computer memory, and the address of it is stored in the list at index 1. Similarly, the data value for the literal 7 is stored in computer memory, but the address of it is stored in the list at index 0. A picture of the situation after the script has completed follows:



Here we have arrows starting from inside an object (not a frame), pointing to other objects in computer memory. Note how the diagram starts the arrow inside the list. It reminds us that the arrow is an address, stored in the list.

Dictionaries behave in a similar way. They do not store objects at all. Both the key and the value for each dictionary entry are references to objects in computer memory.

This is the first hint of the value of references to our work in later chapters. Make sure these concepts are well understood, because we will be using them!

### 3.9 Python caches values for efficiency

We said that all data values are stored in computer memory, and that expressions cause data values to be created. This is true most of the time, but there are a few exceptions. Some values are used so often that Python creates objects for them in advance, and never creates a new object with the same value.

For example, the literal value `None` is an object, but it is created only once, and this happens when Python starts up. Every expression whose value is `None` refers to this singleton value; Python makes sure of it! The Boolean values `True` and `False` are also created when Python starts up. Every expression whose value is `True` refers to the one and only `True` value created; similarly for the object `False`.

This technique is an example of a more general technique called *caching*, i.e., computing something once and keeping it available, because you know it will be needed many times in the future. Caching commonly used values makes Python run faster, without changing how normal programs work. This caching technique will not change correctly written programs, and you will not notice the caching at all, unless you are writing programs to check if Python is caching these values.

Since small integers are very common in Python programs, at startup, Python creates values for a somewhat small range of integers. Any expression whose value is one of these integers refers to the one and only object stored in memory for that value. The exact range might be different on different computers or different versions. Note that small floating point numbers are not treated the same as integers; every floating point expression results in a new floating point object.

We should always use `==` for numerical comparisons. We do not care about addresses (references) when we are dealing with numeric values. For values that are not cached, `x is y` will be false when `x == y` is true. On the other hand, comparison to `None`, `True`, and `False` can be done with either `==` or `is`; but comparing addresses is more efficient, so prefer `is`.

#### 3.9.1 Mutable vs Immutable Data

Numbers, strings, tuples, and Boolean values are *immutable*: they cannot be changed. When you add two numbers together, say `3 + 4`, you're not changing the numbers, you're creating a new value, `7`.

But now consider a *mutable* data type, like a list or a dictionary. When we create a list, we are creating a new value; it's a compound data value, but still a value, and it's in computer memory. Mutable data values can be changed. That means if you extend a list, add or remove elements, the list actually changes; however, the address of the list does not change. Consider the following example:

```
copy = []
for v in a_list:
    copy.append(v)
```

The first line has a literal expression, an empty list. Python creates a new empty list, and stores the address of the new list in the frame beside the variable `copy`. Observe very carefully that inside the loop, there is no assignment statement for `copy`. Because there is no assignment statement, the script never changes what `copy` refers to; it refers to the same list object the whole time. But that list grows in size to accommodate all the items in `a_list`. Notice that the method `append()` is used for its effect, namely, adding a value to the end of the list. This method returns `None`, which the script ignores.

The approach Python takes, wherein every variable stores an address, is very versatile. It allows us to do complicated things in Python with very little effort. When we talk about lists of lists, we never actually store a list inside another list. Lists are objects; the list does not store values, but it does store addresses. A list-of-lists is really a list of addresses, and those addresses refer to other lists. We can store addresses in any kind of container: lists, dictionaries, and tuples. We can have a list of dictionaries, which is really a list of dictionary addresses. We can have a dictionary that refers to other dictionaries, or lists, or anything else for that matter. A variable can refer to any kind of data value: integers, lists, dictionaries, whatever. It's the data itself that knows its type; not the variable.

### 3.10 Casual language

As you may have noticed, this chapter introduced the difference between value and object fairly carefully, but then the distinction between these two concepts was left up to the reader to make. We referred to objects as values several times. When we need to be precise, we have to be precise, but when precise language gets in the way, more casual language is preferred.

There are computer scientists who say that abstraction is the fundamental skill, and this fluidity between precision and imprecision is part of what they are talking about. For a beginner, it can get confusing; be as precise as you can be until you master the concepts.

There are other examples of the use of casual language. Most of the time, we'll be happy to say that a variable "stores" a given value, even though it is more precise to say that a variable *refers* to a given value. When it is convenient and helpful, we will gloss over the notion of reference entirely; that's a detail that doesn't always need to be made precise. There are times when the notion of reference is crucial, and it is at those times when we have to be precise.

### 3.11 Summary

In Python, all data values are stored inside the computer as *objects*. An object encodes, or *represents*, a data value inside a computer. Because objects are stored in computer memory, every object also has an *address*. To *evaluate* any kind of expression, Python examines the expression to determine what computations are required, then performs those computations. Any expression involving an operator can be understood as a function call: the operator is the function's name, and the function creates the new object from its operands according to some code that we cannot see directly. In Python, functions and operators do not receive data values directly; instead, arguments and operands are referred to indirectly by addresses (i.e., references to objects). If a Python function returns a value, the value itself is not returned, but a reference to it is returned instead.

In Python, a variable is nothing but a name in a Python script. Using an assignment statement, we can associate an object with a variable; these associations are managed by Python using a table called a frame.

## 4 — Software Design Goals

### Learning Objectives

After studying this chapter, a student should be able to:

- Define design goals of correctness, and efficiency.
- Define implementation goals of robustness, adaptability, and reusability.
- Assess, at a preliminary level, the quality of example code with respect to the design and implementation goals.

### 4.1 Design Goals

We write software to accomplish a task. Sometimes the task is “get credit in a course,” and sometimes it’s “do something fun.” But often, we will write software because we want the computer to do some work we need to get done. Writing software can take a long time and a lot of effort. We invest the time to write software for fairly specific reasons:

- The computer can do the work much faster than if we do it ourselves by hand.
- The work will be done more than once.

Either of the above is good enough reason; they don’t both have to be true.

It would be a pointless exercise to write a program if:

- The program didn’t give the right answer.
- The program gave an answer too late to be useful.
- Writing the program was more expensive than getting the answer some other way.

Given the reasons we write software, it’s important to be concerned with the *correctness* and *efficiency* of our software.

**Definition 4.1** Software is *correct* if it does everything it is intended to do, and nothing that it isn't supposed to do.

A program that is not correct is said to have errors (though sometimes we call those errors “bugs” to make them seem less serious). Very often, a program is partially correct: it gets the right answer some of the time, but maybe not always. For example, a program that simply displays the text “Today is Monday!” is correct on Mondays, but incorrect all the remaining days of the week. As another example, consider a function that returns the square root of any number: if it returns a useful approximation for every non-negative number, it's doing what it should, but if it goes into an infinite loop when the input is negative, it is doing something it shouldn't be doing.

**Definition 4.2** Software *efficiency* is a relative measure of the degree to which software consumes computational resources, including time and memory.

A program that is not efficient may be using an algorithm that takes more time than a better algorithm, or it may be using more memory than it needs to. On modern computers, any given program may be efficient enough if the task is small, or if the program won't be reused. But once the tasks start getting big, inefficient programs may give an answer too late to be useful. On modern computers, programs usually are running alongside other programs (multi-tasking), e.g., your web browser and your video editor may be running at the same time. Programs that do not make efficient use of resources will interfere with other programs by using up resources that ought to be shared.

Correctness and efficiency are our design goals. There are quantitative measurements that we can use to assess them. We will discuss these topics in more detail in Chapters 7 and 18.

## 4.2 Implementation Goals

If software were art, then we'd finish writing the software and it would never change after we're done. But it's not art: software changes every time someone thinks of a way to make it better. People who work with software rarely get to admire it; they are busy trying to get it to do more! So it's also important for software to have other attributes, related to the goals of software developers.

**Definition 4.3** Software is *robust* if it behaves well even when things go wrong.

Even if a program is correct, things can go wrong if the user enters bad data, or if a network connection is lost, or any other kind of external interaction behaves unexpectedly. To assess *robustness*, we consider what kinds of things might go wrong that are beyond the control of the program. A program that is not robust is called *brittle*. Brittle software crashes without warning, causing data loss or other inconveniences. If a program still has errors, a robust program might be able to recognize that an error had occurred, and recover from the error without causing the program to crash. Even if a crash is unavoidable, it would be useful for a program to do something helpful, like saving a document, or reporting that an error had occurred with a helpful error message.

**Definition 4.4** Software is *adaptable* if small changes in behaviour of the software can be achieved by making small changes in the software.

For example, suppose you've written a program to calculate the average of 10 values in a list or array.



An adaptable program would not require a lot of changes to calculate the average of the numbers in a list or array of any size. If you had to do a major rewrite of the code, the code you started with was not very adaptable.

The reality of writing software is that software is never truly finished. There's just a point at which someone stops working on it. Somebody, sometime (possibly you, tomorrow, when you realize you misunderstood something yesterday — oops!), will find errors, or other short-comings. Your client may request changes to the requirements you thought had already been finalized, and to meet those changes, you'll have to adapt the code you've already written. Somebody, possibly you, will realize that the program you wrote 4 months ago to solve problem X can be modified to solve related problem Y. If you design your software with these kinds of possibilities in mind, making the changes will be easier. If you give no consideration to adaptability when you write programs, you are probably making someone else's work in the future more difficult. You might not care about that at all, except that very often, that person will be you.

**Definition 4.5** Software is *reusable* if it can be used more than once.

We already make use of a lot of reusable software. Every programming language provides tools to help the programmer that can be used in any project. How much harder would it be for programmers if Python's `print()` function only accepted a single argument? Think of the Python function `range()`: it was designed to be called with one, two, or three arguments, giving the programmer a variety of control and convenience. Think of all the math functions; think of all the modules! This kind of reusability does not happen by accident. Software designers deliberately design these components to be reusable.

Robustness, adaptability, and reusability are our implementation goals. They are qualitative: there's currently no way to measure them scientifically. For example, software can be more, or less, robust, but we don't have a precise way to measure it.

### 4.3 Example

Consider the short Python program below (which we studied in Chapter 2):

```
1  unsorted = [3, 2, 5, 7, 6, 8, 0, 1, 2, 8, 2]
2  sorted = list()
3
4  while len(unsorted) > 0:
5      out = min(unsorted)
6      unsorted.remove(out)
7      sorted.append(out)
8
9  print(sorted)
```

It creates a list of numeric values whose values are in increasing order. The algorithm can be summed up as follows: *Keep on removing the smallest value from unsorted and adding it to end of sorted until there are no more values left in unsorted.* This is one way to interpret the notion of *sorting* a list. Let's consider the algorithm in light of the design and implementation goals from the previous sections.

In terms of assessing or demonstrating *correctness*, we have to give some reason to believe that the code does what it is supposed to do. The English language summary (repeated above) lends credibility to the correctness of the script, though we must be careful, since a summary omits detail, and it might be the details themselves that are wrong. We could try to highlight the principle upon which the algorithm seems to be based. This requires some insight. For example, when you remove a smallest element  $x$  from a list, the next time you remove a smallest element, say  $y$ , it will be larger than  $x$ , and so it should come after  $x$  in the sorted list. Mathematical tools (such as mathematical induction) can be applied to address the question of correctness. However, at this point in your study of computer science, the main tool you currently have to establish correctness is testing. We will have much to say about testing in Chapter 7.

To assess *efficiency*, we have to be able to measure the quantity of resources used by the program, and compare it to other programs or algorithms. It might be tempting to make a vague statement: taking one element out of a list and putting it into another list seems more or less efficient. You may recall that finding the smallest element of an unordered list is a form of linear search, and to find the smallest, linear search may have to look at every element. The fact is that this algorithm is not terribly inefficient, but it is not as good as some of the algorithms we studied in the previous course. In a later chapter, we will have much to say about how to analyze algorithms and establish their costs in a scientific way, so if you were hesitant to assert strongly that the algorithm is efficient, that is probably a natural feeling to have right now.

To assess *robustness*, we consider what kinds of things might go wrong that are beyond the control of the program. Because there is no user input, no file I/O, and no other interaction with the outside world at all, the behaviour of the program is entirely deterministic, i.e., the behaviour of the code is determined by the code itself. As a result, we can call the program robust, even though we did nothing special to make it so.

In terms of *adaptability*, we have to consider what kinds of changes we might imagine making to a program that sorts lists. For one, we might want to change the order of the sorted list to be decreasing; we could do this by changing `min` to `max`. We might want to sort lists of non-numeric values; fortunately, Python assists us here a lot: the function `min` (and `max`) can be applied to lists containing numeric values, or string values, as long as all the values are the same type. We would not be able to use the algorithm as it is to sort more complicated values, say employee records. From this discussion, we might conclude that the program is somewhat adaptable, but limited. We will see, in later chapters, how to get past limitations that seem to prevent programs like this from being applied to complicated values.

Finally, *reusability*. Could this short script be reused in other programs, and if so, how? The algorithm sorts a list, and does not make any assumptions about the size of the list, so it will not be a problem to apply the algorithm to larger lists of numbers. The algorithm depends on the variables `unsorted` and `sorted`, which is a bit of a problem. Right now, the only way to reuse this program is copy-paste, which is well-known to be a bad strategy for code reuse. The program would be a lot more reusable if the algorithm were encapsulated as a function, and saved as part of a module; that is, of course, why we have built-in Python functions like `sorted()`. It would be fair to say that the example program is not very reusable at all! To restate the theme of this chapter, later on in this course we will see principles to make our programs more reusable.

## 4.4 How to Use the Design and Implementation Goals Effectively

One of the major learning objectives of this course is to begin the development of the intellectual tools to achieve the goals we've outlined above. To develop these skills, you should be keeping these goals in mind as you work. You should see them as concepts that describe ways that you could improve code that you are writing. You can, and should, ask questions such as:

- How do I know this program works correctly?
- Are there situations that might break this code?
- Are there ways to make this program more efficient?
- Have I made coding choices that reduce future adaptability?
- Which parts of this program might be useful in another project?

It is important to work on these skills; try to be mindful of these questions, and try to gain experience in making these judgments. You'll get lots of examples in this course. Perfection is not required at this point in your education. You shouldn't be overly concerned about achieving perfect code that satisfies all of our design goals in your first draft, and you shouldn't use them as excuses to avoid making progress on your work.

## 4.5 Pedagogical Goals

Presumably you're taking this course to learn to write software, so we have to make room for pedagogical purposes that are often at odds with practical purposes. Students write software to gain experience writing software, even if the software they end up with is not completely correct, and even if they take longer than an experienced practitioner would. Experience is a very important pedagogical goal for students! Taking time to gain experience is necessary preparation, so take time to make mistakes, and to learn to fix them. It's not a waste. It's what you're here to learn!

It's also important for professionals to take time to learn new things, so the pedagogical goals don't end immediately after you finish a course or graduate from a program. Another important aspect of this course, and all university courses, is that you learn enough fundamental concepts that you are prepared for a lifetime in which changes happen fast. In all likelihood, you have to learn new things for the rest of your career, at a pace that might have frightened your grandparents.

One of the things you can learn immediately is how to distinguish a principle that will remain useful throughout your career, from a detail that might change as technology changes. For example, the Python language is a detail; if you have a career in software development, you'll have to learn a dozen languages. The new best language ever is just a few years down the road. On the other hand, software efficiency is a principle.



#### Software development processes

Waterfall model

Iterative and incremental model

Requirements and Specifications

#### Software Design Strategies

Version Control

Prototypes

Summary

## 5 — Software Development Processes

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe the difference between software development and software design.
- Describe the waterfall model of software development.
- Describe the iterative method of software development.
- Apply the waterfall and iterative models to assignment-sized problems.
- Discuss the importance of version control.
- Explain the purpose of prototyping.

Producing quality software is expensive in terms of time and money. Software design processes are methodologies to help keep the cost of developing software down. Even though this is an introductory level course, it will be worthwhile to consider your own time as a resource that you want to spend wisely. If you are not following a methodology, and if you are not monitoring your time, and if you have no responsibilities to anyone for how you use your time, you are probably wasting large amounts of time.

In this chapter, we'll discuss some practical strategies for students to apply when building software. The topics of software development and design could easily occupy entire courses, and it is not our goal to learn everything about them. The purpose of this chapter is to give students tools to apply when building software for small to moderately-sized projects, including undergraduate assignments and projects.

### 5.1 Software development processes

Software *development* is usually defined to include the whole life-cycle of a software system, beginning from initial intent, and concluding with deployment and maintenance. Within that cycle, there is a activity called software *design*, which is the process used to plan the software system.

### 5.1.1 Waterfall model

A very basic methodology is called the *waterfall model*. The waterfall model outlines all the activities relevant to software development, giving them a sequential order. It's called *waterfall* to emphasize the commitment placed on the direction of progress; going backwards is really not an option, and for this reason merely saying *sequential* lacks this implied commitment. The waterfall model outlines the following stages:

1. Requirements. At this stage, you figure out what the software is needed for, and what needs to be done. The result of this stage is usually a requirements document indicating what the software must do. The larger the project, the more important this stage is. It is essential that all parties (team members, clients, management) agree on the requirements. We'll say more about requirements as we go.
2. Design. At this stage, the developer(s) come up with a plan describing how to meet the requirements. This is called the design. The design may not include much code, but should present a clear set of objectives and approaches for building the software. The design can be approached as the high level aspects of iterative refinement, without fully completing the application. The output is a specification document, which can be analyzed and scrutinized. The more time spent at this stage, the fewer surprises will come up during the implementation.
3. Implementation. This stage is all about writing the code for the system. Start with the low level functional units, and build upwards. Part of the implementation phase includes testing of the software. The implementation is complete when the design plan is completed.
4. System Testing and Verification. At this stage, the application is checked, to ensure that it meets the requirements.
5. Maintenance. The application has to be delivered, either to a client, to an app store, or put on some kind of storage medium for brick-and-mortar sales. Releasing software into the world is just the beginning of the work! Software to install the application might need to be written. Updates and upgrades need to be planned. Bug reports need to be collected, and bugs need to be fixed.

The waterfall model is straight-forward, and it names all the major tasks in a software development project. The "waterfall" part suggests that there's no going back to a previous stage, which is probably not realistic. Moving along sequentially will work for fairly small projects, but it is not very effective for very large projects. The main problem with the waterfall model is that requirements change very quickly. Another big problem is that for very large projects, simply agreeing all the requirements in advance may be infeasible.

Despite its weaknesses, it is helpful to consider the waterfall model applied to your homework assignments. It gives names to tasks that you didn't even realize you are doing. When you read an assignment question, you are reading an informal requirements document; when you ask questions about it, you are clarifying the requirements. When you think about how to solve the problem, you are doing design and implementation; perhaps you are not in the habit of separating these two tasks, though. Before you hand your program in, you're doing system testing and verification. The waterfall model can help you monitor how much progress you are making, and it can help you to be objective about your efforts.



### 5.1.2 Iterative and incremental model

The *iterative and incremental model* of development is based on the tasks identified in the waterfall model, but is slightly more flexible. Rather than figuring out all the requirements in advance of all development, the iterative model starts with a small subset of the full requirements, and applies the waterfall model to it, repeating steps 1-4 (above), as needed. With a small system built, you can get feedback from your client, and more requirements can be added. The system gets larger incrementally, and when the full set of requirements are documented and achieved, the process can enter the maintenance phase of the waterfall model.

The *iterative and incremental model* allows a large project to be constructed incrementally, and getting feedback from the client or customer is useful sooner than later. The drawback of the incremental model is that the requirements can still change, and without careful planning, requirements added later may be inconsistent with earlier versions of the system.

It may be valuable for students to apply the incremental model to assignments and projects. The key thing would be to understand the requirements well-enough that you can iteratively add requirements by the time it is due, so that you can obtain credit for your work. The risk is that you might iteratively step away from a good product, and not have enough time to correct your development trajectory. Still, students benefit from the advice of always having something working, even if it is not complete; testing as you go, rather than all at once, and don't try to build the whole project all at once.

### 5.1.3 Requirements and Specifications

One of the key tasks in the software development process is establishing the requirements. Software requirements indicate, as precisely as possible, what the application or script is supposed to do. For example, every programming question you've seen in a computer science course starts with a description of what you're supposed to do; your instructors have provided the requirements for you. Outside of homework, gathering the requirements for a project is a normal part of accomplishing the task.

This task almost always starts with an English (or other natural human language) description of the requirements. Requirement specifications written in human natural languages are often slightly vague, and may leave out some information. Even professionals start out with a vague idea about what they are trying to do, and they have to work on it to make it more precise. When a client or a colleague comes to you with an idea, moving from vague statement to precise specification is a key part of the collaboration. Sometimes your colleague or client will make contradictory requests, without realizing it. Sometimes they won't even know what they want.

After the requirements are written, the application needs to be designed. We'll have more to say about this task in the next section. A design document should include specifications for functions or modules (and other abstractions we'll see in this course, such as objects and ADTs). These are sometimes called specifications. A specification is a precise statement about what a function or module should do, including arguments, inputs, outputs. A specification document may suggest an algorithm at a high level of abstraction, but the specification is a document that one uses to assess the implementation, so it should not be expressed in code. A specification is the precursor to writing code.

Again, the natural starting point for specifications is to describe them in natural human language. We can try to be fairly careful, but specifications written in English will always be vague unless they are ponderous and impenetrable. Human natural language is not always the right language for

specifications. For one thing, specifications written in natural languages are difficult to reason about. It is far more appropriate to use mathematics to write specifications. The necessary math to describe computer science specifications is the language of sets, relations, logics, proofs, and mathematical induction. These topics are taught in second year courses; for now, we have to make do with English.

For courses and assignments, you will find requirements and specifications of varying levels of precision. As you no doubt have observed during your computer science courses, written specifications for assignments do not answer all the questions you may have, and sometimes have to be corrected because they contain mistakes. For students, it is easiest if instructors give very precise and thorough specifications. But part of the task of problem solving is to be able to refine a vague specification, and make it precise. Don't flinch when you have to do this (but you may grumble about it, of course, that's normal).

## 5.2 Software Design Strategies

Software development includes the task of designing the software system, or application. In the previous course, we said only a little about software design, leaning heavily on the concept of abstraction and stepwise refinement as a way to move from idea to Python script. Stepwise refinement is the process of moving from an abstract description of a computational task to a more refined description by adding detail. We usually start with a pseudo-code description at a very high level, and in the process of adding detail, the task begins to look more and more like a script or a program.

A very powerful approach to software design is called *test-driven design* (TDD). It can be applied the waterfall model or the incremental model during the implementation stage, that is, after you have a plan for what functions and modules you want to build. The main idea is that test code is written first, before any code is written to provide the functionality. This seems counter-productive, but it is powerful for a number of reasons. First, if you can't write a test case for a function you are thinking about writing, you don't understand it well enough to start. Second, with careful design of your test cases, you can know exactly when your code satisfies the requirements; you won't be writing extra code. Third, TDD keeps developers from writing a lot of code all at once, and adding multiple errors, that need time to be debugged. The time you spend alternating between test code and application code is more than made up for in reduced debugging time.

The TDD process is based on three principles:

1. Test first: Do not write a line of code until you have written a failing test for it.
2. Do not write more of a test than is necessary to cause the test to fail.
3. Do not write more code than is necessary to pass the failed test.

The idea is that you develop the function and the tests incrementally, starting with a trivial test that will fail because you haven't written the function yet. Then you write just enough of the function to pass the trivial test. When it does, you add another test, which must fail because your function doesn't do anything yet. Then you write enough code in the function to pass the test. Keep adding to the testing, incrementally, and keep adding to the function, to satisfy the tests.

Test-driven development can be directly applied to projects of small to moderate sizes. To apply TDD strategies, you must have a clear requirements document already in place. You also have to have a design plan, identifying, at a high level, the functions needed to achieve the requirements; each function must have a specification, written fairly precisely. With this in hand, the TDD approach can be applied to the implementation stage of the development process. Start with the low-level

functions (the ones that do not call other functions in your design). Apply the TDD strategy to the implementation of this function, starting with trivial tests, and moving towards completion. When you've tested every part of the specification, and your implementation passes all those tests without error, you move on to the next function.

### 5.3 Version Control

Software development and design is difficult even for experts. There are design decisions to be made before they can be implemented; it is not rare for decisions to be reversed some time in the future. You've probably experienced the problem yourself. You had a program almost working, then you added or changed some code, and made it worse. Now you are faced with a choice: you can undo your changes and try something else, or you keep those changes, and make more. It's common for both options to feel risky, because of the effort you'll need in the future if you make the wrong choice now. You could make a backup copy of your code in its current state, so that if you had to, you could come back to it and try something else. And if you have done this several times, you are faced with another problem: which of the multiple backup versions do you revert to?

Version control is the name for software tools that allow content providers to manage their project content and all their backups in an organized way. Each backup is called a *version*, which is maintained by the version control software. The entire history of all versions of the project is stored carefully, and each version can be documented with a note that summarizes why the version was created. Version control allows a content provider to switch to any stored version at any time, without searching through folders of backup files. Multiple versions can be compared, to see differences in the content. Version control allows teams to collaborate, and helps them keep consistent versions of the project as they work on different parts of the project at the same time.

Version control tools provide a lot of functionality, but we'll focus on the basics. When you start a new project, you *initialize* version control for that project. Every time you accomplish part of the project, you document what you did with a short note, and you tell the version control to store the current state of your project as a *version*. Version control can help you navigate the history of the versions you stored; you'll want your short notes to be descriptive and helpful. You can use version control to compare two different versions already stored, or the difference between the current state with any previous version. If you ever need to return to a previous version, you identify which version you want, and tell the version control to restore that version.

The value of learning version control is that you can safely record the entire history of every assignment question or project, as you build it. If your project is under version control, you can make experimental changes, and discard them with a single button click, getting back to your previous version. If you are using the incremental and iterative development model, you can earmark every feature you add to the project, and if necessary, roll-back to a previous version. You can add or delete code without hesitation, because you can always call it back at any time, if you have been using version control to record your version history.

### 5.4 Prototypes

A software prototype is an application intended to demonstrate some subset of the requirements of a software project. Prototypes are valuable to the process of software development because they allow the exploration of design decisions (stage 2 of the waterfall model), before committing to them in an implementation (stage 3). They also allow developers to get feedback from the clients

before the final application is complete. Software professionals have wide variety of strategies for using prototypes. A prototype can evolve into a final product by a process similar to the iterative and incremental model discussed above. Sometimes a prototype is developed quickly with the intent to discard it before the implementation of the final project. For this strategy to be effective, a prototype has to be developed very quickly, and the developers have to learn enough from the effort to help them make better design choices for the final implementation.

The value of the prototype concept to us in this course is the idea that it is not foolish to try a few ideas quickly, discover some hidden implications of those ideas, and then throw them away completely to try something new. Students typically commit themselves to a single path of development, and are naturally quite reluctant to make drastic decisions like deleting whole pages of code. With the knowledge that prototyping is a real strategy, and the confidence that version control has the complete history of your work, you should learn to throw code away and try something else.

## 5.5 Summary

There is still a lot more to say about software development and design processes, and those things must be left to more advanced classes. Take to heart the fact that software development is difficult because it is so easy to make mistakes, and you need to allocate time to managing those mistakes deliberately. If you assume that everything will go well, that no errors will occur, that your understanding of the requirements is perfect, or that they are perfectly stated, you will find yourself underestimating the time it takes you to complete your work. By monitoring your progress, by allocating time for things to go wrong, and by careful planning, you can make far better use of your time. By using version control, you have the freedom to try out ideas, without committing to them, and without making a total mess of your project with ad hoc copies and code that's been commented out. You don't need your development and design plans to be perfect at this point; you simply need to practice using them. This course is intended to teach you to develop software; the answers to assignments are not as important as the process by which you answered them.

### Procedural Abstraction

The benefits of procedural abstraction  
Multiple levels of procedural abstraction

### Describing a Function: Interface Documentation

Documenting the interface after defining the function

### As part of the implementation process

### Using Procedural Abstraction in the Design Process

## 6 — Procedural Abstraction

### Learning Objectives

After studying this chapter, a student should be able to:

- Define the concept of procedural abstraction.
- Give an example of procedural abstraction, and explain how it helps meet design and implementation goals.
- Describe the components of a function's interface documentation.
- Write interface documentation for simple functions.
- Describe the role of procedural abstraction as part of the implementation process.
- Describe the role of procedural abstraction as part of the design process.
- Apply procedural abstraction during stepwise refinement.

### 6.1 Procedural Abstraction

Computer science uses abstraction to help manage the complexity of the software we try to build. Without some strategy to manage complexity, software design is prohibitively difficult. Computer scientists in the 1970s were very alarmed that their ambitions for complex software were beyond their ability to produce. There was a concerted research effort to invent concrete techniques for abstraction for the software design process that would allow software developers to create the kinds of software we see today. We will look at two very specific forms of abstraction: *procedural abstraction* in this chapter, and *data abstraction* in Chapter 8. Both of these concepts are useful to help achieve the design and implementation goals described in Chapter 4.

When we define a function, we are creating a *procedural abstraction*. A function encapsulates an algorithm, i.e., a sequence of instructions that achieve a purpose. More importantly, a function definition separates the implementation (its instructions) from the function's use (its function call). Procedural abstraction allows us to look at software from two different perspectives: implementation,



i.e., *how to do it*, and purpose, i.e., *what it's for*.

### 6.1.1 The benefits of procedural abstraction

Procedural abstraction enhances *reusability*: a function containing encapsulated instructions can be called multiple times without copying the instructions. Novice programmers might prefer to copy a block of code, rather than invest the time to define a function to encapsulate that block of code, and make it reusable. But copying blocks of code carries the risk of copying errors that you have not detected yet. Your work increases if you are forced to debug multiple copies of the same code. On the other hand, re-using a function multiple times may cause the same error to happen multiple times, but if you fix the function, you've fixed it once-and-for-all, which is far better for you!

Procedural abstraction enhances *adaptability*: changes to the encapsulated code can be made often without affecting any other part of a program. This also facilitates *efficiency* improvements indirectly: encapsulated code can be improved, or even replaced with something much more efficient, and every call to this function benefits from the improvement! It enhances *correctness*: a function or procedure can be the target of focussed testing; this can often be accomplished without having to test the whole program all at once. It allows us to use variable names in the scope of a function without conflicting with other variables of the same name in other functions. Most importantly, once written, we can use the function by remembering its purpose, without having to keep the implementation in mind.

### 6.1.2 Multiple levels of procedural abstraction

Figure 6.1 gives a visual representation of the different levels of procedural abstraction we work with in Computer Science. The highest level of abstraction is a statement of purpose, and the lowest level of abstraction is the purpose fully realized as behaviour on a computer. It is possible to consider even lower levels, e.g., the quantum behaviour of sub-atomic particles as they interact in the computer, but the author's diagramming skills were not up to the challenge.

Between these two extremes is a continuum, and the diagram gives some key concepts within this continuum. You'll (hopefully) recognize an incomplete snippet of Python code near the middle of the continuum; the precise location is not really important, only that it's between the two extremes. The activity of software design usually operates in the top levels of the continuum, although it is not uncommon for lower-level concerns (about the computer hardware, or the operating system) to have an impact on software design. In this course, we're likely to stay higher up in the continuum; the lower levels are just as important, but we can't learn everything all at once. There are courses that will teach you more about the lower levels later on.

Another important position highlighted in the continuum diagram is labelled *interface*. A *function interface* defines a function's input-output relationships: the function's name, its input, and what it returns. Every function implemented (in any programming language) has an actual interface simply because the programmer has defined the function.

## 6.2 Describing a Function: Interface Documentation

The documentation of the interface for a function should provide enough information for someone to use the function without needing to study the implementation. In this course, we adopt a strict practice, so that the important aspects are always included, and there is no uncertainty about what is required:



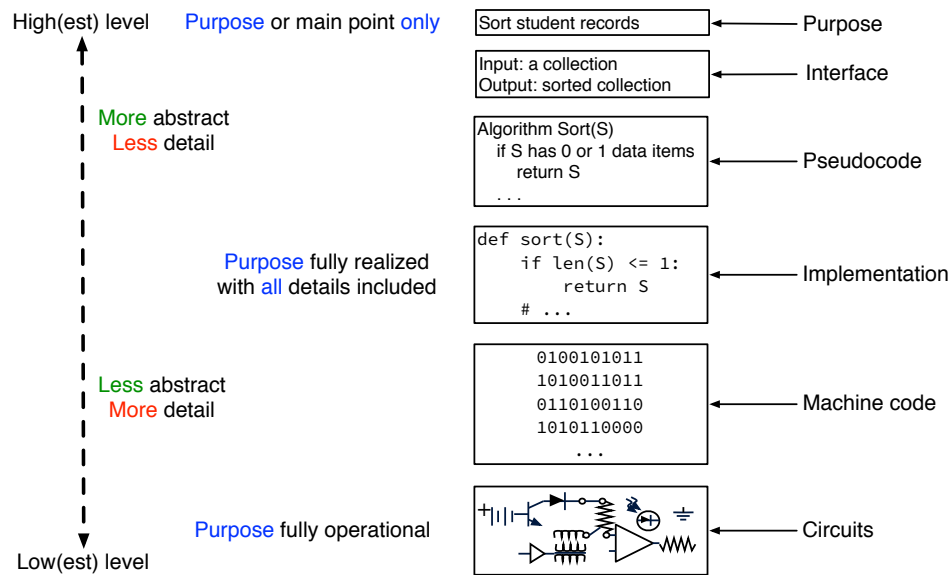


Figure 6.1: Abstraction is a continuum, with many levels. Computer Science can be studied entirely in the top half of the continuum, but it's only a practical science because of the efforts made in the levels represented in the bottom half of the diagram.

**Practice 6.1** Every function's interface documentation will consist of four parts, as follows.

**Purpose:** A brief statement of the purpose of the function, in plain language.

**Pre-conditions:** A description of every parameter for the function, including assumptions or restrictions on their values.

**Post-conditions:** a description of any effect produced by the function that could be detected after the function returns, e.g., changes to data defined outside the function itself, or any kind of I/O; optionally omitted if the function has no effects.

**Return:** a description of the value returned; optionally omitted, if a function returns no value.

Whether you write your function in Python, pseudocode, or any other programming language, these are the aspects that professionals will want to see in your interface documentation. Here's how our convention looks in Python:

```

1 def sqrt(x):
2     """
3     Purpose:
4         Approximate the square root of a given number, x.
5     Pre-conditions:
6         x: A non-negative number
7     Post-conditions:
8         (none)

```

```

9      Return:
10         a number y such that y*y is close to x.
11         Note: if x is negative, the value None is returned
12     """
13     if x < 0:
14         return None
15     else:
16         y = 1.0
17         while abs(x - y*y) > 0.001:
18             y = (y + x/y)/2.0
19     return y

```

The pre-condition gives a restriction on the value for the parameter,  $x$ , indicating that it must be greater than or equal to zero. In the above example, we wrote (none) under the heading *post-conditions* to indicate that the function does not have any effect apart from returning a value. A function with no effect other than returning a value is sometimes called a *pure function*. Notice that the return value is described informally, but very carefully, and in a way that allows the return value to be checked.

In Python, it is convention to put the interface documentation in the doc-string of the function. This will allow anyone to access your interface documentation in Python's interactive mode, as in the following example:

```

>>> help(sqrt)
Help on function sqrt:

sqrt(x)
    Purpose:
        Approximate the square root of a given number, x.
    Pre-conditions:
        x: A non-negative number
    Post-conditions:
        (none)
    Return:
        a number y such that y*y is close to x.
        Note: if x is negative, the value None is returned

```

As the example shows, the `help` function finds the function definition, and extracts the doc-string, printing it to the console. The format of our documentation is a convention, and you will be expected to use it for every function you write, unless instructed otherwise, e.g., in certain examination questions. Every computer science course may have slightly different expectations, and every software development project team will have different formats. Once you get in the habit of documenting your work, the format is a detail.

### 6.2.1 Documenting the interface after defining the function

In a course such as this one, it is very common for assignments and tests to ask students to implement a function to achieve some purpose, either to practice some concept, or to demonstrate their mastery of a concept. What often ends up happening is that the assignment question usually does a pretty

good job of describing the function, and its interface. It is also quite natural for students to implement the function as the primary intellectual activity, and then, as an after-thought, and perhaps a bit resentfully, go through the trouble of documenting the interface. To this, the author says “Fair enough!” But remember that assignments and tests are artificial pedagogical devices, not real-world problems. We have to demand professional practices in artificial environments so that you can bring those practices with you in more natural, less constrained environments. Be diligent in documenting your interfaces, and learn to do it first, not last (and certainly not at the last moment), and learn to do it quickly. This is a practice that will help you a lot in the future, even though it may seem to be a burden right now.

### 6.3 As part of the implementation process

One of the ways you can use procedural abstraction is to take some code that you have already written, and make it the body of a (new) function. For example, we could start with the code from Chapter 2, and we might simply wrap a function definition around it:

```
1 def selection_sort():
2     """
3     Purpose:
4         Displays the list [0, 1, 2, 2, 2, 3, 5, 6, 7, 8, 8]
5         to the console.
6     Pre-conditions:
7         none
8     Post-conditions:
9         A list appears on the console.
10    Return:
11        none
12    """
13    unsorted = [3, 2, 5, 7, 6, 8, 0, 1, 2, 8, 2]
14    sorted = list()
15    while len(unsorted) > 0:
16        out = min(unsorted)
17        unsorted.remove(out)
18        sorted.append(out)
19    print(sorted)
```

The interface to this function is very simple: no inputs or return, and a very specific effect, namely some output appears on the console.

It’s not a good function, because the programmer makes too many decisions. Most importantly, the programmer decided which list to sort (see line 13); we say that the data is *hard-coded* into the function. At the end of the function, we see the decision to display the sorted list to the console. In between, the programmer decided which algorithm to use to sort the list, which is really the only decision that makes any sense in a function called `selection_sort`.

With a few edits, we can take this as a starting point to generalize the function and make it more useful. First, we can add a parameter to the interface, so that the function’s input is the list to be sorted. This allows the decision about which list to sort to be made outside the function. Second, we

can return the sorted list, which allows the decision about what to do with it to be made outside the function as well.

```
1 def selection_sort(unsorted):
2     """
3     Purpose:
4         Create a list with values in increasing order.
5     Pre-conditions:
6         :param unsorted: a list of numbers
7     Post-conditions:
8         All values are removed from unsorted.
9     Return:
10        a list of all the values from unsorted,
11        but in increasing order.
12    """
13    sorted = list()
14    while len(unsorted) > 0:
15        out = min(unsorted)
16        unsorted.remove(out)
17        sorted.append(out)
18    return sorted
19
20
21 data = [3, 2, 5, 7, 6, 8, 0, 1, 2, 8, 2]
22 data_sorted = selection_sort(data)
23 print(data)
24 print(data_sorted)
```

The interface documentation is careful to describe a significant aspect of the function: namely that the function empties out the original list! You might need to review Chapter 3 to understand how this side-effect happens. In short, this function's documentation tells us that "You can't sort your data and keep it, too." This change to the original list is an example of a *side-effect*. A side-effect is any action that has an effect outside the function; we typically don't include actions that only affect the function's local variables. Side-effects include actions like modifying a global variable, printing to the console, writing data to a file, adding data to a list or dictionary, among others.

**Practice 6.2** If a function has a side-effect, it is absolutely necessary that the side-effect be documented as part of the *Post-Conditions*. ■

Imagine using this sorting function (as part of a module, say) and not being told of its side-effect. You might easily be faced with the mysterious disappearance of all your data, and it might take you quite a long time even to suspect that a sorting function was causing it.

A few more edits, and we can arrive at the following variation:

```
1 def selection_sort(alist):
2     """
3     Purpose:
4         To put the values of alist in increasing order.
5     Pre-conditions:
6         :param alist: a list of numbers
7     Post-conditions:
8         Values in alist are in increasing order
9     Return:
10        (none)
11    """
12    sorted = list()
13    while len(alist) > 0:
14        out = min(alist)
15        alist.remove(out)
16        sorted.append(out)
17
18    alist.extend(sorted)
19
20
21 data = [3, 2, 5, 7, 6, 8, 0, 1, 2, 8, 2]
22 selection_sort(data)
23 print(data)
```

Notice that the documentation indicates that original list is being modified, and the code verifies this. You might need to review Chapter 3 to understand why we used the list method `extend()` rather than an assignment statement like

```
alist = sorted
```

There are several lessons we can learn from this example. First, making a function by encapsulating code that you've already written is a perfectly normal implementation activity. In fact, it is so common that many interactive development environments (IDEs, such as PyCharm) have *refactoring tools* that assist developers in making these kinds<sup>1</sup> of changes. Second, the example emphasizes procedural abstraction as a technique for generalization; a function is more general if more of the decisions about what should happen inside the function are made outside of the function. For example, a function that returns the sorted list is more general, and therefore more useful, than a function that displays the sorted list to the console. Third, it suggests a very simple problem solving strategy: when faced with a computational problem, try to solve a very concrete instance of the problem, just by writing a script. When the script works, use procedural abstraction to make a general tool out of it.

---

<sup>1</sup>There is a whole application menu devoted to refactoring tools in most IDEs. Check it out in PyCharm.

## 6.4 Using Procedural Abstraction in the Design Process

In the previous course we introduced a simple design process called *stepwise* or *iterative refinement*. Starting with a very abstract statement of the purpose of the application, the iterative refinement process keeps adding detail to abstract statements, either by adding information to a statement, or by breaking an abstract statement into a sequence of somewhat more detailed statements. Procedural abstraction can play a valuable role in this process. Each abstract statement in the iterative refinement process could be replaced by a function call to a function that has not been implemented yet. The power of procedural abstraction comes into play because you can talk very intelligently about a computation in terms of its interface (what its purpose is, and how to call it) without having an implementation. The first cut of an implementation is to write down the steps you think you'll need to accomplish the purpose. For many development problems on the level of homework assignments, you can design an application by starting with the names and purposes of a few very high level functions, without having to implement them. You can design the interfaces of these functions by deciding how they will work together, without having to decide on the algorithm for each one first. The process of iterative refinement can give you direction when you're not sure where to start, and it can help you focus on one task at a time.

**Principle 6.1** The definition of an interface is an essential part of any design process. ■

To design an interface, even if we're not exactly sure what the algorithm should be, we basically have to write the function's interface documentation. We start by writing a statement of purpose, clearly and concisely. Then we can add pre-conditions, post-conditions, and a description of what the return value should be, again, as carefully as possible, without ever worrying about what the algorithm will be. Remember that this is a design process, and that you can treat it like a work-in-progress: include everything you can think of now, and as you spend time working on it, you may realize you have to make changes.

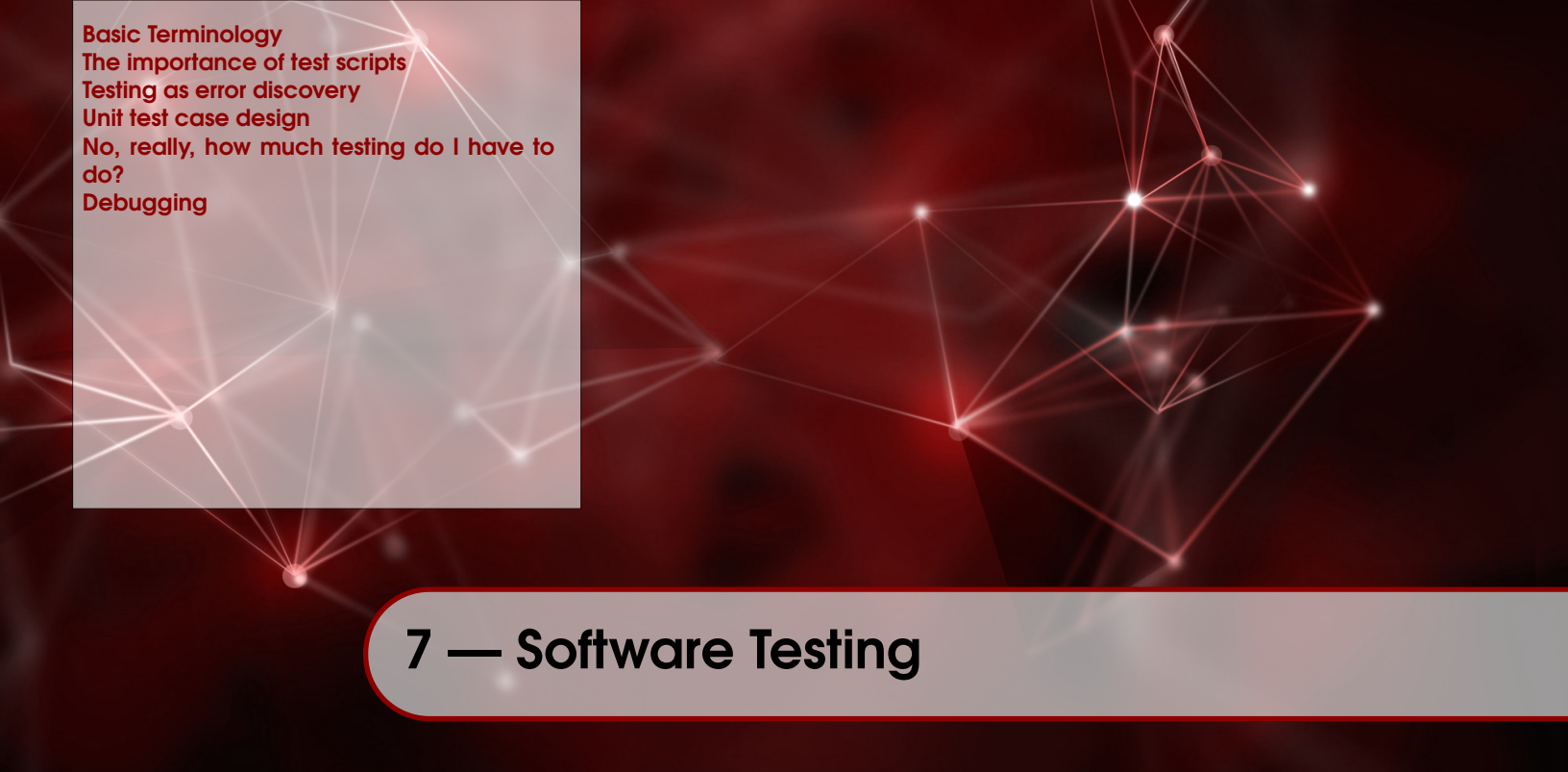
**Practice 6.3** When applying the iterative refinement design process, define function interfaces before you define function bodies. ■

When complete, the interface documentation should provide enough information to allow a software developer to implement the function or procedure, without requiring any further information about how it would be used. It should also provide enough information to allow a software developer to know how to use the function properly, without studying the encapsulated instructions. And of course, you should be able to write black-box test cases from the interface documentation alone. See Chapter 7.

Good interface documentation serves a number of other purposes, beyond communication. First, it is a good indicator for understanding the problem. If you cannot write down the purpose and the pre-conditions, then you do not understand what you are trying to do, and you need to think a bit more carefully before you start writing code.

Second, designing an interface by writing the interface documentation helps you focus on a single task. You can make a number of decisions about your function just at the interface level, and these will guide the rest of your problem-solving. Even if you haven't worked out the interface completely, focussing on it will help. For this reason you should treat the interface as part of the problem-solving.





Basic Terminology  
The importance of test scripts  
Testing as error discovery  
Unit test case design  
No, really, how much testing do I have to do?  
Debugging

## 7 — Software Testing

### Learning Objectives

After studying this chapter, a student should be able to:

- Provide examples of software errors resulting in observable faults.
- Define the terms error, fault, specification, and oracle.
- Compare black box testing and white box testing.
- Explain the differences between unit testing, integration testing, and system testing.
- Identify unit test case equivalence classes for a given function.
- Explain why the creation of correct program components is important in the production of high-quality software.
- Apply a variety of strategies to the testing and debugging of simple programs.
- Construct, execute and debug programs using a modern IDE and associated tools such as unit testing tools and visual debuggers.
- Construct and debug programs using the standard libraries available with a chosen programming language.
- Explain the importance of test coverage.

Software testing is a professional obligation that arises from the fact that normal human behaviour includes making mistakes. A single spelling mistake in an English essay is an inconsequential accident; a single grammatical mistake in the same essay could lead to a slight misunderstanding of the writer's argument, so that's a tiny bit more consequential. But there are serious consequences for software applications that get released with undetected errors. Software errors can be harmful to people, or may cause damage to equipment or the environment, as demonstrated by example:

- In the 1980s, the **Therac-25** was the first medical radiation therapy machine completely controlled by software. Due to a combination of factors, including the presence of software errors in the input module, several cancer patients received massive radiation overdoses, leading to disability, and at least 3 deaths.

- During the 1996 launch of ESA's \$370M **Ariane 5 rocket**, the autopilot software obtained a very large value from an accelerometer. The software mistakenly rejected the value as impossible, causing the autopilot program to halt abnormally (a software crash). Without an autopilot, the rocket veered out of control and disintegrated 37 seconds after launch.
- In 1998, NASA launched their \$125M **Mars Climate Orbiter**, but it failed to establish a stable orbit around Mars, and crashed. The control software used on Earth sent control signals to the orbiter using imperial units (lbs, feet, Fahrenheit, etc), but the orbiter itself was programmed to receive those signals in metric units (grams, meters, Celsius, etc).
- Researchers at the **Scripps Institute** borrowed some software for the analysis of protein crystal structure from another lab, and published five significant papers in very prestigious journals. Then they discovered that the software had one function whose output flipped the sign of some numeric values, invalidating their entire analysis. The team had to issue public retractions for their papers in these highly respected journals.

It's very important to realize that these examples, while serious and dramatic but few in number, do not represent the full extent of software errors.

All software developers, from novices to experts, make mistakes; there is no level of training that one can reach which guarantees a professional will never make a mistake. Expert developers will make fewer silly mistakes, but they never stop making them. On the other hand, the mistakes that an expert developer will make will sometimes be subtle and hard to find.

To detect software errors in a computer program, the best and only strategy is good testing. Software testing is an important topic in the practice of software development. There are advanced courses in software engineering that delve deeply into the topic of testing. This chapter will give you some tools that will help you test software for courses and projects before you get to those more advanced courses.

Testing skills should be learned and practiced when stakes are low, and programs are small, so that when the stakes are higher, and projects larger, students already possess proven testing skills. The cost of software errors you make in assignments and course projects will be measured in terms of the time spent you spend debugging. For the kinds of small programming assignments you've already seen in first year, a little testing may have seemed good enough. Even if you didn't spend a lot of effort testing, you probably still got your code to work well enough to get a decent grade. When projects are larger, errors can have a disproportionately large impact on your work: a single, small, undiscovered error could prevent your code from working properly. The more time it takes for you to find and fix errors, the more likely it will be that you won't be able to complete your work. The consequences could be a missed deadline, a poor or failing grade, or a very long time spent debugging.

## 7.1 Basic Terminology

It turns out to be rather difficult to define the concept of a software error directly. As you may have already realized, a software error is detectable only from its effects, which allows us to make the following definition:

**Definition 7.1** A *software fault* occurs when the software does not do what it is supposed to do, or when it does something that it is not supposed to do.

A fault is behaviour that can be observed. It's a symptom, an indication that something went wrong,

but it may or may not tell you exactly what caused the fault. In our work, the fault may be that the script halts and prints an error message, or it might be an incorrect value displayed on the console. As a symptom, a fault has a cause, which is what we call an *error*.

**Definition 7.2** A software *error* is any code that causes a software fault.

A single error might cause a cascade of observable faults, one after the other. Also, the code that contains the error could be quite far away from the code that makes the fault. We sometimes call a software error a *bug*.

**Definition 7.3** *Software testing* is the active process of fault discovery.

Software testing is the application of the scientific method to a software design project. The hypothesis we make during testing is that the software is correct. Testing is a form of experiment that tries to prove the hypothesis wrong. Testing is supposed to reveal the presence of errors by creating situations that would result in a fault if there is a certain kind of error. In the previous course, we called these test cases.

**Definition 7.4** A *test case* checks the behaviour of a program or a function for a very specific input or context. A test case consists of:

- The input values (or other context) to try
- The expected or right answer
- A brief description of the kind of thing the test is intended to check.

A test case is designed to cause a fault if the code has an error, by checking the effects of some computation against a known, right answer. When a test case causes a fault, we say that the test case *fails*; if no fault is caused by the test case, we say it *passes*. Sometimes, the right answer or behaviour for a function or program is obvious, based on common sense understanding of the code's purpose. But this is not always true.

**Definition 7.5** An *oracle* is any method that can tell us whether we have the right answer for a given test case.

Computer science uses the concept of an oracle in several interesting ways, but for now, we're only using the term to refer to a source for a right answer during testing. An oracle can be external to the software, like an expert providing expert knowledge, or an assignment description providing a few examples of inputs and correct outputs. An oracle could also be a function whose only purpose is to check the result of a function being tested.

In our first year courses, we present assignment questions (and exam questions) as clearly as we can; these are specifications.

**Definition 7.6** A *specification* is a written statement or document describing what a program is supposed to do.

In light of Chapter 6, the specification of an assignment question could be used as the basis for a function's interface documentation. In a larger software development project, the specification

document is written, perhaps in consultation with a client, before the design and implementation can begin. In a professional software development setting, developing a correct specification document is just as important as developing a program. See Chapter 5. It's important to understand that a specification document cannot be an oracle except for very simple, almost trivial, programs; a specification is an abstraction of all the oracles for all the functionality of a program.

## 7.2 The importance of test scripts

You may already have a practice of informal testing, e.g., running your program on a few inputs and checking your output by hand. As your projects get bigger, informal methods are no longer effective: if your informal testing consists of typing a few commands by hand, you'll have to repeat the typing to repeat the informal testing. It's important to see the computer as a tool for testing, so that you can run a whole barrage of tests very quickly, and you can repeat the testing very easily.

A test script (or test driver) is a program that tests other programs. In the previous course, a test case was designed by considering its components, and implemented using Python. For example, if we have a function named `sqrt()`, we might design a test case, and implement it as follows:

Input(s):	4
Output(s):	2
Reason:	Positive integer with integer result.

```
1 # sqrt(4) == 2
2 result = sqrt(4)
3 # expected output: 2
4 if result != 2:
5     print('Error: sqrt(4) returned', result)
```

We will continue to recommend that test cases only report faults, and remained silent on successful tests. However, this is not a universally accepted practice. Sometimes it is useful to know when a test passes.

The scripts we wrote for the previous course represent the very simplest forms of testing, in which one could easily compare the output of a function to the known, correct value. This form of testing is called *unit testing*, in which you test a single unit of code, e.g., a single function or method. When testing is done on functions or modules working together, it is called *integration testing*. Finally, when completed applications are tested, it's called *system testing*.

## 7.3 Testing as error discovery

Software testing is the application of the scientific method to a software design project. The hypothesis we make during testing is that the software is correct. Testing is a form of experiment that tries to reveal the presence of errors.

The logic of unit test case design goes like this: *If the function is correct, it will return a specific value Y given specific argument X.* For example, if the `sqrt()` function is correct, it will return the value 2 given the input argument 4. If the function returns an incorrect value for a given input argument, then we have to conclude that the function is not correct. On the other hand, it is completely logical for an incorrect function to return the right answer for any single test case. For example, consider the following function definition:

```
def sqrt_a(x):  
    """  
    Purpose:  
        Approximate the square root of a given number, x.  
    Pre-conditions:  
        x: A non-negative number  
    Post-conditions:  
        (none)  
    Return:  
        a number y such that y*y is close to x.  
        Note: if x is negative, the value None is returned  
    """  
    return x / 2.0
```

This function is not correct. However, the function passes the test case from the example above.

A single failed test case is enough to guarantee that a function is not correct, but a single test case that passes is not enough to guarantee that a function is correct. When writing test scripts, it is not good enough to demonstrate that the function is correct; we have to do our best demonstrate, through careful design of test cases, that the function cannot be incorrect.

This example also shows how important the interface documentation is. In context, it is clear that the square root of  $x$  is not  $x/2$ , except for a couple of special cases. The interface documentation gives us enough information to know what the function's purpose is, and what inputs it expects, and what its output is. Consider how we might test the following function:

```
def f(x):  
    return x / 2.0
```

There is no interface documentation, so we don't know what the function is intended to do. All we can see is what the function actually does. Testing a function without any interface documentation is useless, because all we can hope to establish is that the function does what it does.

**Practice 7.1** A function cannot be considered correct without proper interface documentation.

## 7.4 Unit test case design

In the previous course, we presented the concept of *black-box testing*, where test-cases are designed based entirely on the interface of a function (inputs, outputs, and effects); or more generally, on the purpose of the function. The term *black-box* implies that the test-cases were designed without looking at the code that is being tested. Similarly, we presented the concept of *white-box testing*, where test-cases are designed based entirely on the lines of code that make up the function or program. White-box test cases are designed to force execution of certain lines of code, or certain paths through the code, such as checking complex conditionals or nested loops. When you look at a test case, it may not be clear whether it is an example of *black-box* or *white-box* testing. These terms do not define what a test case looks like; rather, they indicate strategies for coming up with test cases.

Most functions or programs can be applied to infinite combinations of inputs, and you can't test all possible combinations. Many complex functions or programs will have a very large number of

possible paths through the code, and it may be impossible to attempt to identify them all, much less test them all. Even if you test all the paths through your current program, your program may be incorrect because of code that is missing from your function.

To address the the problem of a seemingly infinite number of unit test cases, we introduce the idea of test case *equivalence classes*. In general, an equivalence class is a set of things that are equivalent from a certain point of view. In our case, we're trying to come up with a set of inputs that are similar from the point of view of testing a given function. Once we have identified that set, we can choose a couple of examples from the set that are representative of the whole set.

For black-box testing, we can make some progress by analyzing a function's specification. Equivalence classes frequently express a range of test values, depending on the data type of a function's parameters. For example, you might start with these ideas:

- Numeric parameters: E.g., positive, zero, negative; even, odd; integer, floating point.
- String parameters: E.g., non-empty strings with no blanks; non-empty strings with blanks; empty strings.
- List parameters: E.g., empty; singleton; all the same; increasing order; decreasing order; even size; odd size.

Of course, these examples only show the kinds of equivalence classes you might start thinking about, not the ones that will be right for any given function. If your program or function has multiple inputs, equivalence classes will consist of combinations of ranges.

For example, consider the `sqrt()` function. Starting from the point of view of black-box testing we can identify a number of test case equivalence classes:

- The set of numbers  $x$  for which  $x = \text{sqrt}(x)$ . Examples: 0, 1.
- The set of numbers whose square root has no fractional part. Examples: 1, 4, 9.
- The set of numbers  $x > 1$ ; for these numbers,  $x > \text{sqrt}(x)$ . Examples: 4, 23.7.
- The set of numbers  $x < 1$ ; for these numbers,  $x < \text{sqrt}(x)$ . Examples: 0.25, 0.771331.
- The set of negative numbers. Examples:  $-1$ .

You may well ask how these sets were determined. In part, from understanding the meaning and purpose of the square root operation. But also from the understanding of how numbers are represented by a computer, as well as experience coming up with test cases. It's important to understand that we're not trying to teach how to decide the perfect set of test cases, or equivalence classes, but to give students a starting point from which to start. There's no definitive right answer, and if there is a wrong answer, it's that you didn't work hard enough to find test cases. If there were a definitive right answer, someone would have designed a software script that could automatically design a full set of test cases for any given function, and which would write the test script, run it, and report the results. For now, you have to learn to design tests the same way that you learned to write programs: by writing a lot of them, and improving your skills by practice.

Equivalence classes for white-box testing are a bit easier. A function may have multiple paths through it, depending on how the arguments affect conditionals and loops. You can consider all the input values that follow the same path through a function as a single test case equivalence class. To understand these paths, check conditions using `<`, `>=`, etc. Useful equivalence classes for testing a loop include the range of inputs that cause a loop to do exactly one iteration, exactly zero iterations, and more than one iteration. If your function has more than one return statement, make sure you consider the different paths to get to each one.

Another important concept in test case development is *boundary cases* (also called *edge cases*).



Equivalence classes typically express a range of possible inputs, and some of the ranges will be mutually exclusive. More importantly, the ranges will meet at the boundaries of the equivalence classes. Values on or near the boundaries are very important test cases, and you should include test cases at the boundary, and as close as you can on both sides of every boundary. For example, in the test case equivalence classes for `sqrt()`, above, the value 0 is between the positives and the negatives; it's a boundary case. Likewise, we identified a set of numbers  $x < \text{sqrt}(x)$ , and  $x > \text{sqrt}(x)$ ; the value 1 is on the boundary between these sets. As a counter-example, there is no simple boundary between the set of numbers whose square root has no fractional part, and those whose square root requires a fractional part, so looking for boundary values between these two sets is probably not going to be productive.

Finally, use the structure of a test case (*inputs, outputs, reason*) to document where the test case comes from, and why you chose it. The most under-appreciated part of the test case is the *reason* component. When you have a lot of test cases, the reason should be useful itself to remind you what kind of error you're looking for. It's far more useful to you if it can correctly and uniquely identify the test case itself (so you can find it quickly), and it should give a clue about what kind of error occurred. When you fail to give a good description in the *reason* component, the only person who will pay the price is future you.

## 7.5 No, really, how much testing do I have to do?

For good reason, this is the question that all students ask when approaching the task of testing. The answer will depend in part on the software you are testing.

An important testing concept is *test coverage*, which is an objective measurement of the extent of the testing that has been done. Typically, test coverage is measured in terms of lines of code, number of branches, number of functions, or number of modules tested. Naturally, more coverage is better. The law of diminishing returns applies to testing: each additional test adds confidence, but the increase in confidence due to an additional test case gets smaller as you add test cases; the time cost to implement each test is roughly constant. The goal is to design a set of test cases that is small enough to avoid wasted effort, but is still effective at finding errors if they exist.

Once you have identified test case equivalence classes, and boundary cases, you can choose a small number of examples from each class. The number of equivalence classes depends on the function you're testing. A simple function will have a small number of test case equivalence classes. A complicated function will have more. This is as good an argument as there is for preferring a design that is based on a collection of simple functions, and for avoiding a design that is based on fewer but more complicated functions. It is far easier to test a simple function than a complicated function.

**Practice 7.2** Write test cases for every function you write, considering test case equivalence classes and boundary cases, using white-box and black-box strategies. Test cases should be written in a test script that is separate from the function or script you are testing, and should be executable from the command-line (UNIX). ■

## 7.6 Debugging

When your testing reveals a fault, all you really know is that your program did not perform as intended. Knowing you have a fault does not tell you how to fix the error! To find out what your

program is *doing wrong*, you need to know what your program is *doing*. Do not assume that an error has a quick or easy fix. It's a natural human tendency to hope, but the professional attitude is more cautious. Errors are *almost always where you are not looking*, so you must continuously challenge your own assumptions about what is causing the faults you've found.

It is very tempting to assume that all faults were independently caused. But this is not always true, and can cause you more trouble, not less. The problem is that you can almost always change your program to prevent any single fault from occurring without actually identifying the error. Such changes may only cover up a error that occurred earlier.

If you have a failed test case, write similar test cases within the equivalence class and at the boundaries, and see if they fail too. Once you are sure that the code makes the same mistake on the equivalence class, you can focus on the parts of the code that are associated with that class. You should always be careful not to make a change to your code too soon, because you could be introducing another error.

Try to come up with a hypothesis for the error that caused the fault. Use your hypothesis about the error to write more test cases for it. If your hypothesis is right, these new test cases should also fail before you fix the error, and succeed after. Do not make any change to your program unless you have a hypothesis, and you have verified that your hypothesis is at least plausible. You should only change the code to fix the error if the testing confirms your hypothesis.

As long as the functions are small, and their intended behaviour is well-understood, you can step through the functions with an interactive debugger, or you can do a very careful reading of the function, to try to discover the cause of the fault. For faults discovered during integration testing, there may be too much code for you to step through one line at a time using the debugger. For these faults, you need a less time-consuming option than single-stepping line-by-line through your program.

All debuggers allow programmers to set a *break-point*, which is a way of telling the debugger to run the program normally until a chosen line of code, where the execution breaks, and you can inspect the state of your data structures and variables. When you do this inspection, you should be looking for variables whose values are not what they should be, or data structures whose organization is inappropriate, or containing the wrong data. Of course, to use these tools, you need to be quite clear about what the program should be doing; if you are not, you won't be able to tell what's gone wrong!

Another technique called "wolf-fencing" often helps locate where something has gone wrong. The idea is to use print statements throughout your code, to display evidence about what your program is doing. While you can simply blanket your program with print statements, it can also be helpful to use a kind of binary search: put a print statement in the middle of your program, to see if the problem can be detected at that point. If so, you would add another print statement near the middle of the first half of the program; if not, add a print statement near the middle of the final half. Remember, the point is to locate where something has gone wrong, so that you can study a small portion of the code close up, stepping through it with a debugger. An advanced form of this technique is called "logging;" instead of printing information to the console, you'd send the input to a separate data file called a log file. The advantage is that your console is not cluttered by data, and you don't have to write or delete tons of print statements during debugging. The log output always gets sent to the log file; but you don't have to read it all the time. The disadvantage is that log files can get pretty big, and looking through them can be tedious.

Modern languages like Python and Java have runtime errors, which can be helpful in tracking down what has gone wrong. Some older languages are less helpful, and sometimes programs

will simply halt abnormally without warning and without any hint about what happened. For these languages (C, C++ to name just two), wolf-fencing, logging, and break-points are essential. Thankfully, most languages in current use have debuggers with more or less the same set of features (breakpoints, stepping, etc). No matter which language you are using, if you want to be productive and respected, learn the associated debugging tool as well.

Don't under-estimate the value of simply reading your code carefully. Software engineering calls this a *code walk-through*, and it's usually done with a group of team members. Human beings are biased against finding errors in their own code; errors based on a mistaken (or unstated) assumption are going to be invisible to you. That's human nature. Try to explain your code to someone else, even someone who does not understand what you are doing. Use a puppet. <sup>1</sup>

---

<sup>1</sup>"Experience has shown that many of the errors discovered are actually found by the programmer, rather than the other team members, during the narration. In other words, the simple act of reading aloud one's program to an audience seems to be a remarkably effective error-detection technique." Myers, Glenford J. "Chapter 3 - Program Inspections, Walkthroughs, and Reviews". The Art of Software Testing. John Wiley & Sons. © 1979.



## 8 — Abstract Data Types

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe the concept of data abstraction.
- Give an example of data abstraction, and explain how it helps meet design and implementation goals.
- Define the concept of data type in terms of values and operations, and give examples.
- Define the concept of data structure in terms of values and structure, and give examples.
- Define the concept of abstract data type, in terms of what is hidden, and what is exposed, and give examples.
- Identify the data components and behaviours of multiple abstract data types.

In this chapter we will spend some time looking more carefully at the concept of data types. Outside of computer science, data is represented by pixels on the screen, or ink marks on paper, and people generally know that numbers are for math, and words are for sentences, etc. People know that the order of digits in a number cannot be changed without drastically changing the meaning of the number (compare 3.14159 and 1.13459), but that you can move words around in a sentence, but keep the same meaning (compare “I am Sam” and “Sam I am”). People know that you can’t add the letters “ed” to a number to put the number into past tense (for example “The student 110ed.”), and that you can’t add the number 1 to a word to get the next word (for example, in the dictionary, the word after `'aardvark'` is not `'aardvark'+1`). People don’t need data types, because people are pretty smart, and if a person makes an error on a piece of paper, the error is limited to that piece of paper.

Computer programs need data types as a constraint on how data is affected by calculations in an algorithm. By limiting the kinds of things we can do to data, we can increase correctness by limiting silly errors, and we can improve efficiency by implementing fast algorithms that work on specific kinds of data (like fast arithmetic on numbers). We would not have to figure out what a programmer

might mean by writing `'armadillo'+1`). We can say those kinds of expressions don't make sense. Even if the programmer wrote `'110'+1`), which beginners do before they understand data types, we don't need to allow such operations, because they are too hard to get right. Limiting operations to data types makes computation much easier.

## 8.1 Data Types

We have already become familiar with a variety of data types. The simplest examples are numeric data types. There are integer numbers, real numbers, and, in computer science, floating-point numbers. Numbers are just values with a specific literal appearance: mostly digits. We know that numbers can be added, multiplied, squared, cube-rooted, etc; these are operations for numbers. There may be some operations that don't always make sense, such as dividing a number by zero.

**Definition 8.1** A *data type* consists of a set of values, and some operations on those values.

Integers are a data type. Floating-point numbers are a data type. They are different data types, though they have some values and operations in common. Booleans are a data type. There are just two values in that set, and the logical operators (AND, OR, NOT) are the operations.

Data types can be built into a programming language. In fact, most languages have integers, floating-point numbers, and Booleans available for use. The operations provided by the language provide the only means to manipulate these data types.

Data types can also be added to a programming language by writing programs. This turns out to be a very significant part of programming. The basic idea is to figure out what's special about the computational problem you're trying to solve, and design (or use) a data type that helps you write the program that solves it. We'll see many examples of this during the course.

## 8.2 Data Structures

A data structure is a special kind of a data type.

**Definition 8.2** A *data structure* is a data type with compound data values that are organized (or *structured*) in a particular way.

Python strings and lists are data structures: the values are compound, i.e., a collection of some kind, and there is a sequential organization imposed on the values. A Python dictionary is a data structure, too. It organizes the data in a clever way to make indexing efficient; we'll have more to say about this organization in a later chapter.

Data structures, being a special kind of data type, also have operations. You are aware of some of the list operations, at least: indexing using brackets, methods like `append()`, and `remove()`, to name a few. In this course, we'll learn about several typical data structures that are frequently used to solve computational problems. We'll learn what they are, how to use them, and how you can implement them.

The fact that data structures deal with compound data is important. For example, a list contains a sequence of values. In a textbook, or in a Python program, we represent a list using text symbols, e.g., `[1, 2, 3]`. But inside the computer the list has to be stored in some kind of computer-based



organization, so that the operations can do their job. This organization is what makes a data structure what it is.

The concepts of data types and data structures are Computer Science principles. We will see how they work in Python, but the principle extends beyond any particular programming language.

### 8.3 Abstract Data Types (ADTs)

Abstraction is important in Computer Science to manage the organizational complexity of software. We have already seen how the concept of a function can be seen as a form of abstraction, called procedural abstraction. Similarly, we can think about applying abstraction to data. *Data abstraction* is based on the idea that we can look at data from three different perspectives: its purpose, its operations, and the implementation of its structure.

You are already acquainted with Python's lists and dictionaries. We called them data structures in the previous section, but lists and dictionaries are typical examples of abstract data types (ADTs). They are abstract in the sense that we use them without having seen the Python code for the operations, nor the organization of the data inside the object. The implementation and organization details are hidden, and we do not need access to these hidden details to use them. We do know, however, that Python lists and dictionaries allow us to store collections of data; this is their purpose.

**Definition 8.3** An *abstract data type* (ADT) is a data type that encapsulates data (using code) to prevent almost all access to it, except through carefully designed operations.

Abstract data types hide the organizational structure of the data, and the implementation of the operations, but expose the interface of the operations. That means we can use the operations on the data to accomplish a task, simply by knowing what arguments to use, and what values are returned.

There are many benefits to hiding the details of organization and implementation. If the operations for accessing the data are correctly defined, then the data stored cannot be corrupted by using them appropriately. Furthermore, we can test an ADT and its operations very thoroughly, without having to do any system level testing of an application that uses the ADT. In other words, correctness is enhanced by abstraction. Robustness is improved, as the operations can be designed to manage unexpected situations. Adaptability is improved, as the operations can be standardized (like the operations on Python's lists).

The alternative to hiding these details is to allow unrestricted access to data; this leads to software that is not robust, not adaptable, and presents very grave danger of that data being corrupted. In Python, it is difficult to imagine data corruption being much of a problem, but it's very common, especially in languages like the C programming language, which was designed for speed, not safety.

In terms of reusability, there are two important points to consider. The first is application independence: ADTs allow software designers to create a concrete data structure, and provide useful operations for it, independently of any consideration about who is going to use it, and what their overall purpose is. The second is implementation independence: ADTs allow software designers to obtain all the benefits of data organization in their application through operations. Code that uses ADT operations is easier to understand, because operations describe the purpose, without describing the implementation.

The cost of using an ADT is that you have to read some documentation very carefully, and be sure that you are using the ADT and its operations properly.

## 8.4 Examples of Abstract Data Types

In this section, we'll present two very simple ADTs that are not part of Python's built-in toolset, and at the end of the section, we'll give an example of the application of these ADTs. The examples show that an ADT can be designed for some specific purposes. In other words, Python's lists and dictionaries are so useful and valuable, that using them as the primary example for the ADT concept is misleading. Programmers need to solve specific problems, but their solutions are not required to be as universally general as Python lists and dictionaries. The value of an ADT is not its universality, but its ability to simplify a program using abstraction. Universality is a good quality, but not a necessary aspect for ADTs.

ADTs can be implemented in any language, because an ADT is not a language feature; an ADT is a design strategy. All we really need is a way to encapsulate the data, and describe operations on the encapsulated data. In this course, we'll see how to use classes to define ADTs, because it's natural to encapsulate the data in an object, and to use methods for its operations. However, we could just as easily define ADTs using procedural programming concepts only. For example, we could use a dictionary-record whose keys are attribute names, and then write functions that operate on these dictionary-records to implement the operations.

**Principle 8.1** When using an ADT, refrain from accessing the data in the data structure directly: always use the ADT operations, and never take advantage of any knowledge that you might have of the data structure itself. ■

This principle will help improve program correctness, adaptability, and reusability. It takes practice and discipline. For novices, it also seems to take extra time, which always seems in short supply; be assured that the extra time you take to make use of an ADT's operations will pay off in time saved debugging your code later.

### 8.4.1 Registry

A registry is a data structure that stores a fixed number of Boolean values. You might be thinking it sounds a lot like a list or array, which is true. A registry will only store Boolean values, whereas a list can store anything. Lists can be made longer or shorter, whereas a Registry cannot. A Registry is not a Numpy array, either, and will not respond to the kinds of operations you might apply to a Numpy array.

A Registry has four operations. We will describe the operations in object-oriented terms.

- `Registry(n, bval)` creates a new registry object, giving it exactly  $n$  elements, each initialized to the Boolean value `bval`. The new registry object is returned.
- `reg.set(i)` stores the value `True` in the given Registry object `reg` at index  $i$ .
- `reg.reset(i)` stores the value `False` in the given Registry object `reg` at index  $i$ .
- `reg.is_registered(i)` returns the value stored in the given Registry object `reg` at index  $i$ .

A Registry object stores Boolean values, but the above description does not say how they are stored. The operations describe their effects on the Registry object, but the detail of the implementation of the operations is hidden. We can infer that a Registry object is a mutable data type, because its operations modify the data stored in the Registry. Apart from these 4 operations, there's nothing else we can do to a Registry object sensibly.

The Registry ADT is a data type, because its possible values are well-defined (the set of values is the set of all possible sequences of  $n$  Boolean values), and it has methods that permit specific

operations. It's an ADT because it encapsulates the data, and allows access only through the operations. The Registry ADT is not as flexible as a list, or as general as a Numpy array. In some applications, a Registry object gives exactly the right level of utility, neither too limited nor too flexible.

### 8.4.2 Statistics ADT

The purpose of this ADT is to compute the *mean* and *variance* of any numeric data set. The mean (also known as *average*) gives us an idea of what the typical data value is like. Another way to think about the mean is as a way to describe the data set's central tendency; data values are spread out around the mean value. The average is calculated by summing all the data values, and dividing by the number of data values. The variance tells us how much the data values vary from each other. If it's close to zero, then all the data values are very close together near the mean (on a number line); if the variance is large, then the data set has values that are quite spread out around the mean. Mathematically, the "variance" is defined as the average of all the squared differences between each data value and the mean value.

If the final sentence in the previous paragraph seems confusing to you, that's okay, because English is not the right language for it. There are two immediately better languages: mathematics, and a programming language (any one of them). We use math to describe a calculation precisely, and maybe to re-express it one way or another by symbolic derivations. We use a programming language to *implement* the calculation described by the mathematics.

The mean of  $n$  numbers is traditionally given the symbol  $\mu$ , and the formula for the mean is simply written as follows:

$$\mu = \frac{\sum_{i=1}^n a_i}{n}$$

Here, we're using the symbol  $a_i$  to represent each data value. In math, a subscript is used to *index* data values; in Python, we might write `a[i]`, implying that the data is stored in a list. The  $\sum$  symbol indicates a summation, i.e., a repeated addition of the term to its right. Below the  $\sum$  symbol is a variable  $i$  and an indication that this variable has a starting value, 1. The value written above the  $\sum$  symbol indicates the last value (inclusively) that  $i$  can take. So the formula says, in English, *add all the data values together, and divide by the number of data values*.

The variance of those  $n$  numbers is traditionally given the symbol  $\sigma^2$ . The formula for variance is as follows:

$$\sigma^2 = \frac{\sum_{i=1}^n (\mu - a_i)^2}{n}$$

Look carefully at the right-hand side of the formula. You'll see the symbol  $\mu$ , which is the mean, as described above. It's the central tendency for the data values. The expression  $\mu - a_i$  expresses the difference of any data value  $a_i$  from the mean; it can be positive or negative, depending on whether  $a_i$  is to the right or left of  $\mu$ . This quantity is squared, primarily because squaring makes differences positive; a secondary consideration is that a formula involving a squared term is more amenable to analysis than a formula involving an absolute value. So we can see that the variance adds up all the squared differences, and divides by the number of data points; in other words, it's an average squared difference. The square root of the variance is called the *standard deviation* (whose symbol is simply  $\sigma$ ).

These mathematical formulae can be expressed in Python as follows, assuming that the values are stored in a list named `data`.

```
mu = sum(data)/len(data)
sigsq = sum([(mu - x)**2 for x in data])/len(data)
```

Now that we understand the computations necessary to calculate mean and variance, we're going to hide the calculations behind an ADT that can be reused without memorizing formulae.

A Statistics ADT has four operations, which we describe in object-oriented terms:

- `Statistics()` creates a statistics record. The mean and variance are undefined until data is added. The new record is returned.
- `stat.add(val)` adds the data value `val` to the given Statistics object `stat`.
- `stat.mean()` returns the mean of the data previously added to the given Statistics object `stat`.
- `stat.var()` returns the variance of the data previously added to the given Statistics object `stat`.

The Statistics ADT is a data type, because its values are well-defined, and it has well-defined operations. It uses data values to determine the statistical summaries *mean* and *variance*. The operations describe how data can be added to the calculations, and how to ask for the statistical summaries. It's an ADT because the details of the implementation are hidden, and the only way to interact with the data.

### 8.4.3 Application example

In the following example, we use the Registry ADT to keep track of the outcomes of simulated die rolls. We're also using the Statistics ADT to find the average number of die rolls until we see the first repeated number.

```
1 import random as random
2 import Registry as R
3 import Statistics as S
4
5 sides = 6
6 stats = S.Statistics()
7 trials = 100000
8
9 for i in range(trials):
10     die_record = R.Registry(sides, False)
11     repeat_observed = False
12     roll_count = 0
13
14     while not repeat_observed:
15         roll = random.randrange(6)
16         roll_count += 1
17         if die_record.is_registered(roll):
18             repeat_observed = True
19         else:
20             die_record.set(roll)
21
```

```
22     stats.add(roll_count)
23
24     print("Rolled the", sides, "sided die", stats.mean(),
25           "times (on average) to observe a repeat event.")
26     print("The variance in the data is:", stats.var())
```

On lines 2 and 3, we import the modules that implement the ADTs; we haven't seen that implementation yet, but because the interface was described above, we can use the operations without knowing the implementation. On line 6, we create a Statistics object, and give it the name `stats`. On line 10, we create a Registry object; note that we create a new Registry object for every value of *i*. This registry has 6 Boolean values, all to them `False`. On line 17, we check to see if the registry object has already recorded one of the 6 possible rolls of the die. If not, on line 20, we use the registry object to record the first observation of the current roll. On line 22, we tell the Statistics object how many rolls were needed to observe a repeated roll in the current trial. Finally, on lines 24 and 25, we report the mean and variance of the number of rolls. For a 6-sided die, the average number of rolls to get the first repeat is about 3.8, and the variance is about 1.5.

The example shows that we can import modules that define the ADTs, create initial values, and use their operations to solve a computational task. It's no more complicated to use an ADT than any other Python object you've learned to use, like lists or dictionaries. But ADTs are different because programmers design them to solve a particular problem within some application. These ADTs are not part of the Python language.

#### 8.4.4 Discussion

We noted above that a Registry seems to be nothing more than a list of Boolean values. Since you are already familiar with lists, and Boolean values, you might well question the utility of the Registry ADT. You could replace each of the operations with Python code using lists. But think about the mistakes you could make with a list that you can't make with a Registry. A registry only stores Boolean values. It would be a mistake to store any other value in a registry, and the registry interface simply does not allow that kind of error to happen. The operations also go some distance towards documenting the code: the operations are descriptive of the purpose of the code, and using plain old Python list operations is less descriptive.

Any application using a Statistics record could also be rewritten to use a couple of Python variables, and a line or two of code to compute an average of a data set. It would be slightly more effort to replace the calculation of the variance; for one thing, you'd have to look up the formula, which you may not know. Then you'd have to check that your calculations are correct. Using the Statistics record is easier, because you don't have to look up the formula, and you don't have to test your implementation of the formula.

If you dig around the code (Appendix A), you can discover how the Statistics object works. You might have thought that it stores all the data in a list, and then calculates the mean and variance from that list. But the code shows that it's possible to implement a Statistics object so that you don't have to store all the data; you can calculate the mean and variance from 3 variables, without storing the whole data set inside the record. This can save a ton of space, especially for large data sets! This is the true value of ADTs: someone can make some clever design decisions and you just need to know how to use them. Furthermore, as long as the operation interfaces do not change, someone could improve the ADT by redesigning its organization, or any of the algorithms that implement the operations, without having to change any of the software that uses the ADT, because all of the

changes are hidden behind the interface.

The investment that you make in using and designing ADTs pays off more when your programs get bigger. If you don't practice using ADTs when programs are small, you will not have the experience to use them when you really need them!

As a final comment, remember that these two examples are just examples to demonstrate a concept. In other words, the examples we used in this chapter are details to demonstrate the principles of abstract data types. If you think ordinary Python code is more direct and simpler than these two examples, that's an acceptable point of view, because our example application is fairly small. There are more interesting ADTs to study in the coming chapters, and their benefits will be a little more convincing.

## 8.5 Summary

In this chapter, we defined the terms data type, data structure, and abstract data type. A data type is a set of values, and a set of operations that can be applied to these values. Computer Science needs the concept of data type, because computers do not (yet) have common sense enough to know when some operations don't make sense. A data structure is a special kind of data type, one that organizes multiple data values. Lists, strings, and dictionaries are examples of data structures.

An abstract data type hides organizational and implementational detail. This directly addresses several of design and implementation goals from Chapter 4. ADTs enhance correctness, because these operations can be tested thoroughly. Furthermore, since the data is encapsulated, and because the data cannot be changed except through well-tested operations, a programmer cannot corrupt the data accidentally. ADTs are adaptable, because once the interface to the operations has been defined, the implementation can be modified (usually to improve efficiency), but the changes are limited to the ADT definition. It will not be necessary to modify any application that uses the ADT, because the interface has not changed. ADTs are also highly re-usable.

We did not show how to create a data type in this chapter. Our ADTs will be defined by objects and classes. We will see how to do that in the next few chapters. It is important to emphasize that objects and classes are sufficient to define ADTs, but not necessary. Programmers can design ADTs without objects and classes; in fact, ADTs came first historically.



## 9 — Objects and Classes

### Learning Objectives

After studying this chapter, a student should be able to:

- Explain the differences between Procedural and Object Oriented Programming (OOP).
- Explain the difference between a object and a record.
- Explain the difference between a class and an object.
- Explain what attributes and methods are in terms of object oriented programming.
- Define simple classes, including data and methods, in Python.

### 9.1 Introduction

So far, the programs and scripts we've been writing have been based on a programming paradigm known as *procedural programming*. This approach to software emphasizes the technique of using functions (also known as procedures) to encapsulate algorithms. In the procedural programming paradigm, the data abstraction (dictionaries or structs) is enabled by functions that operate on the data, as we have seen in recent chapters. Since this approach is the only one we've presented so far, it would be natural to conclude that this is the only way to design software.

One of the benefits of this procedural approach is that it is minimalistic. Almost every programming language has these tools, and so you can apply the design techniques no matter which language you end up being required to use for any given project. Another benefit is pedagogical: we can teach procedural ADT design without introducing any new language features, and therefore can focus on concepts, not details. The procedural programming paradigm arose in the 1970s, when software engineering was in its infancy, and it was seen as a great advance over what had been the common practice before hand: completely unstructured and undesigned software (rather like what you'd see modern students do if we didn't teach things like functions and dictionaries).

In this chapter, we will introduce the object oriented programming paradigm. So we have been



making use of objects for quite a long time. What we will learn in this chapter is to define our own classes, and instantiate objects from our own classes. We'll see some important advantages over procedural programming, but we will pay a price: we have to learn new programming language details. Many if not all of the object oriented programming concepts we learn using Python will carry over to other object oriented languages, like Java or C++, but they may vary slightly in detail and appearance.

The primary concept of object oriented programming is that we will extend the idea of procedural abstraction, and use it for data abstraction, so that a single entity, which we call an *object*, encapsulates data and the procedures needed to operate on the data. This will allow us to design new data types that store data and its operations in the same entity. Object oriented programming languages, like Python, Java, and C++, have language features that make this idea as easy as possible, but also provide extra functionality that we would not be able to replicate very easily using the procedural programming paradigm. The object oriented programming paradigm is by far the most common approach to the design of very large software systems.

An object consists of data stored in the object ("encapsulated"), in much the same way that a record stores data; an object also encapsulates functions that operate on the data. The data in an object are called *attributes*, and the functions are called *methods*; this terminology encourages the idea that an object is an entity with a behaviour and some importance, as opposed to a container with stuff in it.

It's important to understand that an object is a value, in the same way that a dictionary is a value. In previous chapters, we've talked about objects somewhat informally, but we've tried to keep the use of the word consistent with the meaning we take here. An object is stored in Python's memory, and we can use variables to refer to objects. References to objects can be passed to functions as arguments, and can be returned from functions as return values. We can store references to objects in lists and dictionaries.

A *class* is a formal description of attributes and methods contained by an object. An object is created according to a class description; we say that the object is an *instance* or an *instantiation* of a class. When we say that an object is a value, we also have to say that its type is the class that it instantiates. Many objects can be instantiated from the same class, and the objects will have the same methods, but can have different attributes.

We already know a bunch of Python objects. For example, strings, lists, and dictionaries are all objects. A string is an instantiation of the string class, and has two primary attributes: the sequence of characters, and the length of the sequence. The string class has methods such as `split()` and `rstrip()` among many others. The list class has methods `append()` and `extend()`.

Recall that to call a method, we use the dot notation. For example, we can call `alist.append(9)` provided that `alist` is a name that refers to a list object. The dot-notation combines the name of the object and the name of the method; this is the way that Python knows which object should be affected by the method. To call the method, Python looks for the `append()` method associated with the object named `alist`. If it exists, it gets called; if it is not there, Python invokes an `AttributeError`, which tells us that the named object does not have the named method we tried to use.

## 9.2 Defining classes and instantiating objects

The Python keyword `class` begins a class definition. The name of the class follows, and the name is usually capitalized. In brackets following the class name, we write (`object`) followed by a friendly colon `:` (the box is only used here to draw attention to to colon). Every remaining part of the class

definition is indented relative to the keyword `class`, and the class definition looks somewhat like a function definition. Particularly, the `(object)` looks like a function's parameter, but it is not. We'll say more about this later. For now, accept it as a standard part of a class definition.

Every class definition must define an `__init__()` function; its purpose is to define what the object's attributes are, and to give them initial values. It's usually the first method defined inside a class definition. This method is responsible for setting up and initializing a new object. The use of the underscore characters in its name implies that this is not a method a programmer calls directly (usually); Python makes sure this method is used at the time an object is instantiated. The `__init__()` function is almost a normal function, but it is subject to the following conditions:

- It must have at least one parameter, conventionally named `self`, and it is always the first parameter.
- It must return `None`, which happens implicitly when there is no return statement.

The parameter `self` refers to the object itself, and it is not optional. If we think of an object as a thing, then its methods are actions, and the parameter `self` is the object's way of saying *me* or *I*. This method, along with any other method defined in the class, must have as its first argument the parameter `self`.

### 9.3 Simple examples

Consider the following very simple class definition:

```
1  class Hero(object):
2      def __init__(self, nn, pp):
3          '''
4              Purpose:
5                  Initialize the Hero object, self.
6              Pre-conditions:
7                  nn: a string, the Hero's name
8                  pp: a string, the Hero's power
9              Post-conditions:
10                 Initialize the Hero object, self.
11              Return:
12                 None
13          '''
14         self.name = nn
15         self.power = pp
16
17     def say_hello(self):
18         '''
19             Purpose:
20                 Introduces the Hero on the console.
21             Pre-conditions:
22                 None
23             Post-conditions:
24                 None
25             Return:
```

```
26         None
27     '''
28     print('Hello, evil-doers! My name is', self.name)
29     print('My super power is', self.power)
```

In this example, we define a class with two methods only. The class name is `Hero`. The first line of the definition uses the word `object` in parentheses; we'll say more about this line in the next chapter.

In the example, we can see that the `__init__()` function (lines 2-4) has three parameters. The first parameter is `self`, which is a necessary part of a method definition. On lines 3-4, the object's attributes `name` and `power` are being created and initialized using the dot-notation and the object `self`. This example basically says that a `Hero` object stores a name and a string as attributes.

Lines 17-29 define a method for the `Hero` class. This method doesn't do much except display some text to the console. As with all methods, the first parameter is `self`, and because the method doesn't do much, it has no need for more parameters. In this method, the attributes `name` and `power` are being accessed for their values. Notice how the dot-notation is used to access the object's attributes. Attributes are named variables stored inside the object. Methods are named functions stored inside the object.

Now let's see a simple application using this simple class:

```
1 import Hero as H
2
3 diana = H.Hero('Wonder Woman', 'super strength')
4 diana.say_hello()
5
6 bruce = H.Hero('Batman', 'martial arts')
7 bruce.say_hello()
```

Line 1 imports the class definition, and gives it an abbreviated name. Lines 3 and 6 show how to instantiate an object from a class definition. It's very easy to forget to use the name of the module, because very often, the module's name is the same as the name of the class defined in the module.

Notice that the class name is used, as if it were a function. Behind the scenes, Python creates a blank object (this is the object in the class definition statement above). Then Python calls `Hero's __init__()` function, giving three arguments: the new object, and the two arguments we see on line 3. When the class's name is used to create an object like this, we call it a *class constructor*.

This example created 2 new objects; both are instantiations of the `Hero` class. They come from the same class, and they both have attributes `name` and `power`, but the values for these attributes are different. Both `Hero` objects have the method `say_hello()`, but the effect of this method is slightly different, because the values of the attributes are different.

Notice that the method calls on lines 4 and 7 uses dot-notation, and that no arguments are given. In the class definition (above), the parameter `self` is required, but in the method call on line 6, there is no value given explicitly. This may seem confusing, but when we write `object.method(a1, a2)`, Python rearranges things behind the scenes, calling the method `method(object, a1, a2)`, inserting the object as the first argument. Thus, when the Python creates a frame for the method call, the parameter `self` refers to the object specified in the dot-notation. This is the Python designer's way of keeping Python similar to other object-oriented languages, which all use dot-notation.

Hopefully, this example is familiar to you, because of your previous experience using Python objects like lists, strings, and dictionaries (to name a few). The method calls on lines 4 and 7 use the same dot-notation structure as methods like `append()` or any other list method. The class constructor `Hero` used on lines 3 and 6, is analogous to using the function `list()` to create a new list object. In fact, functions `list()`, `dict()`, `int()`, `float()`, and `range()` are nothing more than class constructors for objects; these classes are provided by Python, but their use is consistent with the concepts outlined in this chapter. If you see a function whose name matches the name of a data type, then it's a class constructor.

## 9.4 Encapsulation and access control

A class contains attributes and methods; we call this *encapsulation*. By designing a class, we hide the data inside the object, allowing access only through the operations defined as methods. By restricting access to the data, we prevent the typical kinds of bugs that might occur when programmers have unrestricted access to data.

Among the concepts that object oriented programming brings to software design, the concept of *access control* is crucial. A class can indicate that some attributes should have restricted access, and others can have less restricted access. This flexibility allows programmers to protect what needs protecting, and to loosen the restrictions when doing so makes certain behaviours more convenient.

Python has a convention that governs access control for attributes. An attribute can be *public*, *protected*, or *private*.

- A *public* attribute can be accessed by any script, without restrictions. We use public attributes when allowing access makes other scripts and programs more convenient, and when having unrestricted access to the data will not put the data at too much risk. In Python, a public attribute starts with a normal lower-case letter; unlike private or protected attributes which start with at least one underscore character (`_`).
- A *private* attribute cannot be accessed by any script outside of the class definition. We use private attributes when limiting access to the data is essential to the correct behaviour of the object. Any attempt to access a private attribute outside the class definition will raise a run-time error. In Python, a private attribute starts with 2 underscore characters (`__`).
- A *protected* attribute is neither public nor private. It can be accessed outside the class definition, but only in related scripts designed specifically for the class, where insider knowledge will not put the data at risk. In Python, a protected attribute starts with exactly 1 underscore character (`_`).

For example, we've revised the `Hero` class above to make use of different levels of access control.

```
# define the class
class Hero(object):
    def __init__(self, name, power, sid):
        self.name = name
        self._power = power
        self.__secret = sid

# a simple script to use the class
bruce = Hero('Batman', 'martial arts', 'Bruce Wayne')
print(bruce.name)
```

```
print(bruce._power)
print(bruce.__secret)
```

In this example, the attribute name is public, the attribute `_power` is protected, and the new attribute, `__secret` is private. The main script below the class is able to display the Hero's name and super power, but a runtime error is raised as a result of trying to access the private attribute.

In some languages, for example, Java and C++, access to the attributes is very strictly enforced. In Python, the access control is more or less a suggestion. If you try hard enough, you can always access any attribute of any object of any class. That's part of the flexibility of Python programming. Still, use access control to indicate what level of access you think is essential, and trust that other Python programmers will have the wisdom to respect your design. Certainly, no one will have any sympathy for you if you ignore access control for Python objects and end up with a lot of buggy code because you were making changes to data you should not have been changing.

You can always provide methods as operations that returns the value of an attribute, or changes the value of an attribute, even if those attributes are protected or private. In object oriented programming, these kinds of operations are called *getters and setters*, for obvious reasons. By providing an operation, you can do things like check the validity of the operation, or hide details that might change with future designs. But you are also increasing the effort required to use your objects. Choosing whether to make attributes public, private, or protected is a matter of experience.

Finally, we can now explain why `__init__()` is named the way it is. It's a private method. Python's convention for access control for attributes extends also to methods. You can define public methods, that can be called by any application code, or private methods that can only be called within the class definition (but never in an application). Choosing whether to make methods public, private, or protected is a matter of experience.

## 9.5 Using classes to define ADTs

Perhaps you noticed that the previous sections on classes and objects sounded similar to some of the discussion in Chapter 8. One of the best uses for classes and objects is to define ADTs. In this section, we'll reveal the implementation of the Registry ADT, defined as a class.

```
1 class Registry(object):
2
3     def __init__(self, size, value):
4         """
5         Purpose:
6             Initialize a registry of a give size filled with the
7             given value
8         Pre-conditions:
9             size: number of elements in the registry
10            value: the default initial value for all elements
11         """
12         self.__reg = [value for i in range(size)]
```

The class definition starts with the name of the class, and the `__init__()` method. Note that the size of the Registry object, and the initial value, are parameters for the function. The `__init__()`

method creates one attribute, named `__reg`, which is a private attribute; its value is a normal Python list of the required size, containing the given value.

```
1  def set(self, i):
2      """
3      Purpose:
4          Set the registry element at i to True
5      Pre-conditions:
6          i: an index, in the correct range for reg
7      Post-conditions:
8          the registry value at i is set to True
9      Return:
10         (none)
11     """
12     self.__reg[i] = True
```

The `set()` method puts the Boolean value `True` into the list, using normal list indexing. Notice the combination of dot-notation and list indexing. The `reset()` is similar, so we will skip over it.

```
1  def is_registered(self, i):
2      """
3      Purpose:
4          Returns the value stored at registry element i
5      Pre-conditions:
6          i: an index, in the correct range for reg
7      Post-conditions:
8          (none)
9      Return:
10         the value stored at element i
11     """
12     return self.__reg[i]
```

The `is_registered()` method returns whatever value is stored in the registry.

This example encapsulates a list of fixed size within an object. The methods access the list in very limited ways. There is no way for an application to change the size of the list, and there's no way for a programmer to put any non-Boolean value into the list.

### 9.5.1 Perspective

You might think that this class is too simple to be useful, or you might think that using a normal Python list would be good enough for any application. That's a legitimate opinion for this particular ADT. But the purpose of this presentation is not to prove that a Registry is a prize-winning invention, but that it demonstrates the concept of an ADT, and how the ADT can be implemented using classes and objects. It should be very clear how the object encapsulates the data, and provides only limited access to the data through operations. We'll see many more ADTs, and their value won't be as dubious.

Earlier we said that the object oriented programming paradigm is by far the most common approach to the design of very large software systems. But that is not to say it is the only modern



programming paradigm. Functional programming languages, like Scala and Haskell, are showing very great potential, and modern languages like Python and Java are increasingly incorporating functional programming ideas. The point of this comment is that one should not expect software development to remain static. There is still much that computer science has to learn about software design, and many great ideas about software have not even been thought yet. Computer Science is still too young to become too rigid about programming!

## 10 — Stacks and Queues

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe what a queue is, and its basic operations.
- Describe the FIFO protocol.
- Employ a queue as a solution to a software design problem.
- Describe what a stack is, and its basic operations.
- Explain the LIFO protocol.
- Employ a stack as a solution to a software design problem.
- Describe common applications for stacks and queues.

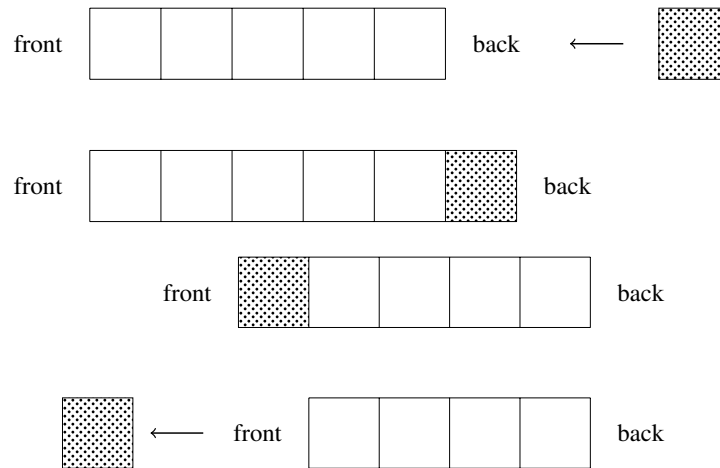
In Chapter 8, we defined an ADT called a *registry*. The operations allowed elements of the registry to be set or reset in any order. That kind of “any order” access is typical of Python lists, and of arrays in many languages. The term *random access* is often used to describe the idea that any element of a collection can be accessed quickly and efficiently. It may surprise you to learn that we don’t always need access to every element of a collection. It may also surprise you to learn that there are data structures that are useful but which do not allow random access.

In this chapter, we will present two kinds of collections that have limited access to their elements. The ADTs we look at here restrict the access to the data in very specific ways, and these are key to their definition. These ADTs are very common in computer science, and have a wide variety of applications. Functionally, these are very useful tools. Abstractly, they set the foundation for a useful way of designing all kinds of tools.

### 10.1 Queues: First-In First-Out (FIFO)

A queue is a data structure used to ensure that data is processed in the order it is received. Real-life examples of queues include the order that customers are served at stores and coffee shops, or the

Figure 10.1: Top: An *enqueue* operation adds a new data element (shaded) to the back of the queue. Bottom: A *dequeue* operation removes a data element (shaded) from the front of the queue.



order of cars waiting at a red-light at an intersection. In a coffee shop, customers enter the queue by standing behind the person who was previously last in the queue. The customers in the queue get served in the order that they arrived.

The organization of a queue is its key aspect: items in the queue are processed in the order that they arrived. This behaviour is commonly called *first-in first-out*, and when this behaviour is encoded by rules (such as a computer program), it's called a *protocol*. A queue has a *front* and a *back*. The data at the front has been in the queue the longest amount of time, and the data at the back has been in the queue the shortest.

**Definition 10.1** A *queue* (also called a *FIFO queue*) is a compound data structure in which the data values are ordered according to the first-in first-out protocol. A queue has a *front* and a *back* as attributes.

Data in a queue is accessed in limited ways. Every data element is added to the queue at the back. The common name for this operation is *enqueue*. A data element is removed only from the front of the queue. The common name for this operation is *dequeue*.

The primary operations provided for Queue ADTs are as follows:

- `Queue()` creates an empty queue.
- `queue.is_empty()` returns `True` if the given queue has no data in it.
- `queue.size()` returns the number of data values in the given queue.
- `queue.enqueue(value)` adds a given value to back of the given queue.
- `queue.dequeue()` removes the value at the front of the given queue, and returns it.
- `queue.peek()` returns the value at the front of the given queue, without removing it from the queue.

Below is a simple demonstration of the key operations mentioned. In it, a queue is created, prior to enqueueing the values 1-5. The value 1 is enqueued first, and 5 is enqueued last. All the data is dequeued one element at a time, and displayed. The data is pulled out of the queue in FIFO order:

1 is dequeued first because it was enqueued first; 5 is dequeued last because it was enqueued last. Note how the while loop uses the `is_empty()` operation to control how often values are dequeued.

```
import Queue as Q

data = Q.Queue()

for v in range(1,6):
    data.enqueue(v)

while not data.is_empty():
    d = data.dequeue()
    print(d)
```

## 10.2 Queue Application: Buffering communications

Queues are extremely useful, and play a significant role in many applications. We've already briefly mentioned situations in which you might have encountered the FIFO protocol in real life: the coffee shop, for example.

Queues are used as *buffers* to smooth communications. The most obvious example is the transmission of audio or video data, from websites and services such as YouTube, NETFLIX, or SoundCloud.<sup>1</sup> Data transmission rates are highly variable on the internet, and computers themselves can be faster or slower. Queues can be used to moderate the effects of variable rates.

A general model for this application has a producer service sending data to a consumer. If the producer and consumer communicate directly, sending and receiving one data element at a time problems can arise. If the producer sends data faster than the consumer can make use of it, the producer has to wait for the consumer. If the consumer can make use of the data faster than the producer can send it, the consumer has to wait for the producer. If the transmission rate varies, then both the producer and consumer can be waiting at different times.

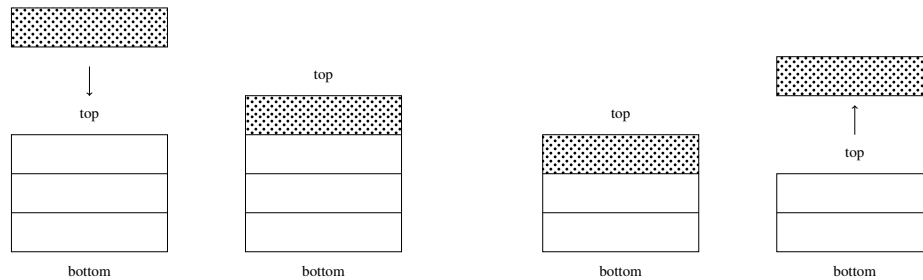
To smooth out the amount of time either side spends waiting, a buffer is used (either at the consumer side, or both sides). The producer can send data to the consumer's buffer as quickly as possible at first, and the buffer stores it temporarily. Once there is enough data in the buffer, as determined by the application, the consumer can start using the data in the buffer without waiting for it. Then the producer can send data to the buffer in bursts, or a steady rate. As long as there is data in the buffer, the consumer doesn't have to wait.

A buffer is essentially a queue. The data in the buffer is processed in the order it was sent, which is the order the consumer is expecting it. You are probably aware of buffering on video-on-demand services, such as YouTube or NETFLIX. But buffers are an essential component underlying the software that enables all internet communications, not just video. Buffers are used in modern graphics hardware, queuing up video frames to be displayed on your monitor. Even the console I/O that we do in Python is smoothed by input and output buffers. Disk drives, too.

---

<sup>1</sup>These names are trademarks owned by their respective owners.

Figure 10.2: Left: An *push* operation adds a new data element (shaded) to the top of the stack. Right: A *pop* operation removes a data element (shaded) from the top of the queue.



### 10.3 Stacks: Last-In First-Out (LIFO)

A stack is a data structure used to ensure that the most recent data is processed first. In real life, we tend to accumulate stacks of things like plates, books, exam papers. Other examples of stacks include the layering of sedimentary rocks: the top of the layer is the youngest layer, and the bottom layer is the oldest. These real life examples are over-simplified, since you could lift a stack of books to access one near the bottom, and sedimentary rocks do experience phenomena such as folding and subduction.

We use the expression *last-in first-out* to describe a stack's protocol. It implies that the most recent data is always processed first. Organizationally, a stack has a top, and a bottom. The data at the bottom of the stack has been in the stack for the longest amount of time, and the data at the top of the stack has been there the shortest. Operationally, data is added and removed only at the top of the stack, so the *top* is the only attribute we need to define a stack. Adding to the top of a stack is commonly called a *push*, and removing an item from the top of a stack is called a *pop*.

**Definition 10.2** A *stack* (also called a *pushdown stack*) is a compound data structure in which the data values are ordered according to the last-in first-out protocol. A stack has a single attribute, called a *top*.

The primary operations provided for Stack ADTs are as follows:

- `Stack()` creates an empty stack
- `stack.is_empty()` returns `True` if the given stack has no data in it
- `stack.size()` returns the number of data values in the given stack
- `stack.push(value)` adds a given value to top of the given stack
- `stack.pop()` removes the value at the top of the given stack, and returns it
- `stack.peek()` returns the value at the top of the given stack, without removing it from the stack.

Below is a simple demonstration of the key operations mentioned. In it, a stack is created, prior to pushing the values 1-5. The value 1 is pushed first, and 5 is pushed last. All the data is popped one element at a time, and displayed. The data is pulled out of the queue in LIFO order: 5 is dequeued first because it was pushed most recently; 1 is popped last because it was pushed first. Note how the while loop uses the `is_empty()` operation to control how often values are popped.

```
import Stack as S

data = S.Stack()

for v in range(1,6):
    data.push(v)

while not data.is_empty():
    d = data.pop()
    print(d)
```

## 10.4 Stack Application: Backtracking

You’ve almost certainly encountered stacks of plates and books and other physical items that can be literally stacked.

You have also encountered stacks in your web-browser. When you visit a page, the URL for that page gets pushed on a stack, which we will call the *back-stack*. The *back* button pops the back-stack to re-display the page you the page you were most recently looking at. Most web-browsers also have a *forward* button, which also uses a stack, which we could call the *forward-stack*. When you go back, your current page’s URL is pushed onto the forward-stack, and one URL is popped from the back-stack and displayed to you. Likewise, when you use the forward button, the URL of the current page is pushed onto the back-stack, and a URL is popped from the forward-stack and displayed. That’s how you can go back and forth through your recent browser history without losing anything.

Stacks are also useful in applications like word processors and other kinds of content creation applications. These applications typically have an *undo* button, and a *re-do* button. These are very similar to the browser buttons back and forward. Every time you modify a document, data describing the modification<sup>2</sup> is pushed on an *undo-stack*. When you use the undo button, the undo-stack is popped, and the application uses the information to undo the action. The modification data can be pushed on a second stack, which is controlled by the re-do action button. The only tricky thing with undo and re-do is encoding the modifications in a way that allows the application to store the changes on the stack, so that they can be performed.

This concept is so common that it has a name: back-tracking, i.e., returning to a previous situation so that you can try something else. We’ve seen how stacks enable a person to back-track to a previous state in the above examples. The concept of back-tracking is also applicable to the design of algorithms, as we will see in Chapter 26. In a nut-shell, if a program uses back-tracking, it has been programmed to explore a *space* of possible solutions, looking for the right solution that might be in the space. If the algorithm is faced with a number of options, and there’s no known way to tell which option is best, back-tracking allows the algorithm to try them all, without getting lost. Back-tracking gets used as an algorithm strategy sometimes when no better algorithm has been discovered, kind of a last resort, or as we will suggest, a first stab at solving tricky problems. A stack is a crucial part of all back-tracking algorithms.

We will see other applications of queues and stacks in Chapter 11.

---

<sup>2</sup>It might be hard to imagine what form this data might take, but whatever it is, it goes on a stack!





## 11 — Applications of Stacks and Queues

### Learning Objectives

After studying this chapter, a student should be able to:

- Write programs that use stacks and queues.
- Explain the key features of a simple queueing simulation.
- Explain how the queueing simulation makes use of the Queue ADT.
- Explain the key features of bracket checking, using a stack.
- Explain how bracket checking makes use of the Stack ADT.
- Explain the key features of evaluating post-fix expressions, using a stack.
- Explain how evaluating post-fix expressions makes use of the Stack ADT.

### 11.1 Introduction

In this chapter we demonstrate the use of the Queue and Stack ADTs in algorithms that highlight the importance of their key aspects. The algorithms are demonstrative of techniques in computer science as well.

### 11.2 Queueing simulation

A recurring theme in computer science is the use of software to answer questions that would be difficult or impossible to answer using analytical techniques. In the case of the simulation in this chapter, analytical techniques are able to demonstrate several of the properties we describe, but studying the algorithm is useful for the technique it demonstrates.

We will be simulating the situation in which customers arrive (say at a coffee shop) and wait in a queue for service. One of the main assumptions in our simulation is that there will be one queue in which all the customers will wait. This keeps things simple, as we won't have to simulate the decision process of customers faced with multiple queues. The customers arrive, wait in the queue

for some amount of time, get served, and then leave the queue. We're interested in knowing how long the customers wait in the queue on average, and how long the queue gets. To answer these questions, we will demonstrate a program that literally simulates arrivals, waiting, service, and finally departure of customers.

It should be pointed out that this kind of simulation has a long history in computer science. Queueing theory, as it is called, starts with the simple model we show here, and extends to models that are much more sophisticated. Understanding the mechanics of queueing is vital to their use in the real world (as in the coffee shop example, but extending to supply chain management for manufacturing, and other applications), as well as the software services in the modern world. Your computer is full of applications of queues that affect your daily experience, from the queues that are used to buffer the multi-tasking operating system on your computer, to the queues that buffer streaming data from games and multi-media services like YouTube and NETFLIX.

Two of the key aspects of the simulation are the customer *arrival* times, and the customer *service* times. In simulations of this kind, it is standard to assume that the customers arrive randomly, but with an assumed arrival rate. While we have used random numbers generated by Python's `random` module previously, we have always done so by assuming a number is drawn from a given range (e.g.,  $[0, 1]$ ) with equal probability. We could use the same assumption here, but it's standard to assume that arrival times and service times are not equally likely.

The standard assumption for this simple model is that arrival times are governed by something called a *Poisson process*. The Poisson process puts no upper limit on the time between arrivals; a very long period between arrivals is possible, but very unlikely. The Poisson process describing the arrival times is governed by the *arrival rate*, in arrivals per unit time. We won't be doing a real-time simulation, so the time unit can be arbitrary; it's intuitive in coffee shop simulations to imagine arrivals per minute, but in a NETFLIX simulation, data arrivals might be measured in data packet arrivals per millisecond. When we think in terms of arrivals per unit time, we are using an abstraction.

In the Poisson process describing arrivals, the time between arrivals is random, but described by an *exponential distribution*, in which the time between arrivals is more likely to be shorter than longer. The standard model also assumes that the service time is described by an exponential distribution, governed by an assumed *service rate*. As with arrival times, services times are random, and more likely to be shorter than longer.

In Python, we can ask for a random number from an exponential distribution using the `random` module, using the function `expovariate`, as in the following function:

```
def sample_time(x):  
    '''  
    Return a random sample time until an event.  
    In the long run, the events will have a rate of x  
    events per unit time.  
  
    Preconditions:  
        x: the desired arrival rate, per unit time  
    Post-conditions:  
        none  
    Return:  
        a random sample time that obeys the rate x  
    '''
```

```
'''
    return rand.expovariate(x)
```

We call this function with either the service rate or the arrival rate, and we get a random number generated for us that obeys the assumptions about arrivals and service described above. Observe that we're using procedural abstraction here to hide the details about how the samples are obtained; it's a lot of work for a one-line function body, but if we ever decide to change the assumptions about how sample times are generated, we can make those changes in one function, and the rest of the script won't change. That's adaptability!

Let's turn to the mechanics of the queueing model. There will be a single queue, and customers will arrive, enter the queue, and wait for service. When the customer reaches the front of the queue, they are served, which takes some time, before they leave the queue. In our simulation, we will use a FIFO Queue to store arrival times.

```
# customers wait in a service queue
service_queue = Q.Queue()

# we're interested in the average time in the queue
waiting = S.Statistics()

# keep track of when a new customer will arrive
nextArrival = sample_time(arrival_rate)

# keep track of when a service event is complete
nextService = nextArrival + sample_time(service_rate)

# simulate the arrival-service process
while nextService < sim_length:

    # handle all arrivals that occur before service is complete
    while nextArrival < nextService:
        service_queue.enqueue(nextArrival)
        nextArrival = nextArrival + sample_time(arrival_rate)

    # service time is over; serve first customer in the queue
```

When a customer leaves the queue, we can check the current time, and calculate the amount of time the customer spent in the queue by subtraction. We will use the Statistics ADT introduced in Chapter 8 to manage the calculation of the average wait time.

```
# we're interested in the average time in the queue
waiting = S.Statistics()
```

Time is a key aspect of our simulation, and we normally think of time as progressing in uniform steps, like a clock ticking in seconds or milliseconds. But in this kind of simulation, which is called an *event-based simulation*, time advances in leaps to whichever event is determined to be next. To

work through a simulation with arrivals and service, we keep track of two times: the time of the next arrival event, and the time of the next service event.

- The next arrival is an event that places a new arrival in the queue.
- The next service is an event that removes a customer out of the queue; it indicates the time that the service for the next customer finishes.

To start the simulation, the next service event is scheduled to occur after the first arrival.

```
# keep track of when a new customer will arrive
nextArrival = sample_time(arrival_rate)

# keep track of when a service event is complete
nextService = nextArrival + sample_time(service_rate)
```

Our simulation will run for an amount of simulated time, specified by the variable `sim_length`. The simulation will repeatedly advance to whichever event is next. Sometimes the next event is an arrival, so we'll handle the arrival by putting the arrival in the queue.

However, when the service event is next, we'll calculate how long the customer at the front of the queue was waiting, and then schedule a new arrival event. If the queue happens to be empty, then the next service event happens after the arrival of the next customer. On the other hand, if the queue is not empty, the next service event occurs some time after the most recent service event.

```
# simulate the arrival-service process
while nextService < sim_length:

    # handle all arrivals that occur before service is complete
    while nextArrival < nextService:
        service_queue.enqueue(nextArrival)
        nextArrival = nextArrival + sample_time(arrival_rate)

    # service time is over; serve first customer in the queue
    this_arrival = service_queue.dequeue()

    # how long was the customer waiting?
    waited = nextService - this_arrival
    waiting.add(waited)

    # determine when service to the next customer will end
    if service_queue.is_empty():
        nextService = nextArrival + sample_time(service_rate)
    else:
        nextService = nextService + sample_time(service_rate)

print('Average wait time:', waiting.mean())
```

Notice the use of meaningful variables names, and the use of comments. The script could be written without the use of the Queue and Statistics ADTs, but their use makes the code clearer and more

concise. It takes effort to produce code like this, even after you know how the algorithm works. The full listing of the queueing simulation can be found in Appendix B.

### 11.2.1 Discussion

The queueing simulation presented here is also known as an M/M/1 queue, which is studied in applied probability courses, as one of the elementary models for the use of probability theory. Since this is not a probability course, we'll say only a few words about what such studies tell us.

First of all, our intuitions tell us that if the service rate is slower than the arrival rate, the queue will get very big. Formal analysis confirms this. If the arrival rate is slower than the service time, then number of customers waiting in the queue depends on the ratio between rates. To clarify this, formal analysis typically uses  $\lambda$  to represent the arrival rate, and  $\mu$  to represent the service rate. If  $\mu \leq \lambda$ , the queue can grow arbitrarily large, because customers will arrive faster than they are served. But if  $\mu \geq \lambda$ , there may be a few bursts of high traffic, but the service can keep the queue moving: in this case, the average number of customers in the queue is  $\lambda/(\mu - \lambda)$ . For example, if the arrival rate,  $\lambda = 1$ , and the service rate is  $\mu = 2$ , the queue will have on average 1 customer in it. Formal analysis also tells us that the average waiting time for customers is  $1/(\mu - \lambda)$ . So if the arrival rate,  $\lambda = 1$ , and the service rate is  $\mu = 2$ , the average waiting time is also 1 time unit. We can set our simulation to verify these formal results empirically.

## 11.3 Bracket checking

At some point in your study of programming you may have wondered how Python can convert the text you type in an editor and turn it into a thing that the computer does. If you haven't wondered that, perhaps you can take a few seconds now. We'll wait for you.

The process of converting Python code (or any programming language for that matter) into machine code is a complicated process, but we can gain an appreciation for that process by looking at an algorithm for bracket checking. When we write Python expressions that use arithmetic, the expression needs to obey certain rules about the syntax of the expressions. The process of checking syntax is called *syntactic analysis*, or sometimes *parsing*. If the syntax is not correct, Python cannot make sense of the expression, and reports a syntax error.

Bracket checking can be considered a very limited form of syntactic analysis. It demonstrates the use of a stack very effectively; stacks are key to solving the larger problem of more general syntactic analysis. We want our program to recognize the kinds of bracketing we'd expect to see in our Python expressions. For example, we might write an assignment statement using an expression as follows:

```
result = ((3 + 4) * (5 + 6))
```

This example should be accepted as syntactically correct (even though the outer-most brackets are not actually necessary, they are still acceptable). In contrast, the following example should produce a syntax error:

```
result = ((3 + 4 * (5 + 6))
```

because there are more left-brackets than right-brackets.

To simplify the presentation, we will eliminate the arithmetic, and just keep the brackets. This is enough to demonstrate the use of stacks in an understandable, if a bit simplified, application. Without arithmetic, our expressions look a bit weird, and they'll be strings:

```
example1 = '(()())'  
example2 = '((()))'  
example3 = '(()))'
```

We want a program to analyze any string consisting simply of brackets, and indicate if every left-bracket has a matching right-bracket. The key to this task is to realize that a right-bracket matches the most recently seen left-bracket. The phrase “most recently seen” is a strong indication that we need to use a stack. And for reasons that we will discuss later, we will also use a queue. We will use a boolean variable called `unmatched_close` to indicate whether we’ve found a right-bracket without a matching left-bracket, as in the third example above.

```
# create the initial empty containers  
chars = Q.Queue()  
brackets = S.Stack()  
unmatched_close = False
```

The first thing we do is put all the brackets into a queue. This simplifies the matter of stepping through the string, because we don’t need to keep track of the length of the string, or where we are in the string.

```
# put all the characters in the Queue  
for c in example:  
    chars.enqueue(c)
```

The bracket checking process works by examining the queue of characters, one at a time. If the current character is a left-bracket, it is pushed onto the stack. If the current character is a right-bracket, it either has a matching left-bracket on the top of the stack, or, if the stack is empty, no matching left-bracket. The program below also ignores every character that is not a bracket; this allows us to process examples with only brackets, or examples of arithmetic expressions.

```
# brackets match if and only if every '(' has  
# a corresponding ')'  
while not chars.is_empty() and not unmatched_close:  
    c = chars.dequeue()  
    if c == '(':  
        brackets.push(c)  
    elif c == ')' and not brackets.is_empty():  
        brackets.pop()  
    elif c == ')' and brackets.is_empty():  
        unmatched_close = True  
    else:  
        pass
```

There are two reasons we could have left the while loop above. First, we may have detected an unmatched right-bracket, in which case we can report that fact. Otherwise, the reason we exited the loop was because all the characters in the queue were processed. In that case, the stack could be empty, which means every pair of brackets matched perfectly. However, if the stack is not empty, it means there were more left-brackets than right-brackets, and we should report it.



```
# check how the analysis turned out
if unmatched_close:
    print("Found a ')' with no matching '('")
elif not brackets.is_empty():
    print("At least one '(' without a matching ')")
else:
    print('Brackets matched')
```

While simple, this process of stacking up characters is essential to the analysis of text, whether for applications to convert Python programs to machine code, or for processing natural languages (like English) to extract meaning. The programs for those are more complicated, but the bracket-checking program is demonstrative of the principles involved. Text, whether program or prose, is a sequential representation, and keeping parts of it in a stack allows a program to draw local associations, like matching brackets, or noun-verb agreement.

The full listing of the bracket matching can be found in [Appendix B](#).

## 11.4 Post-fix Arithmetic evaluator

Now that we've demonstrated how to use a stack to check for a limited form of syntactic analysis, we will give a demonstration of how to take syntactically correct expression and come up with the correct value. This is called *evaluation*, and a program that does this is called an *evaluator*.

In order to keep the evaluator simple, and to focus on the use of stacks, we have to restrict the kinds of expressions we will evaluate. Normally we write arithmetic expressions by putting the operators *between* the operands, as in  $3 + 4$ . When we have more complicated expressions with several operators in it, we typically use a set of rules called “order of operations” to describe how an expression should be evaluated. For example, the order of operations tells us that multiplication should occur before addition in expressions such as  $3 + 4 \times 5$ . We can dispense with the order of operations if we insist on using brackets for every operator, as in the example  $((3 + 4) \times 5)$ .

There is another way that we can dispense with memorizing order of operation rules, which requires writing expressions differently. Instead of writing the operator between the operands, we will write the operator *after* both operands. For example, instead of writing  $(3 + 4)$ , we would write  $(3\ 4\ +)$ . As another example, instead of writing  $((3 + 4) \times 5)$ , we would write  $((3\ 4\ +)\ 5\ \times)$ . As a final example, instead of writing  $((3 + 4) \times (5 + 6))$ , we would write  $((3\ 4\ +)\ (5\ 6\ +)\ \times)$ . If you take a moment to reflect (we'll wait), you may realize that we actually don't need brackets if we put the operator after its arguments. For example, there is no ambiguity in the expression:  $3\ 4\ +\ 5\ 6\ +\ \times$ . This arrangement of operators and operands is called *post-fix* or *reverse Polish* notation.

To evaluate the expression, we scan it left-to-right, and when we find an operator, we perform the operation on the two immediately preceeding operands. The phrase “immediately preceeding” should be a strong indicator of the use of a stack. The key to correct evaluation is that we replace the

expression with its value, and continue moving to the right. Here is an example of the process:

$$\begin{array}{r} 3\ 4\ +\ 5\ 6\ +\ \times \\ 7\ 5\ 6\ +\ \times \\ 7\ 11\ \times \\ 77 \end{array}$$

In the example,  $3\ 4\ +$  is replaced by 7 (as indicated by underlining), and then  $5\ 6\ +$  by 11. When we reach the  $\times$ , there are two values waiting to be used, and we get the final result, 77.

It may be a little confusing to see the operator last, but once you get used to the idea, it's no more trouble than the traditional organization of operations. In the early days of computer technology, before smartphones and highly portable computers, specialized devices called electronic calculators were used to assist in calculations done on paper. The earliest calculators did not have brackets to help users enter their calculations, because the earliest calculators did not have a stack.<sup>1</sup> The first calculators with a stack were used by scientists and engineers, who mentally translated "normal" expressions into post-fix, so they could use their advanced calculators, much to the annoyance of the scientists and engineers who insisted that slide-rules are the only calculating device a so-called "real" scientist should use.

To evaluate post-fix expressions, we start with a string that contains the expression. We assume that at least once space separates numbers and operators from each other. For example:

```
example = "3 4 + 5 *"
```

We will implement the evaluator as a function, with the following interface:

```
def evaluate(expr):
    """
    Evaluate a postfix expression.
    Pre-conditions:
        expr: a string containing a postfix expression
    Post-Conditions:
        none
    Return:
        the value of the expression
    """
```

As in the bracket checking program, we will use a queue to store the expression, and a stack to assist in the evaluation.

```
# create the initial empty data structures
expression = Q.Queue()
evaluation = S.Stack()
```

The first order of business is to split the string at the spaces, and put the pieces into the queue. Again, using a queue allows us to process the pieces without having to keep track of where we are in the string, or in a list of pieces.

<sup>1</sup>The effects of this deficiency can still be observed in the default Calculator app in Microsoft Windows. Go ahead, open it up and enter  $3 + 4 \times 5$  in exactly that order, and be appalled.

```
expr_list = expr.split()
# put all the items in the Queue
for c in expr_list:
    expression.enqueue(c)
```

The function walks through the expression left-to-right, using the expression queue. We dequeue a piece of the string, and look at it.

- If the piece looks like an operator, we pop the evaluation stack twice, perform the operation, and put the result back on the stack.
- If the piece doesn't look like any operator, convert it to a number, and push it onto the stack.

```
while not expression.is_empty():
    c = expression.dequeue()

    if c == '*':
        v1 = evaluation.pop()
        v2 = evaluation.pop()
        evaluation.push(v1*v2)
    elif c == '/':
        v1 = evaluation.pop()
        v2 = evaluation.pop()
        evaluation.push(v2/v1)
    elif c == '+':
        v1 = evaluation.pop()
        v2 = evaluation.pop()
        evaluation.push(v1+v2)
    elif c == '-':
        v1 = evaluation.pop()
        v2 = evaluation.pop()
        evaluation.push(v2-v1)
    else:
        evaluation.push(float(c))

return evaluation.pop()
```

The full listing of the post-fix evaluator can be found in [Appendix B](#).

The above code correctly evaluates proper post-fix expressions, and makes appropriate use of computational resources. It makes use of procedural abstraction and ADTs, so it takes advantage of reusable components, and itself is fairly reusable. But it's not very good code, because it is very brittle (i.e., not robust). It assumes that the input string is a perfectly correct post-fix expression. The function will cause a run-time error if the string is anything but perfect. The code is also very repetitive, with very small differences between the four cases that handle the different operators. It's good enough to demonstrate the use of stacks and queues, and to teach the concept of post-fix arithmetic, but we can do better. And we will, in [Chapter 14](#)!



Why use Queues and Stacks when we have lists?

#### Adapting Python Lists to Implement Queues

Operation: Initialization  
Operations: `size()` and `is_empty()`  
Operation: `enqueue()`  
Operation: `dequeue()`  
Operation: `peek()`

#### Adapting Python Lists to Implement Stacks

Operation: `create()`  
Operations: `size()` and `is_empty()`  
Operation: `push()`  
Operation: `pop()`  
Operation: `peek()`  
Summary

## 12 — List-based Stacks and Queues

### Learning Objectives

After studying this chapter, a student should be able to:

- Employ data abstraction to implement a queue.
- Employ data abstraction to implement a stack.
- Explain the concept of an adapter in the context of abstract data types.
- Describe the list-based implementation of stacks and queues.
- Explain some of the reasons that wrapping an ADT around a list can improve correctness and adaptability.

### 12.1 Why use Queues and Stacks when we have lists?

Python lists are very versatile, and have a lot of uses. They have all the operations we identified as essential for FIFO Queues and LIFO Stacks in the previous chapter, and lots more besides. While such a useful tool is good to have, it is also useful to apply a certain amount of discipline to complex programming tasks, to improve correctness, and increase adaptability and robustness.

In this chapter, we will see how we can implement the Queue and Stack ADTs using Python lists. Essentially, we will be converting a Queue (or Stack) operation into a list method. This is such a common technique in software development that it has its own name. When we take a data structure or ADT *A* and provide operations so that it looks and behaves like ADT *B*, we have built an *adapter*. It is useful because sometimes a programming task requires a data structure that is not available directly in a language or a module. One way to provide the tool is to adapt something similar to provide the required functionality. This chapter gives the very simplest example of that kind of adaptation. The approach generalizes well beyond what we will demonstrate in this chapter.

Another reason to prefer using an ADT adaptation over Python's versatile lists is the way that it enhances correctness. If your program needs LIFO behaviour, for example, using a Queue documents

this fact in the program, and also enforces the LIFO protocol. Should you need to fix, extend or reuse the program, the fact that you used a Queue will prevent you from forgetting a crucial aspect of the algorithm. If you did not use a Queue, but used a Python list directly, it would be much easier to subvert the purpose of the list, and violate the FIFO protocol. That kind of subversion may allow you to make a quick bug fix in the short term, but might introduce a very subtle bug that would be difficult to find in the long term.

The last reason for this chapter is to set the scene for an important lesson in the flexibility of ADTs. An ADT is defined by the interface of its operations, not by the implementation of the operations. In this chapter, we will implement the Queue and Stack ADTs using Python lists. In a later chapter, we will implement the same ADTs using different data structures. The lesson will be that once an ADT's interface is defined, several implementations may be possible. Any application that uses the ADT's interface can use any one of the implementations. So you can prototype your application using a very simple ADT implementation, and then improve it by replacing the simple ADT with a more efficient one. This is the very essence of adaptability and robustness.

## 12.2 Adapting Python Lists to Implement Queues

The essential nature of a Queue is the FIFO protocol. Python provides two list methods that we will use to implement the Queue ADT.

**append()** This Python list method takes an item and adds it to the end of the current list. We use this method to implement the enqueue operation.

**pop()** This Python list method takes an index and removes the item at that index from the list. If we call `pop(0)` the first element in the list is removed and returned. We use this to implement the dequeue operation.

The list-based implementation simply wraps an ADT around the Python list. The object we are calling a queue is nothing but a list. Anyone who understood this fact could use any of Python list methods on the queue we have implemented. The front of our queue will be at list index 0, and the back of the queue will be at list index  $n - 1$ , where  $n$  is the length of the list. The following sections describe the implementation in a little more detail. The complete implementation can be found in Appendix C.

### 12.2.1 Operation: Initialization

The `create()` function returns a new empty queue. In this implementation, a queue is nothing but a list, hidden behind the Queue ADT's interface. Python `list()` function creates a new, empty list, and we are simply using that function.

```
1 class Queue(object):
2
3     def __init__(self):
4         """
5         Purpose
6             creates an empty queue
7         Return
8             an empty queue
9         """
```



```
10     self.data = list()
```

### 12.2.2 Operations: size() and is\_empty()

The `is_empty()` operation returns True if the queue is empty. The `size()` operation returns the number of elements stored in the queue. Both implementations use Python's `len()` function directly, because in this implementation, a queue is nothing but a list.

```
1     def is_empty(self):
2         """
3         Purpose
4             checks if the given queue has no data in it
5         Return:
6             True if the queue has no data, or false otherwise
7         """
8         return len(self.data) == 0
9
10
11     def size(self):
12         """
13         Purpose
14             returns the number of data values in the given queue
15         Return:
16             The number of data values in the queue
17         """
18         return len(self.data)
```

### 12.2.3 Operation: enqueue()

The `enqueue()` operation adds a new data element to the back of the queue. To implement this operation, we used Python's list method `append()`, which adds a value to the end of the list.

```
1     def enqueue(self, value):
2         """
3         Purpose
4             adds the given data value to the given queue
5         Pre-conditions:
6             queue: a queue created by create()
7             value: data to be added
8         Post-condition:
9             the value is added to the queue
10        Return:
11            (none)
12        """
13        self.data.append(value)
```

### 12.2.4 Operation: dequeue()

The `dequeue()` operation removes a data element from the front of the queue. To implement this operation, we used Python's list method `pop()`, which removes and returns a value from the list. Calling `pop(0)` removes and returns the first value in the list.

```
1  def dequeue(self):
2      """
3      Purpose
4          removes and returns a data value from the given queue
5      Pre-conditions:
6          queue: a queue created by create()
7      Post-condition:
8          the first value is removed from the queue
9      Return:
10         the first value in the queue
11     """
12     return self.data.pop(0)
```

### 12.2.5 Operation: peek()

The `peek()` operation returns a reference to the data element at the front of the queue. To implement this operation, we use Python's list indexing. This method returns the value, without removing it.

```
1  def peek(self):
2      """
3      Purpose
4          returns the value from the front of given queue
5          without removing it
6      Pre-conditions:
7          queue: a queue created by create()
8      Post-condition:
9          None
10     Return:
11         the value at the front of the queue
12     """
13     return self.data[0]
```

## 12.3 Adapting Python Lists to Implement Stacks

The essential nature of a Stack is the LIFO protocol. Python provides two list methods that we will use to implement the Queue ADT.

**append()** This Python list method takes an item and adds it to the end of the current list. We use this method to implement the push operation.

**pop()** This Python list method takes an index and removes the item at that index from the list. If we call `pop()` *without any argument*, the *last* element in the list is removed and returned. We use this to implement the pop operation.

As with the queue implementation, our stack implementation simply wraps an ADT around the Python list. The object we are calling a queue is still just a list. Anyone who understood this fact could use any of Python list methods on the queue we have implemented. The bottom of our stack will be at list index 0, and the top of the stack will be at list index  $n - 1$ , where  $n$  is the length of the list. We could have implemented the stack so that the top of the stack is at index 0. The implementation would be slightly different, but from outside the ADT, no one would be able to tell which end we are using. The following sections describe the implementation in a little more detail. The complete implementation can be found in Appendix C.

### 12.3.1 Operation: create()

The `create()` function returns a new empty queue. In this implementation, a queue is nothing but a list. Python `list()` function creates a new, empty list, and we are simply using that function.

```
1 # Implementation:
2 # This implementation uses Python lists directly.
3
4 class Stack(object):
5     def __init__(self):
6         """
7         Purpose
8             initializes an empty Stack object
9         """
10        self.data = list()
```

### 12.3.2 Operations: size() and is\_empty()

The `is_empty()` operation returns True if the queue is empty. The `size()` operation returns the number of elements stored in the queue. Both implementations use Python's `len()` function directly, because in this implementation, a queue is nothing but a list. This implementation is identical to the queue operations of the same name. Using the same name is actually a helpful idea. It's one fewer detail to memorize, and it hints at the idea that there is a core similarity between queues and stacks.

```
1
2     def is_empty(self):
3         """
4         Purpose
5             checks if the stack has no data in it
6         Return:
7             True if the stack has no data, or false otherwise
8         """
9         return len(self.data) == 0
10
11
12     def size(self):
13         """
14         Purpose
15             returns the number of data values in the stack
```

```

16         Return:
17             The number of data values in the queue
18         """

```

### 12.3.3 Operation: push()

The push() operation adds a new data element to the top of the queue. To implement this operation, we used Python's list method append(), which adds a value to the end of the list.

```

1     def push(self, value):
2         """
3         Purpose
4             adds the given data value to the stack
5         Pre-conditions:
6             value: data to be added
7         Post-condition:
8             the value is added to the stack
9         Return:
10            (none)
11         """
12         self.data.append(value)

```

### 12.3.4 Operation: pop()

The pop() operation removes a data element from the top of the stack. To implement this operation, we used Python's list method pop(), which removes and returns a value from the list. Calling pop() without an argument removes and returns the last value in the list.

```

1     def pop(self):
2         """
3         Purpose
4             removes and returns a data value from the stack
5         Post-condition:
6             the top value is removed from the stack
7         Return:
8             returns the value at the top of the stack
9         """
10        return self.data.pop()

```

### 12.3.5 Operation: peek()

The peek() operation returns a reference to the data element at the top of the stack. To implement this operation, we use Python's list indexing. This method returns the value, without removing it.

```

1     def peek(self):
2         """
3         Purpose

```

```
4         returns the value from the front of stack
5         without removing it
6     Post-condition:
7         None
8     Return:
9         the value at the front of the stack
10    """
11    return self.data[-1]
```

### 12.3.6 Summary

The similarity between the queue and stack implementations is very strong. The names given to the key operations differ, but both ADTs have an operation to add to the collection, and an operation to remove from the collection. The list-based implementation simply provides a discipline on the use of Python's lists in an effort to prevent programming errors. For example, using enqueue means you don't have to memorize which end of the list you're using for the back of the queue. The operation takes care of that detail. And if there is a good reason to change which end of the list is used as the back and which is used as the front, that change can be made in the ADT implementation, once and for all.

Keep in mind that the computer does not define up or down, or front or back the way we do in the real world. Gravity, for example, gives us a preferred orientation for real world stacks. But without gravity, the top of a stack is defined only by the push() and pop() operations. Likewise, for a queue, its front (and back) is defined by dequeue() (and enqueue()).

We should understand this technique of adapting a list as having very strong advantages for large software projects and long term investments in software development. For one-off scripts, and scripts that will get re-written almost immediately, going to the trouble of using an adapter like this is probably not worth the effort. Also, if your short-lived script has a bug that you cannot find, adapting that list you are using as a stack or queue might be a way to ensure that the bug is not in your use of lists.

Finally, the general strategy of taking a data structure that you already have, and making it look like something you need, is very powerful for rapid prototyping. For example, imagine if Python did not have dictionaries. We could adapt a Python list to behave like a dictionary by creating an adapter ADT, but it would probably not be very efficient. Once the ADT operations have been designed and implemented, we can replace the list-based implementation with something better, and every piece of code that used the ADT would run more efficiently, without having to change any of the application code.



## 13 — The Node ADT

### Learning Objectives

After studying this chapter, a student should be able to:

- To describe the concept and structure of a node.
- To explain the operations of the Node ADT.
- To employ Node ADT operations in Python programs.

### 13.1 Nodes

In computer science, the term *node* is used when talking about the way data is connected. In this chapter, we'll use the word node as the name of a data structure that allows us to build connected chains containing data. A node is a simple object with two attributes:

**data** A data value

**next** A node (or the value None)

A node's data value can be any kind of value needed by the application code. However, the attribute `next` can only be a reference to another node; it is a programming mistake to put something else there.

The primary operations provided for the node ADTs are as follows:

- `node(data, next)` creates a node to contain the data value and the next value.
- `get_data()` returns the contents of the data field
- `get_next()` returns the contents of the next field
- `set_data(v)` sets the contents of the data field to `v`
- `set_next(n)` sets the contents of the next field to `n`

Our Python implementation is very simple. For `node()`, we use the *keyword parameter* technique, to provide a default value for the `next` field, which can be over-ridden by calling `node()` with a



keyword argument:

```
class node(object):

    def __init__(self, data, next=None):
        """
        Create a new node for the given data.
        Pre-conditions:
            data: Any data value to be stored in the node
            next: Another node (or None, by default)
        """
        self.__data = data
        self.__next = next
```

The two operations `get_data()` and `get_next()` simply look up the appropriate field in the dictionary.

```
def get_data(self):
    """
    Retrieve the contents of the data field.
    Return
        the data value stored previously in the node
    """
    return self.__data

def get_next(self):
    """
    Retrieve the contents of the next field.
    Return
        the value stored previously in the next field
    """
    return self.__next
```

The operation `set_data()` stores the given value in the node in the data field. Keep in mind that the reference to the value is being stored, not the actual value itself.

```
def set_data(self, val):
    """
    Set the contents of the data field to val.
    Pre-conditions:
        val: a data value to be stored
    Post-conditions:
        stores a new data value, replacing existing value
    Return
        none
    """
    self.__data = val
```

The operation `set_next()` stores the given value in the node in the next field. Keep in mind that the reference to the value is being stored, not the actual value itself. This is especially important for `set_next()`, since its purpose is to store another node (as a reference). As we know, all values are accessed by their reference, and any value at all could be stored in the next field. We have to practice strong discipline to make sure that the value stored in the next field is either `None` or another node.

```
def set_next(self, val):
    """
    Set the contents of the next field to val.
    Pre-conditions:
        val: a node, or the value None
    Post-conditions:
        stores a new next value, replacing existing value
    Return
        none
    """
    self.__next = val
```

## 13.2 Node chains

With the node ADT, we can create a primitive kind of list-like compound data type, which we will call a node-chain. In many programming languages, node-chains are the main compound data type. The first such language, named LISP (for “LISt Processing”), made this kind of node chain central to its operation, so much so that LISP programs themselves look like lists. We use the term node-chains, because it is descriptive and accurate. However, some text books use the term *linked-list*, which would be fine, except that we want to use that term for a higher level, more abstract concept. Because of their ubiquity, it is instructive to see how node-chains work, as we will need the technique in Chapter 20. This technique is so important, and so useful to build programming skills, that it’s common for second year courses to require implementation of node-chains in other languages, like Java, and C++. If your intention is to continue on to second year computer science, you will do yourself a favour if you dig deep here, and really figure them out.

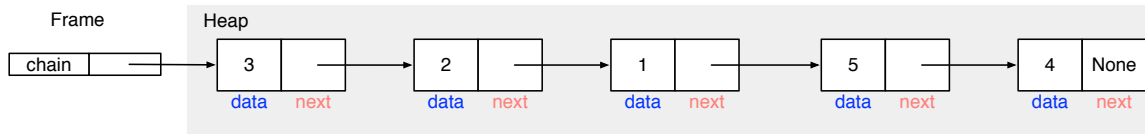
**Definition 13.1** A node-chain is defined as follows:

1. A structure with exactly zero nodes is an *empty node-chain*.
2. A structure with exactly one node is a *node-chain* **only if** its `next` field refers to an empty node-chain.
3. If  $c$  is a non-empty node-chain, then a node,  $n$ , whose `next` field refers to  $c$  is a *node-chain*.

A node-chain is nothing more than a sequence of nodes, one containing a reference to the next. A node-chain has to have an anchor, or a hook, to attach the first node in the chain. Node chains have an implied directionality, namely in the direction suggested by the `next` field. Each node in the node-chain refers to another node, except the last, which has to refer to an empty node-chain.

We will use the Python value `None` to represent an empty node-chain. Note that this value is not actually a node, as defined above. Trying to use any of the node ADT operations on it will result in a run-time error, so we have to be careful to check whether a given node-chain is empty or not, before

Figure 13.1: A node-chain diagram, corresponding to a sequence of node creations.



we do anything to it. A consequence of the use of `None` to represent an empty chain is that the last node in any node-chain has to have the value `None` stored in the `next` field; this is how we know that there are no more nodes further down the chain.

One way to build a node chain is with a sequence of calls to `node()`:

```

1 import node as N
2
3 # create a chain of nodes with data 3, 2, 1, 5, 4
4 chain = N.node(4)
5 chain = N.node(5, chain)
6 chain = N.node(1, chain)
7 chain = N.node(2, chain)
8 chain = N.node(3, chain)
  
```

The first node created (line 4) stores the data value 4; in this function call, we've only given one argument, and the default value `None`, specified in the definition of the function, is used. At that point in the code, the variable `chain` refers to a chain of exactly one node. On line 5, a new node is created, with data value 5, and the current value of `chain` in the `next` field. It is important to understand that on the left-hand-side of the assignment, the value of the variable `chain` is still a single node, containing the data value 3. Creating the new node on line 3 copies the reference stored in the variable `chain` to the data attribute in the new node. Once the new node is created, the reference to the new node is stored in the variable `chain`. At this point in the code, the variable `chain` refers to a node containing 5, and there is a reference from this node to another node containing 3. In other words, `chain` refers specifically to the first node in the chain of 2 nodes, but indirectly, through stored references, `chain` also refers to the whole chain.

A picture of the resulting chain is found in Figure 13.1. Hopefully, you can now appreciate all the work we had to do to explain references. The end of the code sequence above results in a chain of 5 nodes, with 5 different data values. Another thing to keep in mind is that the diagram in Figure 13.1 is an abstraction of the truth: the variable `anchor` is in Python's global scope, which we have not bothered to show, and the node values are on the heap. The biggest simplification in the diagram is that the integers are drawn inside the boxes; the true picture would have each number as a value on the heap, and arrows from the data fields to the corresponding numbers.

The sequence created the node for the value 4 first, and kept putting the new nodes *in front* of the nodes previously created. You should recognize this as being similar to how a stack operates. Be aware that this is our first example, and that it is possible to put new nodes anywhere in the chain, with a little care.

Writing code for node-chains is instructive. As a first example, here's Python code that will display the even data values in the chain:

```
1 # Display the even numbers in the chain
2 anode = chain
3 while anode != None:
4     val = anode.get_data()
5     if val % 2 == 0:
6         print(val)
7     anode = anode.get_next()
```

In this example, we start by making a copy of the reference stored in the variable `chain` from the previous example, and storing it in the variable `anode`. We will use this variable to step through the chain, one node at a time, using a while loop. We must use a while loop for our nodes, because Python has no idea how to step through our chain using a for-loop. The condition for the while loop checks whether we've reached the end of the chain; remember, that the value `None` is used to represent an empty chain. There are no more nodes to look at if the variable `anode == None`. Notice line 7 especially. This is a key technique! Line 7 replaces the reference stored in `anode` with the reference stored in `anode's` next field. Visually, this is stepping from one node to another by following the arrow out of the box labelled next.

It's important to point out that we always use a temporary variable to step through a node-chain. While it is possible to use the variable `chain` for this purpose, it is rarely correct, except for trivial examples. We set up the variable `chain` to be the anchor of our node-chain. If we change the value of this variable, we change the anchor point for the chain, and we lose the original anchor point. If we step through the node-chain by changing the anchor point, and we reach the end of the node-chain, we have effectively unhooked the whole node-chain from the anchor. It's really important to understand that without any anchor point, we have no way to access the node-chain when we're done, and it is absolutely, totally, and in all other ways, inaccessible.

This is not specific to nodes; any data value in the heap is lost if the reference to it is not stored somewhere. Python will simply sweep through the heap reclaiming the memory used by data not connected somehow to a frame; this process is called *garbage collection*, though it is more precisely analogous to recycling.

While we could (and should) write a lot of code to demonstrate the kinds of programs to manipulate node-chains, we will look at just one more example, specifically, removing a node from the far end of the chain. A node-chain can be arbitrarily long, even though our example chain only has 5 nodes in it. To reach the end of the node-chain, we could step through it one node at a time. We know that we've reached the last node of the chain when the node's next value is `None`. To unhook the last node from the chain, we have to set the next field of the second-last node to `None`. However, if we actually step to the last node, there is no way to step backwards across a reference (an arrow in our diagram). Another standard technique for working with node-chains is to use two temporary variables, referring to two adjacent nodes in the chain. If we step both nodes through the chain, when one node is the last node, the other will be the second last node.

Here's the example.

```
1 anode = chain
2 prev = None
3 while anode.get_next() is not None:
4     prev = anode
5     anode = anode.get_next()
```

```
6
7 if prev is not None:
8     prev.set_next(None)
9 else:
10     anchor = None
```

The two variables are `anode` and `prev`; `prev` will always be the node that contains a reference to `anode`. On line 3, the while loop condition checks whether `anode` is in fact the last node in the node-chain. On line 4, we advance `prev`, and on line 5, we advance `node`. Lines 4-5 ensure that `prev` is always one node behind `anode`. On line 6, we check whether `prev is not None`. The only way for `prev` to be `None` is if the first node in the node-chain is the only node in the node-chain, and removing the last node means removing the only node!

## 14 — Programming practice and style

### Learning Objectives

After studying this chapter, a student should be able to:

- To be aware of stylistic concerns in programming, such as the importance of good variable names, program documentation, and clarity of code.

### 14.1 Introduction: Style counts

Overwhelmingly, the focus of first year computer science is on getting the computer to do the thing that the assignment requires them to do. Admittedly, a program that doesn't work is worthless in this context. However, it must be stated with extreme emphasis that a program that cannot be understood by a person, whether that be an instructor, a marker, a colleague, or yourself, is also worthless. In first year, you write a program, hand it in, and then discard it. This pattern is true during your training, but it will absolutely false when you are out there doing real work. As you gain in experience, the exclamation "But it works!" becomes less and less relevant.

Style is important because valuable software will be read by human beings a lot more than you might think. You'll be working with code that someone wrote months or years ago, possibly before you joined the project team. You'll have to revisit your own work countless times as you bring a product to market, or as you prepare your data analyses for your scientific work. A well-written program is far easier to debug than a poorly-written one. And no matter how emphatically you tell yourself "I'll fix it later," *later* rarely turns into *now*.

This chapter is an attempt to inspire the kinds of aesthetics students should be reaching for as they reach the end of their first year of computer science. If you fail to employ the kinds of practices we will present in this chapter, you are increasing the burden on everyone in the future who may be forced to read your horrible, but functional, code. Don't be that person who pretends that *using obscure language features* is an indicator of *cleverness*, and who thinks *it was hard to write* means *it*

*should be hard to read.* Your software, produced under normal non-examination conditions, should be elegant, thoughtful, stylish, clear, concise, and of course, functional.

## 14.2 Style Guidelines

This section is going to be driven by examples. Some of these examples are real, in the sense that they were drawn from real programs. In some cases, the examples are contrived to make a point.

### Practice 14.1 Choose good variable names.

You're tired of hearing this, no doubt, but it's so important, that's where we always start. A good variable name is as brief as possible but still conveys the purpose of the variable. A good name reminds the reader of the purpose, so that reading becomes effortless. In contrast, a bad name adds an extra burden to the reader, by forcing them to memorize or deduce the purpose of the variable. Remember, the person being forced to work hard to figure out what your program does may be you yourself!

A good variable name uses words, or short phrases. Here are some examples:

```
clauses = # ...
left_subtree = # ...
number_of_elements = # ...
```

Different programming languages have conventions about how to write multi-word variable names. In Python, the convention is to use *snake\_case* for variables, with lower case words separated by the underscore character. In other languages, the conventions may be different. Organizations and project teams may also have a set of conventions they use; use them!

Poor variable names use acronyms, ambiguous abbreviations, or just arbitrary characters. For example:

```
cs = # ...
l = # ...
n = # ...
```

When you're just blasting out code, you may find that using short variable names helps you get your ideas expressed quickly. That's fine; a very common piece of advice for all varieties of writing is to write the first draft quickly, and with minimal care, then go back and do your editing. When you use this strategy, make good use of your computer! Most Integrated Development Environments (IDEs, like PyCharm) have a *Refactoring* menu, which will allow you to make changes to variable names all at once. This tool is slightly more advanced than your word processor's search-and-replace, because the IDE knows the rules of your programming language. The point here is that you can't use the excuse of "It takes too long" to turn poor variable names into good ones.

### Practice 14.2 Use short names for loop control variables.

This is an exception to the previous advice. If your variable only has a short lifetime, say as part of a loop, a short name is preferred. Names like *v* or *val* are good in for-loops, and variables like *i* and *j* are good for indexed while-loops. Here's a trivial example:



```
for v in alist:
    # do stuff with v
```

When your loop control variables' names are too long, it actually increases the burden on your reader. For this purpose, long variable names tend to obscure the purpose of the loop:

```
for value_found_in_alist in alist:
    # do stuff with value_found_in_alist
```

You have to read the for-loop a lot more carefully to figure out what the variable is, and what it's ranging over.

**Practice 14.3** Avoid using variable names *l* (ELL), *I* (upper case i), and *o* (OH). ■

This advice is just self-defence. Depending on the font you're using in your editor, the characters *l*, *I*, and *l* may be similar enough to cause confusion. Here's an example:

```
for l in clauses:
    # do stuff with l
```

The same is true for *o*, *O*, and *0*. The best thing to do avoid these two single letter variable names entirely. Here's a better version:

```
for cl in clauses:
    # do stuff with cl
```

**Practice 14.4** Don't over-comment. ■

Adding comments to your code can be immensely helpful. Keep in mind that as you gain experience coding, you will find the code itself to be increasingly informative on its own. The kinds of comments that would have helped you when you were just learning about for-loops won't be needed now that you've mastered them.

Comments should be about what's *not obvious* in your code. Adding comments about obvious concepts simply adds clutter, and distraction. Here's an exaggerated example of cluttered commenting:

```
i = 0 # list index for summing the list
total = 0 # running total for summing the list
while i < len(some_list):
    total = total + some_list[i] # add to the running total
    i += 1 # increment i
```

The problem with the comments above is that they concern the individual statements, but say nothing helpful about the larger purpose. The comments in the body of the while loop are superfluous, as they say in English exactly what the Python says. This is not helpful to anyone, and should be avoided. Here's a better version.

```
# calculate the total cost from the list of costs
i = 0
total = 0
while i < len(some_list):
    total = total + some_list[i]
    i += 1
```

Use comments to explain the purpose of a function, or, if the algorithm is tricky, the intuitions about the algorithm. Another thing you should do is include citations to books, articles, or websites if you are using an algorithm that you got from someone else. Assuming this is permitted by your assignment requirements! A reminder will save time repeating the research you did to find it in the first place.

**Practice 14.5** Write clearly – don’t be too clever. ■

The following example makes use of Python’s list comprehensions. The problem is that the code does too much in one line. Experienced Python programmers encourage the use of comprehensions when they make the code easier to understand, but that’s not the case here:

```
new_cl = [[l for l in c if (-1)*flip*first_lit != l]
           for c in clauses if flip*first_lit not in c]
```

The following rewritten code accomplishes the same calculations.

```
new_clauses = []
for clause in clauses:
    if flip * first_lit not in clause:
        new_clause = []
        for literal in clause:
            if (-1) * flip * first_lit != literal:
                new_clause.append(literal)
        new_clauses.append(new_clause)
```

The important point here is not about comprehensions. It’s about writing code that does what it says clearly. Code that’s too cleverly constructed may contain bugs that are even more cleverly hidden.

**Practice 14.6** Break complex expressions into smaller pieces. ■

The example above is also demonstrative of this advice. In Python, the Python runtime system analyzes the code before it is executed, and finds ways to make your code run quickly on the computer you are using. In fact, all programming language systems do this, sometimes more, and sometimes less. The point is you don’t have to worry about making little decisions about what the computer should do all at once. You can afford to make your program clear, and let the system make instruction-level optimizations.

**Practice 14.7** Replace *repeated* or *copied sequences* by calls to a *common function*. ■

This is the essence of procedural abstraction, about which we've said a lot already. The most egregious demonstrations of repetitive expression come from the use of copy/paste, which is very tempting. But copy/paste is a terrible tool because it duplicates errors, and introduces opportunity to add errors. Here's an example, in which a simple while loop is used to calculate a small number of factorials:

```
1 a = 10
2 b = 12
3 fa = 1
4 for i in range(1,a):
5     fa = fa * i
6 fb = 1
7 for i in range(1,b):
8     fb = fb * i
9 fab = 1
10 for i in range(1,a+b-1):
11     fab = fa * i
12 print('Mario paths:', fab/(fa*fb))
```

No doubt you've spotted the bug on line 11, which was introduced to demonstrate the kind of bug that gets introduced to code when it is copied/pasted. You may also note that the three loops don't quite calculate factorial, even when corrected. Is this a bug, or is this deliberate? It's far from clear unless you understand the purpose of the calculation on line 12.

It's far better to replace repetitive expressions to calls to a common function.

```
1 def fact(n):
2     f = 1
3     for i in range(1,n+1):
4         f = f * i
5     return f
6
7 a = 10
8 b = 12
9 print('Mario paths:', fact(a+b-2)/(fact(a-1)*fact(b-1)))
```

With a single function replacing repeated code, any bugs in the algorithm are localized, and constrained (for factorial, at least). Also notice how using the function helps to clarify the intention of the calculation: For whatever reason, the expression on line 9 requires calculation of factorials that are one less than  $a$ ,  $b$  and  $a + b - 1$ . It seems deliberate, whereas the previous implementation seemed almost accidental.

## 14.3 Improving your code

This section is intended to provoke some ideas about how students can improve their code, without changing underlying algorithms, or worrying about efficiency. The suggestions here can improve robustness, but the basic unit of improvement is readability.

**Practice 14.8** Turn long conditionals into lookups. ■

In our courses, we presented conditionals long before any of the compound data structures like lists. But those long conditionals that you wrote when you were first learning them are not the only way, nor even a good way, to structure your code. Take for example the following partial example of an conditional about month names:

```
if month == 1:
    print('Jan')
elif month == 2:
    print('Feb')
# etc, 12 months
```

It doesn't matter what the code does; we're looking at its structure. A long conditional adds a lot of code without actually contributing clarity. And while no one expects the number of months to change, there are applications of long conditionals that might be a little more dynamic.

Using a list to replace a long conditional is consistent with the concept of encapsulation. It essentially converts code (the conditional) into data (a list or dictionary):

```
month_names = ['Jan', 'Feb', 'Mar', # and the rest
print(month_names[month-1])
```

The long conditional is reduced to a simple list access, and no clarity or purpose is lost. When the look-up is based on data that's not an integer, we can use a dictionary instead:

```
month_dict = {'Jan':1, 'Feb':2, 'Mar':3, # and the rest
month = month_dict[month_name]
```

**Practice 14.9** Factor common statements out of loops and if-statements. ■

This advice takes advantage of the mathematical nature of algorithms. There are patterns in code that can be simplified, in the same way that expressions in mathematics can be simplified. The first example demonstrates the use of sequence in a conditional:

```
if x >= 0:
    y = x
    print('y is positive')
else:
    y = x
    print('y is negative')
```

In this example, Notice that  $y = x$  is the first thing done in both branches, and that it is independent of the condition  $x \geq 0$ . For that reason we can revise the code as follows:

```
y = x
if x >= 0:
    print('y is positive')
```

```
else:  
    print('y is negative')
```

We moved the common statement so that it occurs before the conditional because it appeared first in sequence. The revision does not change the behaviour of the code. It makes the code more concise, and probably clearer.

Similarly, we can transform the next example by moving the common statement after the conditional:

```
if x >= y:  
    x = x / 2  
    print(x)  
else:  
    x = x * 2  
    print(x)
```

Here, the value printed depends on the change that is being made to  $x$ , so moving the print statement before the conditional would change its behaviour. But moving it after would not change the behaviour.

```
if x >= y:  
    x = x / 2  
else:  
    x = x * 2  
print(x)
```

A similar thing can be done with statements in loops. For example:

```
for v in range(10):  
    x = 0  
    print(v*x)
```

Again, the statement  $x = 0$  does not depend on  $v$ , so it is repeated unnecessarily. Here's an obviously better version:

```
x = 0  
for v in range(10):  
    print(v*x)
```

Here, we have to move the statement  $x = 0$  before the loop, because the loop depends on the value of  $x$ . A symmetric case can be made for moving code from inside a loop so that it is executed after the loop.

The point of this advice is that code has behaviour that can be reasoned about mathematically. The concepts you learned in your math class can be applied, if indirectly, to improving the quality of your computer program. Sometimes, mathematics is not about math at all, but about creating an environment where certain patterns of reasoning can be taught, so that those patterns of reasoning can be applied outside of mathematics. In other words, *abstraction*.

**Practice 14.10** Don't check for end-of-loop inside your loops! ■

It's fairly common for novice programmers to worry about doing too much all at the same time. For example, the following code is a counted loop from 0 to  $N$ , but contains a conditional that executes statement  $A$  when  $i = N - 1$ .

```
i = 0
while i < N:
    statement B
    if i == N - 1:
        statement A
    i += 1
```

The condition  $i = N - 1$  is only true on the very last iteration of the loop. In English, this loop says “Keep doing  $B$ , and at the end of the loop, do  $A$ .”

The main problem with this implementation is that it obscures the task unnecessarily. A far better structure is to realize that we don't need to keep checking if the loop is over. We can simply move statement  $A$  outside the loop, as follows:

```
i = 0
while i < N:
    statement B
    i += 1
statement A
```

As an added bonus, the cost the body of the loop is reduced. The initial implementation checked for the end of the loop twice for every value of  $i$ , but the improved version only checks once per value.

**Practice 14.11** Don't use conditionals when a Boolean expression will do. ■

Because Boolean and relational expressions are required for conditionals and loops, it is common for novice programmers to believe implicitly that the only way to use them is in the context of a conditional or a loop. One of the key insights in a novice programmer's progression towards expertise is the realization that Boolean expressions can be used anywhere that any other kind of expression can be used. For example, consider the following function:

```
def function(a, p):
    if a == 3 or p < 0.005:
        return True
    else:
        return False
```

It returns a Boolean value, based on a Boolean expression. A far more concise implementation of this function simply uses the value of the expression as a return value:

```
def function(a, p):
    return a == 3 or p < 0.005
```

Boolean expressions represent calculations that evaluate to Boolean values, which can be assigned to a variable name, returned from a function, or used as arguments to other functions. Boolean expressions are not second-class citizens in computer science; they are absolutely equal in status to any other kind of expression. They simply evaluate to a different data type.

**Practice 14.12** Don't compare a Boolean expression to a Boolean literal. ■

When novices first learn about if-statements and loops, they see lots of examples of conditions comparing numbers using relational operators. It is natural to over-generalize and use the same pattern when dealing with Booleans. For example, consider the following if-statement:

```
if Stack.is_empty(ystack) == True:
    print('Stack is empty')
```

This condition can be simplified by recognizing that the expression `Stack.is_empty(ystack)` is a Boolean value that does not need to be compared to a literal in the context of an if-statement. We can use the expression all by itself to replace a comparison to `True`:

```
if Stack.is_empty(ystack):
    print('Stack is empty')
```

This simplifies code, making it more readable.

The same kind of transformation can be done when the comparison is the the Boolean `False`, but we need to use the operator `not` as well:

```
while not Stack.is_empty(ystack):
    print('Popping', Stack.pop(ystack))
```

We can avoid the over-complicated condition that states `Stack.is_empty(ystack) == False`.

**Practice 14.13** Don't use a while loop to step through a list, even if you need indices! ■

This advice is specific to Python. Suppose you're implementing a simple linear search through a list of items, and you need to know the index of the item:

```
while i < len(some_list):
    if some_list[i] == value:
        print('Found', value, 'at index', i)
    i += 1
```

The use of a while loop gives us control of the index variable, *i*. This is the way we taught you to do this kind of thing, but only because we can't teach you every thing all at once. The better way to write a loop like this is to use more Python, specifically, the function `enumerate`, which takes a sequence, and returns a sequence of tuples, namely the index and the value. Here's the same loop, rewritten:

```
for i, v in enumerate(some_list):
    if v == value:
        print('Found', value, 'at index', i)
```



This makes your code more concise, and eliminates common sources of error: forgetting to increment the while-loop counter (or incrementing it incorrectly), or indexing the list incorrectly. In more primitive programming languages, this kind of abstraction is not available, or it's not core to the language.

**Practice 14.14** Don't *patch* bad code – *rewrite* it. ■

Sometimes there's no way to *nudge* a program to fix it or make it better. While it is contrary to human nature, it is important for programmers to understand that sometimes you have to treat your first implementation as a *prototype*. The thing about prototypes is that they are useful only to help you understand what the algorithmic issues are. In real software development teams, prototypes are deleted and *replaced completely* all the time. Learn when *you personally* need to delete and start again from scratch. Throwing away code seems wasteful, especially when an assignment deadline is looming. But if you start early enough, you can learn a lot from a quick and dirty prototype, and then you can restart the development with a little fore-knowledge of what will be needed. And if you are taking advantage of version control, deleting your code is completely safe. You always have access to your previous versions!

## 14.4 Making your program faster

We have already mentioned the concept of efficiency, and we will spend some time to study an analytical technique to determine the efficiency of our algorithms. And as we saw above, sometimes a simple modification to code can improve its speed. This is all well and good, as long as you don't worry about this kind of thing too soon. Before you do any optimization, improve your code for clarity and robustness. If you can improve the clarity of your code, you might be able to find ways of making it more efficient that might not be possible when the program is convoluted or obscure.

**Practice 14.15** Make sure your code is correct *before* you make it faster. ■

There's no point in speeding up code that doesn't do the right thing. An inefficient implementation that works is far better than an efficient implementation that doesn't work. *Keep your code correct right* when you make it faster. Test your inefficient version thoroughly, too. Make sure that your test results do not change once you've started optimizing. Finally, make sure your function is *important enough* to make faster. If your function only gets called once, optimizing it may not contribute much value to your program. On the other hand, optimizing a function that gets called a lot will have a far greater effect on your program.

## 14.5 How to improve your programming skills

To develop good programming style, you have to read code. Study examples given in class, and in books. Browse [Python.org](https://python.org) and Stack Overflow. Read as much code as you write. Learn the idioms of the programming language you are using. Always practice coding. Write little programs every day. Review your programs from past work, and try to make them a little better. When you start new assignments, imitate the style of your instructors. You don't have to produce a work of art every

---

time you produce a program, but you can make it a habit to assess the quality of your own work. Keep a journal of the errors that caused you grief, or took a long time to debug.

Learn to use all the functionality of your IDE. Don't wait for a class to teach you these things. There are really valuable tools you can learn just by exploring the IDE menu. Learn to touch type. Stop using the mouse and menu. Use and memorize key-commands for your IDE or your text editor.



## 15 — Defensive Programming

### Learning Objectives

After studying this chapter, a student should be able to:

- To describe practices for defensive programming, such as assertions and exceptions.
- To describe different techniques for avoiding or handling floating point equality.
- To describe different techniques for avoiding or handling divide by zero.

### 15.1 What could possibly go wrong?

Way back in Chapter 4, we mentioned the goal of robustness. We said that software is *robust* if it behaves well even when things go wrong. In this chapter we'll give some tools and concepts that we can use to avoid things going wrong, and what to do when it can't be avoided.

Things can go wrong if the user enters bad data, or if a network connection is lost, or any other kind of external interaction behaves unexpectedly. To assess *robustness*, we consider what kinds of things might go wrong that are beyond the control of the program. A program that is not robust is called *brittle*. Brittle software crashes without warning, causing data loss or other inconveniences. If a program still has errors, a robust program might be able to recognize that an error had occurred, and recover from the error without causing the program to crash. Even if a crash is unavoidable, it would be useful for a program to do something helpful, like saving a document, or reporting that an error had occurred with a helpful error message.

### 15.2 Defensive Programming

Defensive programming is an attitude that programmers practice when they stop pretending everything is fine. As a novice, assuming the best lets us learn the basic principles of programming, but if anything can go wrong, as Murphy's Law states, it will. As you gain more experience programming,

you will pick up habits that will help you avoid things going wrong, but to get you started, here are a few pieces of advice to help prevent errors.

Fundamentally, a defensive attitude ensures that a program protects itself against incorrect or illegal data. Preconditions and return values need to be checked; conditionals need to include a branch for those events that you assume can't possibly happen (because they will happen).

Many programming languages, including Python, have a tool for checking conditions that might identify if something has gone wrong. The Python command `assert` can be used to check a condition, and if the condition is false, it will halt the program with a run-time error. We will see a number of examples of its use. It is useful in situations where it is not clear what should be done if something has gone wrong.

**Practice 15.1** Enforce pre-conditions using assertions within a function. ■

We have placed great importance on the specification of a function, to include a description of what values the function accepts as arguments. We assumed that the code making the function call will play nicely, by sending the right kind of information. But we should, when possible, use an assertion to check, as in the following example:

```
def factorial(n):
    """
    Purpose: Compute the factorial.
    Precondition:
        :param n: a non-negative integer
    Return
        :return: a non-negative integer
    """
    assert n >= 0, 'invalid input to factorial'
    if n == 0:
        return 1
    else:
        return n*factorial(n-1)
```

The assertion checks whether  $n$  has an appropriate value, and if `factorial()` is ever called with a negative number, the assertion will *fail*, and Python will halt execution of the program, sending the given message to the console as part of the error report:

```
Traceback (most recent call last):
  File "fact.py", line 16, in <module>
    factorial(-5)
  File "fact.py", line 9, in factorial
    assert n >= 0, 'invalid input to factorial'
AssertionError: invalid input to factorial
```

**Practice 15.2** Check post-conditions in test scripts. ■

You can use assertions to check for impossible situations outside of functions as well. They are very good for wolf-fencing and test scripts too. For example, here's how we might use an assertion in a test script:

```
input = [1, 3, -5, 7, 9]
expected = -5
result = min(input)
assert result==min, 'Expected'+str(expected)+'found'+str(result)
```

The assertion states what should be true, and replaces an if-statement controlling a print statement.

**Practice 15.3** Use assertions liberally in prototypes, test scripts, and programs that are currently under development. ■

An assertion that fails brings the whole program to a complete halt. This behaviour is fine when the program is in development, but it's a bit less ideal in a completed application. For this reason, Python allows you to turn off assertion checking, so that you don't have to edit the code and delete the assertions yourself. For your own protection, we won't show you how to do that yet.

In production code (i.e., code that you've delivered to a client), when coming to a complete halt would ruin the user's experience, it's not appropriate to use an assertion to check for faults. A more versatile tool that can be used when things go wrong as called an *exception*. An exception is like an assertion, in that it can bring a program to a halt. But unlike an assertion, an exception doesn't always bring a program to a halt; it can be *caught*, which allows programs to recover from some kinds of error situations. We will not say any more about exceptions for now.

**Practice 15.4** Annotate parameters with Python type hints. ■

In Chapter 8, we tried to convey the importance of data types in software development. We noted that Python's types were implicit, and type information for each object is stored in the object itself. We also noted that many computer scientists and software developers consider this one of Python's most serious weaknesses. Many other languages (e.g., Java, C++) require the programmer to declare type information as part of the program.

To address this weakness, the inventors of Python added a feature that allows Python programmers to make similar annotations in Python programs. Consider the following function definition, especially the first line:

```
def factorial(n : int) -> int:
    """
    Purpose: Compute the factorial.
    Precondition:
        :param n: a non-negative integer
    Return
        :return: a non-negative integer
    """
    assert n >= 0, 'invalid input to factorial'
    if n == 0:
```

```
        return 1
    else:
        return n*factorial(n-1)
```

The parameter `n` is annotated with a *type hint*, which consists of a `:` followed by a type name. Also, notice the `-> int` that appears at the end of the first line, which indicates that the function is intended to return an integer. These annotations are called hints, because IDEs like PyCharm can use them to highlight possible errors in their function calls; however, as of this writing, Python takes no action when these hints are violated. In future versions of Python, we can expect Python to refuse to execute a script where the data is not consistent with the type hints, just as languages like Java and C++ do by default.

### Practice 15.5 Don't shadow Python functions. ■

As we learned in earlier chapters, variables in Python are just names, and almost anything can be attached to a name. This includes functions! A Python function has a name, but the name is simply a variable whose value refers to an encapsulated algorithm. For this reason, Python programmers can make the mistake of creating variables whose name is the same as a function.

```
sum = 0
for val in my_list:
    sum += val
```

Normally, this is an mistake, though sometimes it has no bad consequences, so it may not be as serious as an error. In the above example, the name `sum` is assigned the value 0, and the value 0 *replaces* its previous value, which is the well-known function capable of adding stuff up. There is nothing special about names in Python; the values are special, though. A function can be called, but an integer cannot be called.

A more subtle problem related to replacing a value is called *shadowing*.

```
def total(my_list):
    sum = 0
    for val in my_list:
        sum += val
    return sum
```

Here, the function `total` includes a variable named `sum`, which has the same name as the Python function `sum`. But the scope of the local variable is the function `total`, not the whole program. In this case, the use of the name `sum` refers to the local variable, and the Python function is not accessible; it is *shadowed*. The word shadow describes the situation well: the function `sum` cannot be seen because local variable `sum` gets in the way.

### Practice 15.6 Avoid modifying list structure while iterating over it. ■

Python lists are mutable, which means they can be changed, either by changing the contents, or by changing the structure. While changing the contents of a list is pretty intuitive, changing its structure requires more care. When you delete or pop an item, the indices of some of the elements may change.



For example, if we delete an item at index 5, the item formerly at index 6 is re-indexed, and takes position at index 5. Every item after index 5 is re-indexed. It's important to emphasize that deleting an element behaves *as if* all the items are shifted forward, but in Python, the lists don't actually shift; but their indices do change! It's important to remember this behaviour when you are trying to change the structure of a list using a loop. For example:

```
a = list(range(10))
i = 0
while i < len(a):
    a.pop(i)
    i += 1
```

Because pop causes elements to be re-indexed, this loop only removes the even numbered items. For this reason, it's probably a good idea to avoid modifying the structure of a list while you are iterating over it. It is far less trouble to create a new list with the structure you want, from the old list.

### Practice 15.7 Watch out for equality of floating point.

Floating point calculations have tiny errors due to the fact that floating point numbers have finite precision. Every calculation introduces a little more error, but some calculations add a lot more error than others. As a result, two floating point values are hardly ever equal, even if the math tells us that real values should be equal.

One of the biggest misconceptions about floating point numbers is that numbers that look simple in decimal notation are simple when stored in a computer. Here's an example:

```
x = 0.3
y = 0.1 + 0.1 + 0.1
if x == y:
    print('Equal')
else:
    print('Not equal!')
```

In this example,  $x$  and  $y$  are not equal, because in binary, 0.1 has an infinitely repeating representation, much like  $1/3$  is in decimal. As a result, the floating point value corresponding to 0.1 gets chopped off after 50 or so bits, and adding three of these truncated values no longer adds up to the floating point value corresponding to 0.3.

When your values are in normal, every day ranges, you can avoid odd behaviour by using a smaller value as an absolute threshold, like 0.000001:

```
x = 0.3
y = 0.1 + 0.1 + 0.1
if abs(x - y) < 0.000001:
    print('Equal enough')
else:
    print('Not equal enough!')
```

The key behind this technique is that the threshold has to be smaller than the values you are comparing. If they are not, small but significant differences may be misinterpreted as unimportant using an

absolute threshold.

Another technique is to use a threshold whose size is relative to the size of the values you are comparing:

```
x = 0.3
y = 0.1 + 0.1 + 0.1
if abs(x - y)/max(abs(x), abs(y)) < 0.000001:
    print('Equal enough')
else:
    print('Not equal enough!')
```

Here, we're comparing the relative difference to a threshold, which works, as long as  $x$  and  $y$  are not both zero.

### Practice 15.8 Watch out for division by zero. ■

The way to handle division by zero depends on your situation. If you have control over the denominator, you can use a conditional. For example, in the previous example about checking equality using a relative threshold:

```
x = 0.3
y = 0.1 + 0.1 + 0.1

if x == 0 and abs(y) < 0.000001:
    print('Equal enough')
elif y == 0 and abs(x) < 0.000001:
    print('Equal enough')
elif abs(x - y)/max(abs(x), abs(y)) < 0.000001:
    print('Equal enough')
else:
    print('Not equal enough!')
```

This can be done because the case for a zero denominator tells you something important: if both values are zero, you don't need to divide at all to check for equality.

If the division is part of a formula, and the denominator value comes from some other part of the program, you should just let the division go through:

```
def ratio(x, y):
    """Compute x/y for some reason.
    Preconditions: y != 0
    """
    return x/y
```

You might be tempted to use an assertion here, but division already raises an error when zero is used as a denominator, so checking yourself is extra work. Also, Python uses an exception for divide by zero, and it's possible that some other part of the program will catch this exception, so using an assertion will interfere with the exception handling.

Finally, if the application domain has special requirements, you may be able to handle a division by zero using those requirements. For example, in some scientific calculations, division by zero can be handled by returning a particular value, maybe even zero:

```
def ratio(x, y):  
    """Compute x/y for some reason.  
    If y == 0, ratio returns 0 because SCIENCE.  
    Preconditions: x, y are numbers  
    """  
    if y == 0: return 0  
    else: return x/y
```

Obviously, this is not a general solution. But zeroes headed towards a denominator don't always mean trouble.



### Why use nodes to implement Queues and Stacks?

#### Using Nodes to Implement Stacks

Operation: Stack()  
Operations: size() and is\_empty()  
Operation: peek()  
Operation: pop()  
Operation: push()

#### Using Nodes to Implement Queues

Operation: Queue()  
Operations: size() and is\_empty()  
Operation: peek()  
Operation: dequeue()  
Operation: enqueue()

#### Summary

## 16 — Node-based Stacks and Queues

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe the node-based implementation of stacks and queues.

### 16.1 Why use nodes to implement Queues and Stacks?

We started the course by reviewing lists, and we've already seen an implementation of Stacks and Queues based on Python lists.

In this chapter, we will see how we can implement the Queue and Stack ADTs using the node data structure we introduced in Chapter 13. As we mentioned, the node structure is not an esoteric detail. The idea of using a data structure like a node is one of the key advances in computer science. Prior to the invention of nodes, computer scientists could only store collections of data in finite, contiguous ranges of bytes, commonly called arrays. There were no lists before nodes were invented. Finally, it is important to keep mentioning the idea that we can completely change the implementation of an ADT like a stack or a queue, but as long as the interface stays the same, the behaviour of any application using the ADT is unchanged. This is the very essence of adaptability and robustness.

### 16.2 Using Nodes to Implement Stacks

The essential nature of a Stack is the LIFO protocol. A node-based stack will use a node chain to contain the values in the stack. The front of the node-chain will be the top of the stack. Since the LIFO protocol requires access only to the top of the stack, the methods only need to keep track of one end of the node chain. The bottom of the stack is conceptual, but we won't need the implementation to worry about it.

The attributes we will be keeping track of are as follows:

**top** A reference to the first node in the node chain that stores the stacked values. If the stack is empty, this should be the Python value `None`.

**size** An integer representing the number of elements currently stored in the stack. This value should be non-negative. A size of zero implies that the stack has no nodes, and no data.

The following sections describe the implementation in a little more detail. The complete implementation can be found in [Appendix D](#).

### 16.2.1 Operation: `Stack()`

The `Stack()` function returns a new empty stack. This operation simply sets all the attributes to the values that indicate an empty stack.

```
class Stack(object):

    def __init__(self):
        """
        Purpose
            creates an empty stack
        """
        self.__size = 0          # how many elements in the stack
        self.__top = None       # the node chain starts here
```

For example, using the `Stack()` operation in a script such as

```
import Stack as S

astack = S.Stack()
```

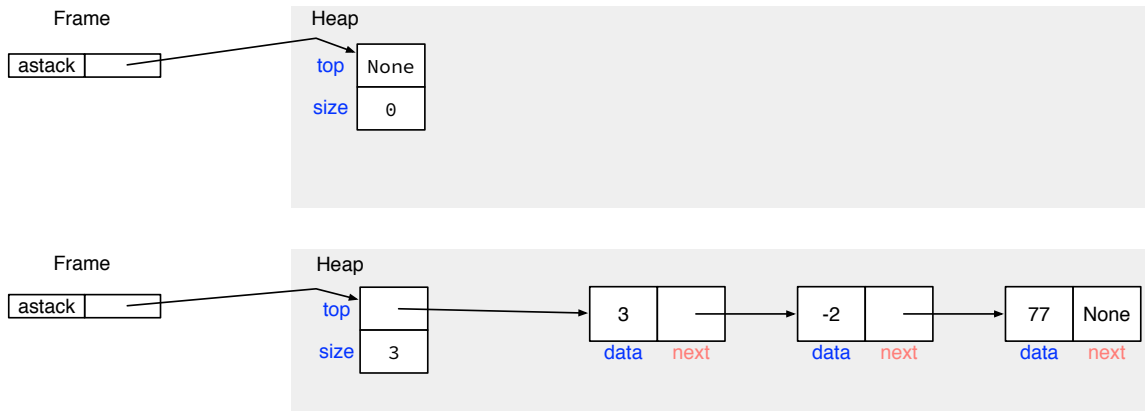
can be visualized as in [Figure 16.1](#). The diagram shows the attribute names in coloured text beside the rectangles. To keep our diagrams as simple as possible, the values are represented as being stored within the `Stack`, but in reality, the object for `None` and `0` are in Python's memory, and the stack stores references to them within the `Stack`. Eliminating one set of arrows helps to focus on the important aspects of the implementation.

### 16.2.2 Operations: `size()` and `is_empty()`

The `is_empty()` operation returns `True` if the queue is empty. The `size()` operation returns the number of elements stored in the queue. Both implementations below simply use the value of the `size` attribute in the object. This is much faster than counting! This implies that the operations for pushing and popping data have to change the `size` consistently.

```
def size(self):
    """
    Purpose
        returns the number of data values in the stack
    Return:
        The number of data values in the stack
    """
    return self.__size
```

Figure 16.1: The top diagram shows an empty Stack, created on the heap by the Stack() operation. The bottom diagram shows what a queue would look like after several values have been pushed.



```
def is_empty(self):
    """
    Purpose
        checks if the stack has no data in it
    Return:
        True if the stack has no data, or false otherwise
    """
    return self.__size == 0
```

### 16.2.3 Operation: peek()

The peek() operation returns a reference to the data value currently at the top of the stack, without removing it from the stack. If the given stack is empty, the operation causes a run-time error, because the assertion will fail.

```
def peek(self):
    """
    Purpose
        returns the value from the front of given stack
        without removing it
    Note: the stack cannot be empty!
    Post-condition:
        None
    Return:
        the value at the front of the stack
    """
    assert not self.is_empty(), 'peeked into an empty stack'

    first_node = self.__top
```



```
result = first_node.get_data()
return result
```

#### 16.2.4 Operation: pop()

The `pop()` operation removes a data element from the top of the stack. In this implementation, the value is removed from top of the stack, but in our implementation, the top of the stack is the first node of the node chain.

If the stack is empty when popped, the operation will cause a run-time error, because the assertion will fail. If the stack is not empty, we make a copy of the reference to the top node, and remember the data stored in that node. Whatever follows the removed node will be the top of the chain, and we decrease the size of the stack.

```
def pop(self):
    """
    Purpose
        Removes and returns a data value from the stack.
        Note: the stack cannot be empty!
    Post-condition:
        the first value is removed from the stack
    Return:
        the first value in the stack, or None
    """
    assert not self.is_empty(), 'popped an empty stack'

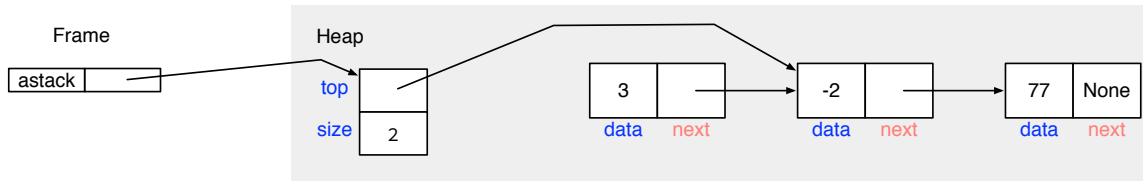
    prev_first_node = self.__top
    result = prev_first_node.get_data()
    self.__top = prev_first_node.get_next()
    self.__size -= 1
    return result
```

For example, if `pop()` were called on the Stack drawn in the bottom half of Figure 16.1, the result would be as shown in Figure 16.2. Notice that the reference stored in `top` now indicates the node containing -2. The diagram shows that the node containing 3 is still present in the heap, but nothing in the node chain refers to it, so it is not part of the Stack. If there is no way to access this node from a frame, Python will reclaim the memory it uses for some other purpose at some time in the future.

It's pretty important to think about how `pop()` works with an empty stack. Popping an empty stack is like dividing by zero, or asking for the logarithm of a negative number. The operation cannot have a right answer. Because there's no right answer, we have to think about what the method should actually do. In the given implementation, we're treating the situation as a serious application error: if the stack is empty the assertion brings the whole program to a halt. Therefore, the application programmer is responsible for checking the stack to ensure it is not empty before the stack is popped.

Programmers have to give consideration to this kind of thing during the design process, and it's easy enough when the situation is familiar, like divide by zero. On the contrary, when you're designing a data structure like a stack, you have to push yourself a little to try to think about what situations might arise that are not typical or expected, like popping an empty stack. We call such situations *special cases*. If you are not already the kind of person who can point out nit-picky details

Figure 16.2: The effects of a pop operation on a non-empty stack. The reference stored in top changes, and the size decreases by one.



about an idea, it would be good practice to start developing those skills. While these skills are not much fun at parties, they are very helpful for programming.

### 16.2.5 Operation: push()

The `push()` operation adds a new data element to the top of the queue. In this implementation, a new node object is created to store the new value. The new node will always be placed at the top of the stack, which is the front of the chain. A new node is created for the new value, and this node has a next reference to the node that used to be the front of the stack's node chain. A reference to the new node itself is copied to the stack's `top` attribute. The size of the stack is increased by one, as well.

```
def push(self, value):
    """
    Purpose
        adds the given data value to the stack
    Pre-conditions:
        value: data to be added
    Post-condition:
        the value is added to the stack
    Return:
        (none)
    """
    new_node = N.node(value, self.__top)
    self.__top = new_node
    self.__size += 1
```

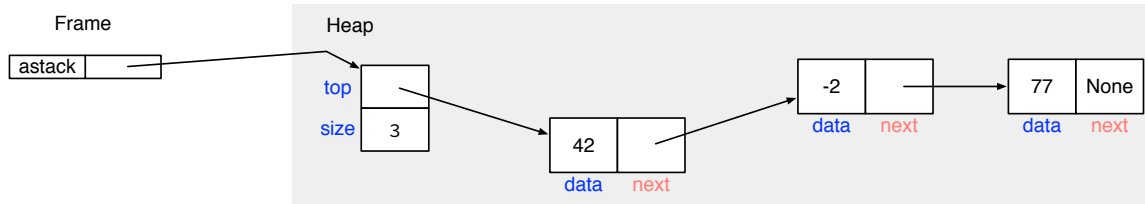
For example, if `push()` were called to add the value 42 to the Stack drawn in Figure 16.2, the result would be as shown in Figure 16.3. Notice that the reference stored in `top` now indicates the node containing 42, and the new node now refers to the node that used to be at the top.

## 16.3 Using Nodes to Implement Queues

The implementation in this section is very similar to the implementation of node-based stacks. Most of the differences will arise because a Queue has a front and a back, both of which are important to the operations. The implementation of Stack only had to manage one end of the node-chain.

The essential nature of a Queue is the FIFO protocol. A node-based queue will use a node chain to contain the values in the queue, and to maintain the sequential ordering essential to the

Figure 16.3: The effects of a push operation on a non-empty stack. A new node is added to the chain, and *top* changes. The size increases by one.



FIFO protocol. The *front* of the queue will be the first node in the node chain, and the *back* of the queue will be the last node in the node chain. We will use the Node ADT to implement the Queue operations, but from the outside, only the Queue operations will be visible, and the reliance on the Node ADT will not be evident.

The attributes we will be keeping track of are as follows:

**front** A reference to the first node in the node chain that stores the queued values; the front of the queue. If the queue is empty, this should be the Python value `None`. This attribute acts as the anchor for the queue's node chain.

**back** A reference to the last node in the node chain that stores the queued values; the back of the queue. If the queue is empty, this should be the Python value `None`.

**size** An integer representing the number of elements currently stored in the queue. This value should be non-negative. A size of zero implies that the queue has no nodes, and no data.

The following sections describe the implementation in a little more detail. The complete implementation can be found in [Appendix D](#).

### 16.3.1 Operation: `Queue()`

The `Queue()` function returns a new object representing an empty queue. This operation simply sets all the attributes to the values that indicate an empty queue.

```
class Queue(object):

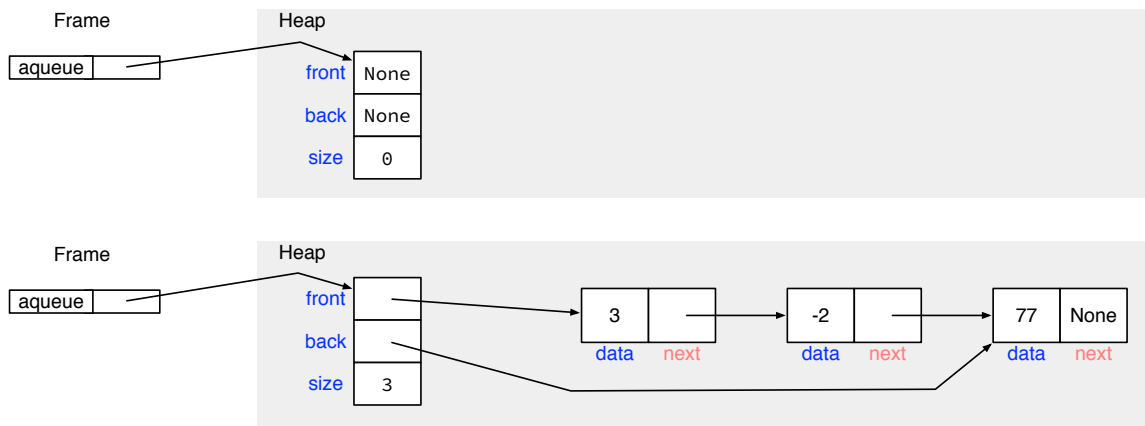
    def __init__(self):
        """
        Purpose
            creates an empty queue
        """
        self.__size = 0          # how many elements in the queue
        self.__front = None      # the node chain starts here
        self.__back = None       # the node chain ends here
```

For example, using the `Queue()` operation in a script such as

```
import Queue as Q

aqueue = Q.Queue()
```

Figure 16.4: The top diagram shows an empty Queue, created on the heap by the Queue() operation. The bottom diagram shows what a queue would look like after several values have been enqueued.



can be visualized as in Figure 16.4. The diagram shows the attribute names in coloured text beside the rectangles. To keep our diagrams as simple as possible, the values are represented as being stored within the queue, but in reality, the object for None and 0 are in Python's memory, and the queue stores references to them within the queue. Eliminating one set of arrows helps to focus on the important aspects of the implementation.

A Queue containing a few values is also shown in Figure 16.4 (the bottom diagram). This demonstrates the kind of consistency that all Queue operations need to maintain. The attribute `front` refers to the first node in the node-chain; the attribute `back` refers to the last node in the same node-chain. The attribute `size` always stores the number of nodes in the node-chain. It's very important to understand that, because we're using references, the attributes `front` and `back` refer to different parts of a single node chain.

### 16.3.2 Operations: `size()` and `is_empty()`

The `is_empty()` operation returns True if the queue is empty. The `size()` operation returns the number of elements stored in the queue. Both implementations below simply use the value of the `size` attribute in the object. This is much faster than counting! This implies that the operations for enqueueing and dequeueing data have to change the `size` consistently.

```
def size(self):
    """
    Purpose
        returns the number of data values in the queue
    Pre-conditions:
        queue: a queue created by create()
    Return:
        The number of data values in the queue
    """
    return self.__size
```

```

def is_empty(self):
    """
    Purpose
        checks if the given queue has no data in it
    Return:
        True if the queue has no data, or false otherwise
    """
    return self.__size == 0

```

These two operations should look familiar, as they are identical to the Stack operations of the same name.

### 16.3.3 Operation: peek()

The peek() operation returns a reference to the data value currently at the front of the queue, without removing it from the queue. If the given queue is empty, the operation causes a run-time error, because the assertion will fail.

```

1  def peek(self):
2      """
3      Purpose
4          returns the value from the front of queue
5          without removing it
6          Note: the queue cannot be empty!
7      Post-condition:
8          None
9      Return:
10         the value at the front of the queue
11     """
12     assert not self.is_empty(), 'peeked into an empty queue'
13
14     first_node = self.__front
15     result = first_node.get_data()
16     return result

```

### 16.3.4 Operation: dequeue()

The dequeue() operation removes a data value from the front of the queue. In this implementation, the value is removed from the front of the node chain stored in the Queue record. The dequeue operation has three situations to be handled:

1. If the queue is currently empty, nothing can be dequeued. In this implementation, a run-time error is invoked.
2. If the queue has exactly 1 value stored, dequeuing a value would result in an empty queue. This situation requires special care. If the last value is dequeued from the queue, then the assignment statement should set the attribute front to the value None. However, when this value is dequeued, the attribute back should also refer to the value None, but that won't happen

automatically. If we didn't check this situation, and deliberately assign `self.__back = None`, the front of the queue would think the queue is empty, but the back of the queue would think there is still something in the queue. This is an example of the kind of detail that causes a lot of bugs. Human brains can be trained to try to watch for these kinds of errors, but we make mistakes when we don't try to think of every possible thing that needs to be done.

This situation is easy to figure out if you remember that we're using two attributes to manage the node-chain. It's like using two hands to manage a string of balloons. When we let go (dequeue) each balloon, one at a time, we have to remember to grab hold of the new front of the string before we let go of the dequeued balloon. When we dequeue the last balloon, we have to remember that we were using two hands for this one balloon, and we have to let go with both hands. If we let go with one hand, but not the other, we haven't really done the job right.

3. If the queue has more than one value stored, then only the first node in the chain is affected; the information about the last node in the chain does not change.

To remove the first node in the chain, we have to save a reference to the data value stored there, and then change the queue's `front` attribute so that it refers to the node after the first node in the chain. By changing the reference of the first node in the chain, we are omitting the node that was previously the front in the chain. This is enough to remove the value and the node from the queue

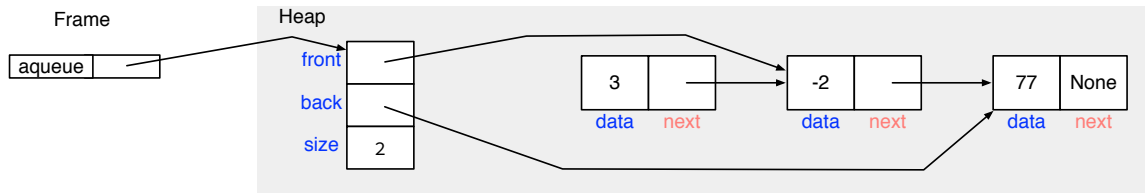
In both cases where a value can be dequeued, we can store a temporary reference to the node to be removed, and temporarily store the value to be returned. We can also change the queue's `front` attribute, by making a copy of the node's next reference. Note that if the node to be removed is the very last one in the chain, its next will be Python value `None`, and that is the proper value for the first node of an empty chain. We can also decrease the size of the queue.

```
def dequeue(self):
    """
    Purpose
        removes and returns a data value from the queue
        Note: the queue cannot be empty!
    Post-condition:
        the first value is removed from the queue
    Return:
        the first value in the queue, or None
    """
    assert not self.is_empty(), 'dequeued an empty queue'

    prev_first_node = self.__front
    result = prev_first_node.get_data()
    self.__front = prev_first_node.get_next()
    self.__size -= 1
    if self.__size == 0:
        self.__back = None
    return result
```

For example, if `dequeue()` were called on the Queue drawn in the bottom half of Figure 16.4,

Figure 16.5: The effects of a dequeue operation on a non-empty queue. The reference stored in `front` changes, and the size decreases by one.



the result would as shown in Figure 16.5. Notice that the reference stored in `front` now indicates the node containing the value -2. The diagram shows that the node containing the value 3 still present in the heap, but nothing in the node chain refers to it, so it is not part of the Queue. If there is no way to access this node from a frame, Python will reclaim the memory it uses for some other purpose at some time in the future.

### 16.3.5 Operation: `enqueue()`

The `enqueue()` operation adds a new data value to the back of the queue. In this implementation, a new node is created to store the new value. The new node will always be placed at the back of the queue, which is the end of the chain. There are two cases that require special care.

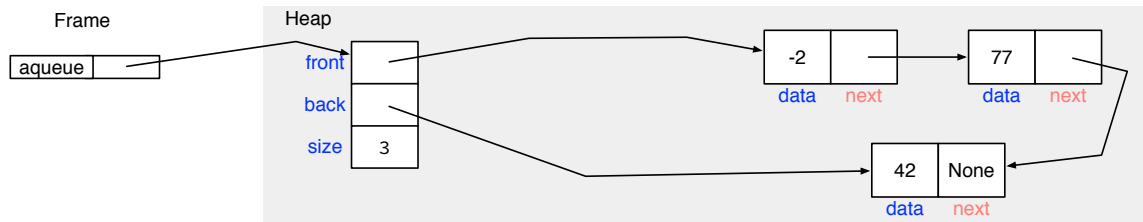
1. If the queue is currently empty, the new value will be stored in a new node that is both the first node and the last node in the chain. For that reason, we have to change the Queue record so that `front` and `back` both refer to the same new node.
2. If the queue is not currently empty, then only the last node in the chain is directly affected; the information about the first node in the chain does not change.

To add a new node at the end of the chain, we can use the reference to `back` directly. To prevent too much confusion, let's use the term *previously last node* to refer to the node that was at the back of the chain prior to a new value being enqueued. When the new node is attached properly, the previously last node should contain a reference to the new node, and the queue's `back` attribute should refer to the new node.

In both cases, the size of the queue increases by 1.



Figure 16.6: The effects of an enqueue operation on a non-empty queue. The new node is connected to the end of the chain, and the reference stored in back changes. The value stored in size is increased by one.



```
def enqueue(self, value):
    """
    Purpose
        adds the given data value to the queue
    Pre-conditions:
        value: data to be added
    Post-condition:
        the value is added to the queue
    Return:
        (none)
    """
    new_node = N.node(value, None)

    if self.is_empty():
        self.__front = new_node
        self.__back = new_node
    else:
        prev_last_node = self.__back
        prev_last_node.set_next(new_node)
        self.__back = new_node

    self.__size += 1
```

For example, if `enqueue()` were called on the Queue to add the value 42 to the Queue drawn in Figure 16.5, the result would be as shown in Figure 16.6. Notice that the reference stored in `back` now indicates the node containing 42, and the node that used to be at the back now refers to the new node as well.

## 16.4 Summary

The similarity between the queue and stack implementations is very strong. The stack was a bit simpler, because we only had to keep track of one end of the chain.

One of the most important lessons to learn in this chapter is the idea of special cases. When a textbook tells you that “the dequeue operation has three cases to handle,” it might seem obvious. But coming up with those cases on your own when you’re solving any problem is not always so easy. Humans typically focus on a typical case, for example, the case where the node chain has a few nodes in it; some people call this kind of situation a *sunny day* case. Some people will also notice the case of an empty node chain, but hardly any novice programmer will realize that the case of a node chain with exactly one node in it is also a special case for the dequeue operation; these are sometimes called *edge* cases. Nothing about Python (or any language) tells us about what the special cases are, or how to handle them; the only way to discover them is to understand your problem so thoroughly that you’ve figured them out. It is essential that you get into the habit of trying to figure out all the cases for your particular problem. The only way to find special cases is to keep asking “What if . . . ?” This is not a good for normal human conversation, because it gets annoying, but it’s essential for good programming. For many kinds of computational problems, drawing diagrams and working things out on paper (or a white board) is really the best practice, before you sit down at a computer to start programming.

## 17 — Linked Lists

### Learning Objectives

After studying this chapter, a student should be able to:

- Employ data abstraction to implement a list.
- Describe the node-based implementation of lists.

### 17.1 Linked List ADT

We started the course by reviewing lists, suggesting that they are very useful tools that help us solve many programming tasks. But we understand Python's lists as an ADT; that is, we have an understanding of their purpose, and we've used many of their operations, but we have not seen their implementation.

In this chapter, we will see how we can implement an ADT very similar to Python's lists using the Node ADT that we presented in Chapter 13; this kind of implementation is known as a linked-list. Our ADT will be based strongly on the implementation of node-based queues in Chapter 16. We'll describe the ADT operations in some detail, but the implementation will be left as a homework assignment.

The term *linked-list* is common in computer science. Frequently, it is used to refer to the structure we called a *node-chain*. For pedagogical purposes, we decided to help avoid confusion by introducing the term *node-chain*, which is helpfully descriptive. This allows us to use the term *linked-list* to name the ADT that is the topic of this chapter. For us, the Linked List is at the same level of abstraction as the Queue and Stack ADTs of the previous chapter. Like the Queue and Stack ADTs, the Linked List ADT will be a record encapsulating a node-chain. Unlike the Queue and Stack ADTs, which provide only a few operations, a Linked List will be far more flexible, providing operations not needed by Stacks and Queues.

The exercise of building a linked list ADT from the node ADT has several purposes. First of all,

it's very good programming experience in general. It is important to point out that many languages use references like Python provides, and the practice allows us to transfer that knowledge to those languages. Second, the practice of using references is necessary, since we will need them again in Chapter 20, for something a little more complicated. Believe it or not, there are more interesting and complicated data structures than stacks, queues and lists that need this concept.

The linked list ADT we develop in this chapter is not how Python lists are implemented. Python lists are even more advanced.

## 17.2 Linked List Record

Like the Stack and the Queue, a Linked List is a record that refers to a single node-chain. The first node in the node-chain is called the *head* of the list, and the last node in the node-chain is called the *tail* of the list. Be careful to understand that there is exactly one chain, but we have two separate references to nodes on the chain.

**head** A reference to the first node in the node chain that stores the queued values. If the queue is empty, this should be the Python value `None`. This field acts as the anchor for the queue's node chain.

**tail** A reference to the last node in the node chain that stores the queued values. If the queue is empty, this should be the Python value `None`.

**size** An integer representing the number of elements currently stored in the queue. This value should be non-negative. A size of zero implies that the queue has no nodes, and no data.

## 17.3 Linked List Operations

Here we briefly outline a few of the operations.

### 17.3.1 Operation: `create()`

The `create()` function returns a new empty linked list record. This operation simply sets all the fields to the values that indicate an empty list.

```
class LList(object):

    def __init__(self):
        """
        Purpose
            creates an empty list
        """
        self.__size = 0      # how many elements in the list
        self.__head = None  # node chain starts here
        self.__tail = None  # and ends here; initially empty
```

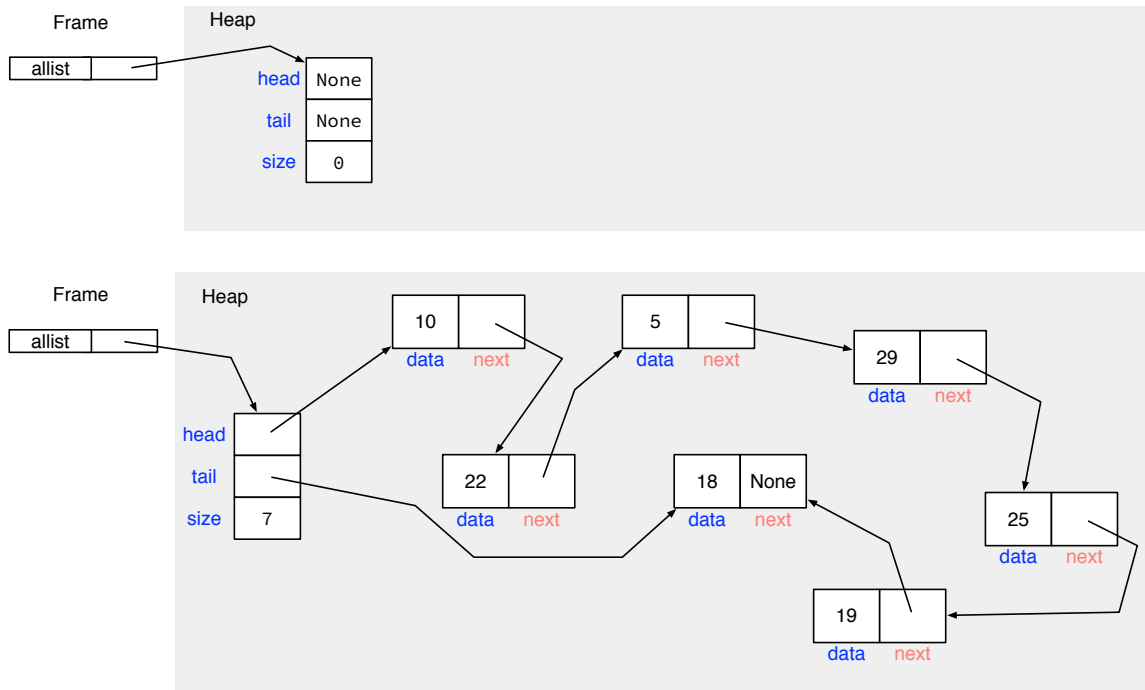
This is very similar to the node-based queue in Chapter 16.

For example, using the `LList()` operation in a script such as

```
import LList as L

alllist = L.LList()
```

Figure 17.1: The top diagram shows an empty List, created on the heap by the `create()` operation. The bottom diagram shows what a linked list would look like after several values have been added to it.



can be visualized as in Figure 17.1. The diagram is simplified a bit, as the Linked List data structure is in reality a dictionary, but the diagram shows the keys in blue, and the values in rectangles. As well, the values are represented as being stored within the Linked List, but in reality, references to None and 0 are stored in the Stack.

### 17.3.2 Operations: `size()` and `is_empty()`

Because queues, stacks and linked lists are all based on the Node ADT, the operations `size()` and `is_empty()` are basically all the same, and therefore do not need to be presented here.

### 17.3.3 Summary of Linked List Operations

In this section, the Linked List operations are described by their purpose. In the descriptions, no mention of the node-chain is made. However, the operations must be implemented using the node-chain stored in the linked list record, in the same way that the Queue and Stack ADT operations were implemented using node chains.

Several operations use the concept of an index. As with Python lists, the index 0 refers to the first value in the list, and the last valid index for a list of size  $n$  is  $n - 1$ .

**`add_to_front(val)`** Insert the given value `val` into the given Linked List `alist` so that the new value is at the front of the sequence of values.

- remove\_from\_front()** Removes and returns the first value in the given Linked List `alist`.
- add\_to\_back(val)** Add the given value `val` to the given Linked List `alist` so that the new value is at the end of the sequence of values.
- remove\_from\_back()** Removes and returns the last value in the given Linked List `alist`.
- get\_data\_at\_index(idx)** Return the value stored in `alist` at the index `idx`. This function does not change the sequence; it simply reports what the value is stored at the given index.
- set\_data\_at\_index(idx, val)** Store `val` into `alist` at the index `idx`. This operation does not change the structure of the list. It simply replaces the value currently stored at `idx` with the given value.
- value\_is\_in(val)** Check if the given value `val` is in the given list `alist`. Returns `True` if `val` is anywhere in the sequence, and `False` otherwise.
- get\_index\_of\_value(val)** Report the index of the given value `val` in the given list `alist`. If `val` appears more than once, the index of the first occurrence is reported. This function returns the tuple `(True, i)` if the given value appears in the list, where `i` is the index of `val`. If the value is not in the list, this function returns `False, None`.
- insert\_value\_at\_index(val, idx)** Insert `val` into `alist` at index `idx`. This operation changes the structure of the list by adding a new value into the sequence, provided that `idx` is a valid index. For simplicity, we'll assume the index is non-negative, and in the range 0 to  $n$ , where  $n$  is the length of the list. If the index given is equal to the size of the list, the new value is added to the end of the sequence.
- delete\_item\_at\_index(idx)** Delete the value at index `idx` in the given list `alist`. This operation changes the structure of the list, by removing a value. For simplicity, we'll assume that a valid index is non-negative, and in the range 0 to  $n - 1$ , where  $n$  is the length of the list.
- delete\_value(val)** Delete the value `val` from the given list `alist`. This operation changes the structure of the list, by removing a value. If the given value appears more than once, on the first occurrence is removed. If the given value does not appear in the sequence, the list remains unchanged. The function returns `True` if a value was deleted, or `False` if not.

#### 17.3.4 Implementation details

The implementation of many of the above operations is similar to Stack and Queue operations given in Chapter 16. The operations that are different usually involve an index argument. Unlike Python lists, which can be indexed directly using normal list indexing syntax, e.g., `alist[2]`, the Linked List ADT cannot use indexing syntax. Because the values are stored in a node-chain, index for a Linked list operation tells you how far along the node chain you have to step to find the value indicated by the index.

The operations for indexing require the use of while loops. Python's for-loops are useful for Python lists, but not so much for node-chains.

#### 17.3.5 Perspective

The purpose of this chapter is to reveal one important technique as the basis for the list data structure given to Python programmers in the core Python language.

Functional programming languages, e.g., LISP, Scheme, SML, Haskell, use node chains as their essential compound data structure. Object oriented languages such as C++ and Java have standard modules for Linked lists very similar to what we have seen here. More primitive languages, such as C, do not have lists at all, and one of the "rites of passage" for C programmers was building their

own Linked List ADT.

While Python lists are not currently linked lists the way we've seen here, they probably were linked lists like this very early in Python's development as a language. As more people used Python, greater importance was placed on making list operations as efficient as possible, so a more advanced data structure was introduced. It is very important to realize that a better implementation of the list operations is hidden behind the ADT interface, and applications already written using Python did not need to be rewritten.

There are lots of operations we could have added to the Linked List ADT. For example, in Python we can *slice* a list, e.g. `alist[2:7:2]`. We could write a Linked List operation that does the same thing. The point is that every Python list method is implemented as a function that is in some sense similar to the Linked List operations in this chapter.





### The importance of understanding costs

Abstraction #1: Time  
Abstraction #2: Input size  
Abstraction #3: Categories  
Abstraction #4: Worst case behaviour

### Asymptotic analysis

Basics of Working with Big-O  
Counting Steps in Algorithms

### Examples

Step analysis  
Method 2: Simplified Steps  
Logarithmic loops  
Nested loops with independent limits  
Nested loops with dependent limits

### Warnings and Advice

## 18 — Algorithm Analysis

### Learning Objectives

After studying this chapter, a student should be able to:

- Explain the value of algorithm analysis.
- Describe the use of abstraction to quantify time and space in analysis.
- Explain the value and purpose of best-case, worst-case, and expected-case analyses.
- Employ best- and worst-case analyses for time and space for algorithms involving loops, conditionals and function calls.
- Explain the relationships between basic polynomial complexity classes.
- Explain what is meant by “best” and “worst” case behavior of an algorithm. [Partial, AL-BA-1]
- Determine informally the time and space complexity of simple algorithms. [Usage, AL-BA-3]

### 18.1 The importance of understanding costs

When you’re using a computer to solve a problem, you should always be aware of the costs of the solution. In this chapter, we’ll focus on the costs associated with running a program: namely, the amount of time and memory required by the program while it is running. Concerns about time and memory costs really only come into significance when the programs get larger, when datasets get larger, and when the problems get more difficult to solve.

Being able to assess costs is essential in making design choices: if you have more than one algorithm to solve a given problem, you will want to choose the one that has lower costs. Computers are fast, but they are not infinitely fast. The programs we write in a first-year course may run in milliseconds, but that’s not something we should take for granted. The software applications, and Internet services, we use are fast, precisely because designers understood the costs, and figured out

ways to reduce them. A lot of the speed we expect comes from the hardware, but even miraculous hardware cannot always save designers from poor software design choices.

We'll hone our analysis skills on Python programs, but it's crucial to understand that these skills are transferable to any programming language. In fact, we will even be able to apply these skills to algorithms expressed in more abstract ways, like pseudocode and flowcharts. The key point here is that we want to understand the costs of a program while it's still in the design phase, without implementing it first. If we can assess the cost of an algorithm, the analysis will apply to any of the programs we might come up with as an implementation of the algorithm.

Our analysis methods will identify the key costs, and express them in a way that is independent of the speed of any particular computer, independent of the programming language, and independent of minor implementation differences. It's such an important skill that we have to start practicing the basics sooner than later. In fact, this chapter will focus almost entirely on time costs. Once you get the hang of the analyses here, you can easily use the same skills to measure memory costs, or network costs, or almost any other computational resource, using a very similar approach.

### 18.1.1 Abstraction #1: Time

It may sound odd, but we will not measure time costs in typical time units like seconds, minutes, or hours. To understand why, consider a single program executed on two different computers, say one from 1967, and another from 2017. It's obvious that the program would be doing the same calculations on both computers, and the computers would be doing the calculations at different speeds. We want our analysis to say more about the algorithm than the computer it's running on.

We will measure time indirectly, in terms of computational steps. For us, a step will be any computational action that cannot be broken down into other steps. Expressing the costs in terms of the number of steps allows us to assess the algorithm, independent of the computer running the program. We will say more about what counts as a computational step a bit later, but it will include such things as arithmetic operations, assignment statements, comparisons, etc. There is nothing surprising waiting for us.

We assume that a step will take a constant amount of time. This is not always true, but it's very useful in many cases. One reason is that computer programs typically use relatively small numbers. Any floating point number that can be stored on the heap counts as relatively small. Any integer that is likely to come up in normal calculations counts as relatively small.<sup>1</sup>

We also assume that all the different kinds of steps take roughly the same time on a given computer to do. This is also a bit of an approximation. For example, we are suggesting that adding two numbers is a step, and multiplying two numbers is a step, but multiplication requires addition (think about the algorithm to multiply  $2345 \times 6845$ ), so multiplication is obviously a bigger step than addition. For many common purposes, we are content to say that addition and multiplication both count as the same size step because the numbers involved are relatively small, and the time it takes to multiply is so small that time differences between adding and multiplying doesn't really matter. The key assertion here is this: We don't (usually) lose a lot of decision-making power by assuming each step is performed in about the same amount of time.

---

<sup>1</sup>In some languages, the range of numbers that can be stored using the basic language constructs has precise limits. In Python 3, the floating point numbers are limited this way, but integers are not limited at all. So strictly speaking, integers in Python 3 are not always considered small in this way. But normal calculations usually don't involve numbers so big that they represent extraordinary costs.

### 18.1.2 Abstraction #2: Input size

Consider a program that analyzes some data. It is reasonable to assume that the program would need a larger number of steps to analyze more data. Our accounting for steps has to be flexible enough to allow analysis no matter what amount of data is used.

To do this, we employ a technique called “parameterization.” The word “parameterization” means that we give a mathematical name to a quantity that is either unknown in advance, or may change depending on things beyond our control.

Using parameterization, we will assign a mathematical name to the size of the problem; typically, we use  $N$  or  $n$ . For example, we might want to sort a list of size  $N$ . Or we might want to add numbers having  $N$  digits. Or maybe we need storage for a video whose length is  $N$  seconds. There are also situations where the size of a problem may require multiple parameters, for example, applying an algorithm to a spreadsheet with  $N$  rows and  $M$  columns.

Our analysis will express the costs of a program in terms of the parameters describing the size of the input. In this way, we have an analysis that does not depend on any given example data set. For example, an algorithm to find the smallest value in a list of  $N$  values might require  $3N + 4$  steps. The cost is being expressed as a function of the size of the input, and the input size is being expressed as a parameter,  $N$ . We’ll see how to do this analysis; right now, we need to understand the idea of input size and parameterization.

### 18.1.3 Abstraction #3: Categories

In biology, scientists study plants and animals, assess their similarities and differences, and use a complicated organizational scheme to express the relationships between species. In computer science, we also study, evaluate, and categorize algorithms. Like biology, our categorization scheme has to be useful for algorithms we already know, and for any that might be invented in the future.

Our categorization scheme will distinguish significant differences, and ignore minor differences. Suppose you are trying to decide between Algorithm A and Algorithm B, both of which solve the same problem. Suppose you determine that Algorithm A takes  $3N + 3$  steps, and Algorithm B takes  $3N + 4$  steps. Obviously, Algorithm A is one step better, but you are on the right track if you are wondering whether a difference of one step is actually significant.

The categorization system will be based on two ideas:

1. Each category will contain a set of algorithms.
2. Each category will be labelled with a simple mathematical function that describes the algorithms in the set.

We’ll say a bit more about the meaning of the categories a bit later.

There are as many categories in our system as there are different kinds of functions, but we’ll start with the easy ones. The categories will be labelled as functions expressed in terms of input size. For example, we’ll have a category labelled as linear,  $N$ , as well as quadratic,  $N^2$ . Any constant power of  $N$  is a category, too, e.g.,  $N^{15}$ . We have categories based on exponential and logarithmic functions, e.g.,  $2^N$  and  $\log_2 N$ . We even have a category for functions that don’t depend on a size parameter  $N$  at all (this is the category labelled with the value 1).

The formal notation for our categories uses something called “Big-O” notation. For example, we

have the category  $O(N)$  and the category  $O(N \log N)$ . A good start to the categories is the following:

$$\begin{aligned} O(1) &< O(\log N) < O(N) < O(N \log N) \\ &< O(N^2) < O(N^3) < \dots < O(N^k) \\ &< O(2^N) < O(N!) < O(N^N) \end{aligned}$$

We're using the symbol  $<$  to mean "increases more slowly than" in the ordering. The order encodes mathematical facts like this: the function  $y = x^2$  increases faster than  $y = \log_2 x$ . We'll sometimes say " $O(N)$  is lower order than  $O(N^2)$ ," which simply means that  $O(N)$  appears to the left of  $O(N^2)$ . The ranking above uses  $k$  to represent any constant expression relative to  $N$ . Finally, don't be misled by the fact that the categories appear on several lines; the whole system doesn't fit on one line!

Since we typically prefer the most efficient algorithms, we want to use ones that fall into categories as far left on the spectrum as possible. The closer to  $O(1)$ , the more efficient the algorithm. The ordering we provided is not complete though. We left out  $O(N^2 \log N)$ , for example, which comes between  $O(N^2)$  and  $O(N^3)$ . You can put any function of  $N$  into this list at the right place.

#### 18.1.4 Abstraction #4: Worst case behaviour

When assessing the cost of an algorithm, there are algorithms that always do the same amount of work, no matter what. For example, taking an average: no matter what, taking an average requires looking at all the data. But there are algorithms whose costs depend on the size, but also the arrangement of the data, or some other factor that does not depend on size.

For example, if we are using linear search to check if a value appears in a list, the program will take fewer steps if the value is near the front of the list, and a lot more if the value is near the end of the list. If the value is the first item in the list, the program has to do very little work. If the value does not appear in the list at all, or if the value is very last in the list, the program has to look at all the values in the list.

We say that the worst case for an algorithm is an arrangement of data that forces the algorithm to do the maximum amount of work for a given problem size. Likewise, the best case allows the algorithm to do the minimum amount of work for a given problem size. Be careful not to misunderstand the role of problem size in best and worst cases. Problem size is independent of best and worst cases. To say that the best case for linear search is when the list is empty is completely wrong!

It would be unrealistic to use the best case as being representative of the possible range of costs for linear search. On the other hand, we often use the worst-case as representative of the costs. This identifies a certain amount of pessimism in our analysis, but when there are factors beyond your control, it is better to be prepared for the worst possible outcome, than to plan for only the best outcome. We are often interested in best and worst cases. Sometimes they describe the range of possible behaviours. Sometimes, and especially in asymptotic analysis, we find the best and worst cases fall into the same category, which reassures us that the categorization is informative.

One naturally thinks of using an "average case" analysis to make a reasonable compromise between the best and worst cases. When this is actually possible, this is a good idea. However, average case analyses are notoriously difficult, because of the complexity of the things we ask a computer to do.

## 18.2 Asymptotic analysis

To put algorithms into categories, we'll be applying a principle called "asymptotic analysis." This means that we'll infer what happens as the problem size gets really big. It will allow us to ignore any effect that is relatively small.

For example, in the previous example with Algorithms A and B, we know that they are different by one step, no matter what  $N$  is, but that one step becomes less and less significant as  $N$  increases. Asymptotic analysis allows us to put both algorithms into the same category, namely  $O(N)$ . We say "Asymptotically, Algorithm A and Algorithm B have the same runtime costs." While the term "asymptotic analysis" may sound deep and imposing, applying it boils down to eliminating details that can be ignored; in other words, it's just abstraction with a fancy name! The name serves to guide us in what details can be thrown away.

It's important to clarify what the asymptotic categorization means. Each category is labelled by a function, like  $N^2$ . If we take any algorithm in the  $O(N^2)$  category, and plot a curve of the number of steps it uses as a function of  $N$ , the graph would never increase faster than the  $N^2$  shape, especially when you get to very large values of  $N$ . More precisely, for any algorithm  $A$  in  $O(N^2)$ , we could always find a positive number,  $k$ , so that the graph of  $kN^2$  is higher than the graph of the costs of  $A$ . An algorithm in  $O(N^2)$  cannot be pushed into another category by increasing or decreasing the costs by a factor of 2 or 10 or 10 million. An algorithm cannot be pushed into another category by adding or subtracting 10 steps or 10 million steps. It's the basic shape of the function that labels the category that is important.

Technically, there is a formal mathematical meaning to our categorization system, which describes the categories as upper bounds on costs. These technical details are important, but we will leave the technical presentation for future courses.

### 18.2.1 Basics of Working with Big-O

Remember our approach is to count steps that a computer would perform carrying out an algorithm. We'll use parameterization to express the number of steps in terms of the problem size. We might end up with an expression like  $3N^2 + 4N + 18$ , which describes the number of steps carried out when the problem size is  $N$ . We'll talk about how to count steps in the next section.

There are a number of simple rules that you can use to perform asymptotic analysis using Big-O notation. These rules are all derivable from first principles, but since I haven't shown you those, we'll simply state the rules.

**Principle 18.1** In asymptotic analysis, the following identities hold for any expressions  $A$  and  $B$ :

$$O(A) + O(B) = O(A + B)$$

$$O(A) \times O(B) = O(A \times B)$$

$$A \times O(B) = O(A \times B)$$

For example, we can use the first identity to simplify an expression as follows. Starting with the previously mentioned example:

$$O(3N^2 + 4N + 18) = O(3N^2) + O(4N) + O(18)$$

**Principle 18.2** In asymptotic analysis:

- If  $A$  is a constant,  $O(A) = O(1)$ .
- If  $A$  is a constant, and  $B$  is any expression,  $O(A \times B) = O(B)$ .

Continuing the example, the three terms we have now are  $O(3N^2)$ ,  $O(4N)$ ,  $O(18)$ . The constant factors are 3, 4, and 18. If we apply the rule, we get the following:

$$O(3N^2) + O(4N) + O(18) = O(N^2) + O(N) + O(1)$$

This rule does not apply to exponents, which are constants, but not coefficients.

**Principle 18.3** In asymptotic analysis, if  $O(A) < O(B)$ , then  $O(B) + O(A) = O(B)$

According to our hierarchy, we know that  $O(1)$  and  $O(N)$  are lower order than  $O(N^2)$ , so they can be dropped. We keep the highest order term only. Therefore, we can write:

$$O(N^2) + O(N) + O(1) = O(N^2)$$

Thus, asymptotic analysis tells us that  $3N^2 + 4N + 18$  is  $O(N^2)$ .

We applied the rules in a specific order, but it turns out that the order you use doesn't matter. We could have dropped lower order terms first, then replaced the constant factors.

## 18.2.2 Counting Steps in Algorithms

At last, we get to the business of the actual analysis, after the long prelude of abstractions. We'll start with a definition of the concept of a computational step. A single step is any computational action that cannot be broken down into other steps; they are the smallest computational actions that we can count by looking at the code. Some authors call these *primitive operations*.

### Definition 18.1 Step

Any one of the following can be considered as a single step:

- A single arithmetic operation ( $\times$ ,  $+$ ,  $-$ ,  $/$ ,  $//$ ,  $\%$ ) provided the numbers are not unreasonably large
- A single Boolean operation (**and**, **or**, **not**)
- Assigning a value to a variable
- Comparing two values using a relational operator such as  $==$ ,  $<$ , etc.
- Indexing a list element
- Accessing a dictionary element by its key.
- A function call (the cost of executing the function is separate)
- Returning a value from a function

Every simple statement in Python (or any other language for that matter) consists of a number of steps. For example, consider the following statements:

```
a = 3*y+4
b = data[2]*data[3]
c = 3*math.sqrt(data[i+1])
```



The first statement has 3 steps; the second has 4 steps, and the third statement requires 4 steps, not counting the number of steps that `sqrt` does. The whole sequence of statements above requires 11 steps (not counting `sqrt`).

Time for another principle stating the obvious:

**Principle 18.4 Sequence**

In asymptotic analysis, the total number of steps for a sequence of statements is the sum of the steps for each individual statement. ■

For a conditional statement, it would be a mistake to add the costs of all branches together to get a “total”.

**Principle 18.5 Conditionals**

In asymptotic analysis, analyze each block in a conditional independently. Add the cost of the conditions to each relevant branch. ■

The best-case uses the branch with the smallest total cost; the worst case uses the branch with the largest total cost. If the conditional has no else block, the best case is an empty block that has zero cost (which is the smallest any cost can get).

The principle of sequential statements can be generalized to a statement about loops:

**Principle 18.6 Loops**

In asymptotic analysis, the total cost of a loop is obtained by adding up the costs from every repetition, and including costs involved with controlling the loop. ■

There is an easy case, and a more complicated case. If the cost of the body of the loop never changes over the repetitions, you can compute the total by multiplication: the number of times the loop body is executed, and the cost of executing the loop body once.

Here’s an example of the easy case. Suppose you eat three apples per day for  $N$  days. The number of apples you eat per day doesn’t change over time. So three apples per day times  $N$  days is  $3 + 3 + \dots + 3 = 3(1 + 1 + \dots + 1) = 3N$  apples. Here, apples represent steps, and the number of days represent repetition.

The more difficult case applies when the cost of the body of the loop changes depending on where you are in the total sequence. If the cost of the body changes, you can’t always use the short cut of multiplication. You have to figure out how to calculate the total using addition.

Here’s a simple example of the more difficult case. Suppose you eat one apple on the first day, two apples on the second day, three apples on the third day,  $\dots$ , and  $N$  apples on the  $N$ th day. The number of apples eaten on a particular day depends on what day it is. The total number of apples is  $1 + 2 + 3 + 4 + \dots + N = N(N + 1)/2$ . This is a formula that we use a lot in computer science; keep it in mind so that when you see it you’ll remember it. Eventually, you’ll memorize it simply due to repeated use!

Keep in mind that the body of a loop can contain statements, sequences, conditionals and other loops. You have to analyze them separately, first!

**Principle 18.7 Functions**

In asymptotic analysis, the cost of executing a function has to be analyzed separately and added to the cost of any statement containing a function call. ■

Above, we stated that calling a function requires a single step. This step involves sending the data to the function, and starting the function off, but does not include the steps performed by the function itself. The cost of the function itself has to be analyzed separately. That cost is added to the cost of any statement the function call appears in.

## 18.3 Examples

### 18.3.1 Step analysis

Consider the following example.

```
1 total = 0
2 for i in range(10):
3     total = total + 1
```

- Line 1: 1 step
- Just line 3, once: 2 steps
- Just line 2, over-all: 10 assignments, 9 additions, so 19 steps
- Line 3 is repeated 10 times, so 20 steps
- Totals:  $1 + 19 + 20 = 40$  steps
- Asymptotically, 40 is  $O(1)$
- Literally, the problem size here is not a variable.

Consider the following example.

```
1 total = 0
2 for i in range(N):
3     total = total + 1
```

- Line 1: 1 step
- Just line 3, once: 2 steps
- Just line 2, over-all:  $N$  assignments,  $N - 1$  additions, so  $2N - 1$  steps.
- Line 3 is repeated  $N$  times, so  $2N$  steps
- Detailed step analysis:  $1 + 2N - 1 + 2N = 4N$
- Asymptotic category:  $4N$  is  $O(N)$ .

### 18.3.2 Method 2: Simplified Steps

Consider the following example.

```
1 total = 0
2 for i in range(10):
3     total = total + 1
```

- Line 1: 1 step, so  $O(1)$

- Just line 3, once: 2 steps, so  $O(1)$
- Just line 2, over-all: 10 assignments, 9 additions, 19 steps, so  $O(1)$
- Line 3 is repeated 10 times, so  $10 \times O(1) = O(1)$
- Totals:  $O(1) + O(1) + O(1) = O(1 + 1 + 1) = O(1)$
- Literally, the problem size here is not a variable.

Consider the following example.

```

1 total = 0
2 for i in range(N):
3     total = total + 1

```

- Line 1: 1 step, so  $O(1)$
- Just line 3, once: 2 steps, so  $O(1)$
- Just line 2, over-all:  $N$  assignments,  $N - 1$  additions,  $2N - 1$  steps, so  $O(N)$
- Line 3 is repeated  $N$  times, so  $N \times O(1) = O(N)$
- Totals:  $O(1) + O(N) + O(N) = O(2N + 1) = O(N)$

### 18.3.3 Logarithmic loops

Consider the following example.

```

1 total = 0
2 i = 1
3 while (i < N):
4     total = total + 1
5     i = i * 2

```

**Question:** How many times does the loop-body execute?

**Insight:** The values for  $i$  are: 1, 2, 4, 8, ...

**Answer:** We can double  $i$  at most  $\log_2 N$  times before  $i > N$

Lines 1,2,4,5 are all  $O(1)$ . Lines 4,5 are executed  $\log N$  times. The total cost is  $O(\log N)$ .

### 18.3.4 Nested loops with independent limits

- Start with the inner-most loop
- Analyze to find the number of steps
- Abstraction: treat the inner-most loop as a statement with the analyzed cost.
- Work outward, one loop at a time.

```

1 total = 0
2 for i in range(N):
3     for j in range(M):
4         total = total + 1

```

Assume  $M$  and  $N$  are separate input size parameters. Lines 3-4 are  $O(M)$ . Repeating the inner loop  $N$  times gives a total of  $N \times O(M) = O(NM)$ .

### 18.3.5 Nested loops with dependent limits

- Start with the inner-most loop
- Analyze to find the number of steps expressed using dependency
- Abstraction: treat the inner-most loop as a statement with the analyzed cost.
- Sum the total costs of the dependent loop
- Work outward, one loop at a time.

```

1 total = 0
2 for i in range(N):
3     for j in range(i):
4         total = total + 1

```

- In general, lines 3-4 are repeated  $O(i)$  times
  - Total cost: sum the costs for  $i$  from 0 to  $N - 1$
  - Total:  $O(0) + O(1) + O(2) + \dots + O(i) + \dots + O(N - 1)$
- To find an answer for this sum, we need to use the following identity:

$$\sum_{i=1}^n i \equiv 1 + 2 + 3 + 4 + \dots + n = \frac{n(n+1)}{2}$$

This is the same one that came up in the discussion of Principle 10.7.

- Total:  $O(0) + O(1) + O(2) + \dots + O(N - 1) = O(N^2)$

## 18.4 Warnings and Advice

None of this material is difficult once you get the hang of it. Don't let the novelty get in the way. It's just counting and abstraction.

There are other kinds of costs associated with algorithms which we won't discuss in this chapter. Some algorithms use a lot of memory, so that is a cost. Some applications require network communications, and the cost of communication needs to be considered. Assessing these costs obey the same principles as the costs we cover in this chapter.

There are costs that the approach in this chapter does not apply to. The most obvious of these is the development cost: how much time it takes to bring the algorithm from design (i.e., when you know how the algorithm will work) to implementation (having correct, working code implementing the algorithm). There's also the cost of research: if you have a problem, but no algorithm, you need to come up with one somehow! It's not always easy. There are other costs as well.

In a previous section, we assumed that all steps take the same amount of time, notably, addition and multiplication. For very complicated mathematical algorithms, especially those that might take a very long time to run, scientists sometimes determine costs by counting additions and multiplications separately, precisely because they take different amounts of time. But, that's a pretty special case of analysis, which you can easily learn after you master the basics.

Finally, this chapter is an introduction to analysis of algorithms. This is where we have to stop, but it's not the final word, and there are other forms of analysis that can be applied to answer slightly different kinds of questions.

## Introduction

Terminology  
The Delegation Metaphor  
Avoid the rabbit hole

## How to Design a Recursive Function

The recursive case

## Examples

Functions, recursion, and the call stack  
Algorithm analysis for recursive functions

# 19 — Recursion

## Learning Objectives

After studying this chapter, a student should be able to:

- Derive recursive formulations for algorithms.
  - Author, test, and debug recursive Python functions.
  - Explain recursion in terms of the delegation metaphor.
  - Describe how recursion is implemented in terms of functions calls and stack frames.
  - Analyze simple recursive functions for time complexity.
  - Analyze simple recursive functions for space complexity, including the use of the system stack.
  - Determine whether a recursive or iterative solution is most appropriate for a problem.
- [Assessment, SDF-AD-6]

## 19.1 Introduction

Recursion is a form of repetition using function calls instead of loops. The main idea is that a function tells you how to perform a task, and so if you need to perform the task again, you just call the function again. There's a little more to it than that, but that's the basic idea.

Anything that can be done with a loop can also be done with recursion, and vice versa. You may wonder why we bother with recursion, since we already know two kinds of loops. It can be said that a loop is the kind of repetition that derives from a practical engineering design, whereas recursion is the mathematical basis for all repetition. Some algorithms can be written very elegantly and efficiently using recursion, but would be significantly more difficult to write with a loop. There are programming languages, such as LISP, Scheme, Haskell, Prolog, whose entire repertoire of repetition rests on the foundation of recursion. Recursion is not some mathematical curiosity; it's a fundamental aspect of computation.

Throughout this chapter, we will refer to an idea of computing or performing a *task*. A recursive function, like any function, is best designed to do one thing well. In a very real sense, recursive functions are highly focussed on a narrowly defined task. Our examples will start with tasks such as summing a series, or calculating a complex formula. We will talk about performing tasks of different sizes, by which we mean that the task depends on a function parameter that dictates how much work to do. We will sometimes use the words *main task* to refer to the general case of the task we are trying to solve, and the word *subtask* to refer to smaller instance or version of the main task.

### 19.1.1 Terminology

All recursive functions have the same general structure. At first glance, they are functions containing a conditional (if-else statement, or a variation).

- At least one of the branches of the if-statement gives a simple answer to a very simple task. This is called a *base case*. In general, if you are designing a recursive function to perform some task involving repetition, then the base case is the result of the most basic variation of the task you can imagine.
- The other branch (or branches) transforms the current task into one or more sub-tasks, and then combines the results of the subtask(s) to form the result of the main task. This is called the *recursive case*. The recursive case always makes a function call (usually one, but sometimes more) to the function being defined.

Let's look at an example. Here, we are adding up all the squares of integers from 0 to  $N$ , that is, we are calculating  $0^2 + 1^2 + 2^2 + 3^2 + \dots + N^2$ :

```
def sumSquares(N):  
    if (N <= 0): # base case  
        return 0  
    else: # recursive case  
        return (N*N) + sumSquares(N-1)
```

The base case here is characterized by  $N = 0$ . When  $N = 0$ , the sum of squares of integers from 0 to 0 is equal to 0, and that is what the function returns, without doing anything else. The base case for a recursive function is almost always so simple that it requires very little problem solving. Usually, common sense tells us the base case without much thought. When writing your own recursive functions, you might be concerned that the base case seems too easy. Don't be. It usually is that easy. We usually write a conditional that checks for the base case explicitly, as above, but as you know, an if-statement can always be reformulated so that one or the other branches is implicit.<sup>1</sup>

The recursive case is characterized by  $N > 0$ . Here, the recursive case is found in the *else* part of the if-statement. It tells us that we can calculate the required sum by

1. Obtaining the sum of squares from 0 to  $N - 1$  using a recursive call
2. Adding  $N^2$  to the sum computed recursively.

In this function, the main task (computing the sum up to  $N$ ) is solved by first solving a smaller subtask (computing the sum up to  $N - 1$ ) using recursion, and the result is combined using addition.

---

<sup>1</sup>In our base case, we also included  $N < 0$ , which is a practical design, in case someone calls the function with a negative integer. It gives the *wrong answer* for negative input (because there is no right answer for negatives, but the function returns 0).

For many of the simplest kinds of tasks, there is one base case, and one recursive case. For more complicated kinds of tasks, there may be more than one base case, or more than one recursive case. You really need to understand the problem very well to decide if you need multiple base cases or recursive cases. The number of these cases is determined by the problem you are solving, not the recursive function you are designing.

### 19.1.2 The Delegation Metaphor

It is helpful to use a metaphor to help reduce confusion. Think of a function call as delegation. That is, no matter what function is involved, a function call is like calling in an assistant, giving him a task, and some room to work, and waiting for him to come back with a result (or the task complete).

With this metaphor, it doesn't matter if you give your assistant the instructions for the same function you are working on, or if you give your assistant instructions for a different function. The assistant takes only the data you give him, and works completely independently. He may create variables with the same names as ones you created, but they are different variables. He may also call an assistant to help with a subtask.

The delegation metaphor is not just a story that sounds nice and easy. It actually reflects the behaviour of the computer when a function is involved, recursive or otherwise. But we'll have more to say about that towards the end of the chapter.

### 19.1.3 Avoid the rabbit hole

A recursive function need not be confusing. Lots of hay is made about recursion as self-reference, but ignore it if you can. Students are often fooled into believing that they have to understand the whole sequence of function calls that lead to the correct answer. Students whose grasp on the concept of recursion is not solid often refer to the "winding" and "unwinding" of recursion. They think of recursion as doing something "all the way down" and then doing something different "all the way back up." I call this the rabbit hole, and it's not helpful.

Recursion can be understood as a relationship, and humans are good at seeking and discovering these kinds of relationships. For recursion, there are two relationships you need to grasp.

1. The base case. It's typically easy to grasp, and therefore not really worthy of long discussion.
2. For the recursive case, all you need to think about is the relationship between the main task and the immediate subtask. You need to be able to say how the parameters for the main task are related to the parameters for the immediate subtask; this will involve a computation or calculation of some sort. You also need to be able to say how a solution to the subtask is related to the solution of the main task.

That's it. It's not confusing or mysterious. That is not to say that it is easy to change your mind. If all you practice is loops, convincing yourself that recursion is worth learning is the true educational battle.

## 19.2 How to Design a Recursive Function

You can always get a start on defining a recursive program by typing the following template:

```
def <function name>(<data>)  
    if <base case test>:  
        <base case task>
```



```

else:
    Calculate <new params> for smaller, related problem
    <subSolution> = Call <name> recursively on <new params>
    <solution> = Combine <subSolution> with <data>
    return <solution>

```

Then you fill in the blanks. Not all at once, and maybe you will make some revisions as you go. And there will be simplifications you can make when you're done.

Regardless, the design of the recursive function requires the following steps. You don't have to derive them in this order, but you do have to come up with something for all of them:

1. Give your function a good name. Not *flower* or *Yoda* or *function*. The name should reflect the computation it's supposed to do. This step helps you focus. Using a name that does not reflect the purpose of the task makes your design work harder.
2. Decide the input and output for your function (pre-conditions and return value, possibly post-conditions). Your input is the data your function needs to determine what to do, and you need to be able to use this data to distinguish between the base case and the recursive case.
3. Determine the test for the base case. It is very common to have tests like  $N == 0$  or  $N == 1$ . But be careful! The right base case test comes from the task, not from Python. If you don't understand the task, you might not have the right base case test. You need to be able to identify the very simplest instance(s) of the task that your function will perform.
4. Write the code for the task that deals only with the base case. Once you have identified the base case, the code you need for the base case is usually very simple. If your function returns a value, then you need a return statement here.
5. Write the code for the recursive case. We can break it down into steps here:
  - (a) Code that calculates the inputs for the subtask using the input information about the main task.
  - (b) Call the function recursively, using the input for the sub-task. The template stores the result in a variable, but that's to be able to make the steps distinct. Simplification can often eliminate this assignment.
  - (c) Combine the result of the sub-task with something about the main task. This is very often a simple step, like addition, or sequence. The combination has to be based on knowledge of the task.
  - (d) Return the result, if appropriate.

### 19.2.1 The recursive case

By far, the recursive case is the part that requires the most thought. There are two key notions to designing the recursive case:

1. You need to be able to say how the parameters for the main task are related to the parameters for the immediate subtask; this will involve a computation or calculation of some sort.
2. You also need to be able to say how a solution to the subtask is related to the solution of the main task.

In our first example, suppose we gave a recursive function to calculate the sum  $0^2 + 1^2 + 2^2 + 3^2 + \dots + N^2$ . This was our main task. We can do some grouping to form a subtask, namely  $0^2 + 1^2 + 2^2 + 3^2 + \dots + (N-1)^2$ . This is a simpler version of the main task. The subtask's parameter

is one less than the main task's. We can give the two sums each a name,  $\text{sumsq}(N)$  for the main task, and  $\text{sumsq}(N-1)$  for the subtask. It's almost too easy to relate the two tasks. If we start with the equations for  $\text{sum}(N)$  and  $\text{sum}(N-1)$ :

$$\begin{aligned}\text{sumsq}(N) &= 0^2 + 1^2 + 2^2 + 3^2 + \cdots + N^2 \\ \text{sumsq}(N-1) &= 0^2 + 1^2 + 2^2 + 3^2 + \cdots + (N-1)^2 \\ \text{sumsq}(N) &= \text{sumsq}(N-1) + N^2\end{aligned}$$

The third line above is the relationship between the two tasks. Notice how the solution to the subtask, namely  $\text{sumsq}(N-1)$  is combined with the value  $N^2$ , which comes from the main task, using addition.

## 19.3 Examples

### Simple Sum

Consider the task of adding up all the integers from 1 to  $N$ , where  $N$  is any positive number. In other words, we want to calculate  $1 + 2 + 3 + \cdots + N$ . This could easily be done with a loop, but let's derive a recursive solution. The word *derive* is chosen carefully. Many people will tell you that writing a program is an art. However, the tools of mathematics are helpful. We literally will derive a function in the same way that you would derive a formula.

The key idea for this problem is to recognize that it does not matter which order you add the numbers; any ordering gives you the same answer as all the other orderings. One of the many possible ways to transform the task is to group all but the last number into a sum, and then add the last value  $N$  to it:

$$1 + 2 + 3 + \cdots + N = (1 + 2 + 3 + \cdots + (N-1)) + N$$

In other words, we have broken the task of adding numbers from 1 to  $N$  into 2 steps: first add the numbers 1 to  $(N-1)$ , and then add  $N$  to that. Clearly, adding up the numbers from 1 to  $(N-1)$  is a smaller task than adding the number from 1 to  $N$ . If we had some way to do that calculation, all we'd have to do is add  $N$  to the result. Fortunately, the function we are designing is just such a function.

Let's use the notation  $\text{sum}(N)$  to represent the sum of integers from 1 to  $N$ . We have to be careful to use this notation only for  $N > 0$ , because otherwise it doesn't make sense. With this notation, we can also say  $\text{sum}(N-1)$  is the sum of integers from 1 to  $N-1$  (as long as  $N-1 > 0$ ). By the property we observed above, we can finally say that  $\text{sum}(N) = \text{sum}(N-1) + N$ , as long as  $N-1 > 0$ . From here, it is a short exercise to write a recursive function in Python.

```
def sum(N):
    if (N == 1):           # base case
        return 1
    else:                  # recursive case
        return N + sum(N-1)
```

The base case comes from the knowledge that  $\text{sum}(1) = 1$ ; the recursive case comes from the equation  $\text{sum}(N) = \text{sum}(N-1) + N$ .

Notice that we were primarily engaged in understanding the task (summing a bunch of integers), and we used a bit of basic math to describe the properties of the task. When we finished, we translated the math into Python. The recursive function, therefore, describes a mathematical truth about the

task. One only has to understand the language of Python and to understand the nature of addition to see that the program is correct for  $N > 0$ .

We can improve the function by considering the base case test. When given a positive integer value for  $N$ , the function works just fine. But if we call this function with the argument 0, or a negative value like  $-5$ , the base case test will never be true, since the recursive case always subtracts 1 from any  $N$ . This is known as infinite recursion. To avoid this, we can change the base case test to  $N \leq 1$ . This prevents infinite recursion, but also gives a wrong answer; there is no well-defined answer to the mathematical problem  $\text{sum}(-5)$ , and any answer a program will give is wrong. A wrong answer is marginally preferable to an infinite loop; a better approach would be to use an assertion here.

### Even or Not?

Here's a slightly different example. We can write a recursive function to determine whether a given positive integer is even or not. There are better ways to do this task, but it's an interesting example for other reasons.

Suppose we are given a positive integer  $X$ , and suppose for the sake of example that  $X > 2$  (we'll ignore negative numbers). There is a property of even numbers that is extremely useful: if  $X$  is an even number, then so is  $X - 2$ ; likewise, if  $X$  is not even, then  $X - 2$  is also not even. In other words, we have identified a relationship between 2 numbers  $X$  and  $X - 2$ : they are either both even, or both odd.

We can make this relationship work for us in the form of a recursive function. Our function will return the Boolean value true if a given  $X$  is even, and false if  $X$  is odd. We also know that two is an even number, but one is not even:

```
def is_even(X):  
    if X == 0:  
        return True  
    elif X == 1:  
        return False  
    else:  
        return is_even(X - 2)
```

This example has two base cases and one recursive case. Notice that  $X$  is the input to the function, and that the recursive step transforms the task into a subtask about the value  $X - 2$ . The recursive case decides whether  $X - 2$  is even or not, and there is no combination here because the answer for  $X - 2$  is the same as the answer for  $X$ .

Again, we motivated the function by explaining a relationship between  $X$  and  $X - 2$ . The function describes this relationship in Python. It is a matter of understanding something about numbers, and a little Python to see that the program is correct for  $X > 0$ .

## 19.4 Functions, recursion, and the call stack

In Chapter 3, we described how Python creates a new frame every time a function is called. The new frame contains all the variables created inside the function, and all the function's parameters. The body of the function is executed in the context of this frame; that is to say that the statements and expressions are executed and evaluated, and if there are variables used any statement, Python first

looks at the new frame, to see if those names are there. If they are, those values are used; if a name is not present in the frame, Python will look outside the function for those names.

The above description should remind you of the delegation metaphor. The delegate's instructions consist of the code in the function. The office that the delegate works in is the frame that gets created when the function is called.

The only thing that's not realistic about the delegation metaphor is that no delegates are "called in" to help. Instead, the computer acts as its own delegate, preparing the frame, going off to do the work, and then coming back to a previous location. It's as if the CEO of the company does all of the work for every delegation. To manage this, the CEO would have to continually make notes to herself, keeping track of what she was doing before going off to the next office for the next part of the job. One of the notes the CEO needs to make is a reminder about which office she has to return to. There may be several function calls, and several offices, and so it's best to write it down in the frame that is newly created.

The analogy makes clear something that, so far, we've said nothing about. In addition to storing variables and parameters, each frame has a reference called a *dynamic link* that points back to the place where the function was called. This is how Python knows where to return the value that the function returns. When a function is called as an expression or a statement, Python remembers which frame was active when the function was called, as well as which statement or expression was being called. When the function returns, Python follows the dynamic link, and returns the value and continues executing whatever code just after the function call.

Because functions are used frequently, we might be in a situation where several functions have been called, and several frames created. Each frame refers to a previous frame using the dynamic link, forming a kind of stack similar to our node-chains. The top of the frame stack is the frame for the function that is currently running; when that function returns, the frame is discarded, i.e., popped from the stack. When another function is called, the new frame is created and the function's code starts executing, i.e., pushed onto the stack.

As we know, a stack is useful for LIFO behaviour, and data gets pushed onto the stack, and popped off the stack. The stack created by function calls represents all the function calls still active since the start of the program. As we described in Chapter 3, every time a function is called, a frame is created and gets added to this stack of frames. While the function is active, the frame remains on the stack. When the function reaches the return statement, the frame is popped from the stack. The LIFO protocol is crucial here. If function A calls function B, then function B's frame is added after the frame for A was added. Function B must return before function A can continue, so when function B is complete, Python follows the dynamic link to the frame for function A, allowing function A to carry on. In terms of our stack, it's as if B's frame has been popped.

With this understanding, we can understand one of the costs of recursion as implemented in Python: every time a function makes a recursive call, a new frame is created. Each frame takes time to create, and uses up some of the available memory. Every function call pays this cost; it's not just the recursive functions. However, for-loops and while-loops do not pay this cost. Only function calls create frames; loops only use the frame created for the function in which they appear. So in that sense, Python loops consume fewer resources (time and memory) than recursive functions do in Python.

Every programming language has this idea of a stack for frames to keep track of details like local variables, parameters, arguments, and return values. In languages like Python, frames are quite flexible and useful, and are actually part of the heap (there is one more topic about Python frames that we have said nothing about so far). The stack is constructed by using the dynamic link from one

frame to a previous frame.

In other languages, like C/C++, the frames are much less flexible, and are not stored on the heap at all. A separate portion of memory is allocated just for these frames, and it's called a *call stack*, or *system stack*. The call stack is allocated to a running program by the computer's operating system, and has a generous amount of room for the complexity of modern software applications. However, an important aspect of the call stack is that it is finite in size, and cannot be extended. The call stack is unlike any of the stacks we have seen in this course for that reason: we have never encountered a stack that we could not extend, either by making the underlying list a little bigger, or by adding a new node data structure.

The finite size of the call stack implies that the computer could, if too many functions were called, run out of space to use for calling more functions. Running out of space on the call stack is known as *stack overflow*. When this happens, the operating system terminates the program immediately, and data could be lost because of it.

Python programs do not use the system stack; all frames are stored in the heap. But the Python interpreter does use the system stack. The Python interpreter reads your scripts, and line-by-line, decides what should be done to carry out the steps described by your program. The Python interpreter has its own functions, which are called in response to the script you are running, and so the interpreter's functions get placed on the system call stack. It's unlikely that the Python interpreter will encounter a stack over-flow, but it will monitor how many times your script's functions call themselves, and terminate your recursion if it feels you have exceeded an arbitrary limit. This is good when you have an accidental infinite recursion, but not so good when your recursion is under control, but needs to be called more often than the arbitrary limit allows.

There is one piece of information that is still needed to complete the picture of how Python frames work, and perhaps you've wondered about it already. If the body of a function uses a global variable, its name and value will not appear in the frame for the function. The simplified story we told earlier indicates that Python looks first to the local frame, and then to the global frame, when it searches for names. However, the actual story is far more interesting and clever. In addition to the dynamic link, which remembers where the function was *called*, a Python frame also stores a *static link* which remembers where the frame where the function was *created*. Usually, Python functions are created in the global frame when Python reads your script, which is why the simplified story works. However, in Python, a function definition can appear nested inside another function definition, and the every time the internal function is called, it creates a frame whose static link refers to the frame of its enclosing function. This gives the internal function access to all the names in the frame of the enclosing definition. This is an extremely clever idea, which allows us to write programs in Python (and other languages that also use static links) that would be impossible to write without a static link. We'll see a little of this idea later.

## 19.5 Algorithm analysis for recursive functions

Analyzing recursive functions to determine the worst-case time complexity in terms of Big-O notation is really not much different from analyzing a loop. Since repetition is being facilitated by function calls, we need to be able to count how many times the recursive function will be called for a given input size. We also need to understand the cost of the body of the function (ignoring the cost of the recursive call).

The combination of these two ideas can be treated very informally. Consider the  $sum(N)$  function from earlier. We know that the recursive step decreases  $N$  by 1, until the base case  $N = 1$ . In other

words, the function is called  $O(N)$  times recursively, total. The body of the function, is basically  $O(1)$ , except for the cost of the recursive function call. So it's pretty clear that the  $sum(N)$  function has worst-case time complexity of  $O(N)$ .

There is a more rigorous way to analyze the costs. It is possible to formulate the costs precisely as a recursive relationship of costs. And using math that we will not cover in this course, it is possible to establish mathematically the runtime costs. We teach the math needed to do these formal analyses in second year, and these analyses apply not only to situations that seem rather obvious like  $sum(N)$ , but also to much more complicated algorithms where the math is actually the only way to get an answer at all.

Earlier we mentioned the fact that a recursive function call builds a chain of stack frames. There is one frame for every function call, and in the case of  $sum(N)$ , there are  $N$  stack frames in the chain. Thus we can see that as well as having worst-case time complexity of  $O(N)$ , the function also uses  $O(N)$  frames. This is the first time that we've talked about the cost of memory in our analyses. In particular, using a loop does not have this kind of memory use: a loop will access variables that are defined in the stack frame, but the loop itself does not create a stack frame every time it is repeated. In other words, a loop does not have a space cost the same way that a function call does.

In many situations, the fact that recursion consumes stack space, whereas a loop does not, is enough of an argument for preferring loops over recursion. It's a good reason, and in normal situations, if you can use a loop instead of recursion, you should. However, there are situations in which we are willing to pay the cost of recursion. We haven't seen any yet, so you'll just have to wait to see them.

For most computational devices, and most programming languages, a few recursive functions will not be greatly detrimental to your application. You should not be afraid to use recursion, especially if a recursive program seems appropriate. Modern computers are generously equipped with memory, though perhaps smartphones and other embedded devices are not so richly endowed as of this writing.

It must also be said that there is a technique, called *tail-call optimization* which can be applied in some programming languages, that does not incur a new stack frame for every recursive call; this technique reuses the stack frame from the previous call of the recursive function. Python does not use this technique, so in Python, recursion always creates new stack frames. But other languages use tail-call optimization, and so can avoid the costs Python pays (along with many other languages like Java or C/C++). The point here is that it is false to conclude that recursion is never a good choice. The value of recursion depends on context.





## 20 — Trees and Binary Trees

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe what a tree is, in terms of its basic organizational principles.
- Explain the terms: node, child, parent, ancestor, descendant, sibling.
- Describe what a binary tree is, in terms of its basic organizational principle.
- Describe basic operations of binary tree ADT.

### 20.1 Trees organize data hierarchically

In computer science, a tree is way to organize data in a hierarchical manner. This is in contrast to a list, which organizes data in a sequential manner. Hierarchical organizations are very common. Companies, governments, and institutions all use hierarchies to organize their activities.

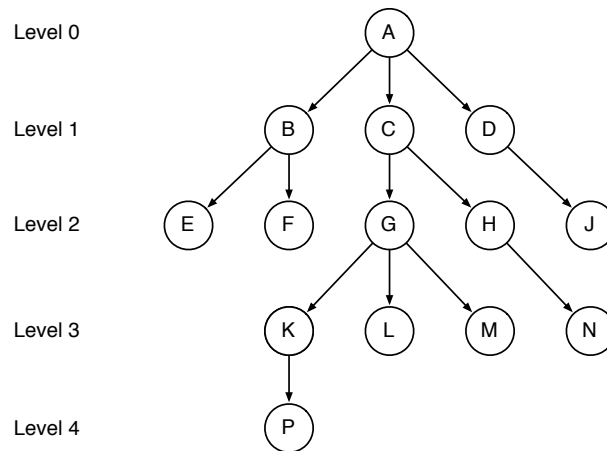
Hierarchies are also used to give structure to information. For example, a textbook consists of a number of chapters, and each chapter consists of a number of sections. Sections contain paragraphs, and paragraphs contain sentences. A table of contents often reflects this organization, from the chapter level down to the section level.

The file system on your computer is organized hierarchically. Folders contain sub-folders. In a UNIX system, there is a special folder, call root, that contains everything stored on the disk. It is not a coincidence that it is called root; this folder is the root of the tree that contains every folder and every document on your computer's disk.

If a computer application has a graphical “drop-down” menu, you will be familiar with menu items that expand into sub-menus. This is hierarchical as well.

Trees are useful in computer science because they can be used to organize data in ways that no other organization could provide, or because using them allows us to write algorithms that are more efficient. We will see examples of both.

Figure 20.1: A tree with 14 nodes, 7 leaf nodes, and height 5.



## 20.2 Informal Terminology

A tree is a collection of data. Each data value is stored in a *node*; this is not a coincidence, as we will define a *tree-node* ADT a little later. In computer science, a tree has a single root. Unlike the root of trees that occur in nature, whose roots have a complicated branching structure, our root is a single node, and in diagrams, we always put the root at the top.

A tree has *branches* connecting nodes. Sometimes the branches are called edges, arcs, or lines. Unlike trees in nature, trees in computer science branch downward; this makes diagrams a bit easier to draw, but also matches the hierarchical organization trees impose on data.

If two nodes are connected by a branch, the node above is called the *parent* of the node below, and the node below is called the *child* of the node above. For example, in Figure 20.1, the node labelled C is a parent of node labelled H; H is a child of C. Sometimes a tree diagram will draw branches with an arrow from parent to child. In general, a parent node may have 1 or more child nodes, however, a child node has exactly one parent. This property is the essence of a tree in computer science. For example, in Figure 20.1, node C has two children, and node H has one child, but no node in the tree has more than one parent.

The only node without any parents is the root; a root is sometimes called an *orphan*, since it has no parent. For example, in Figure 20.1, node A is the root. There may be several nodes that have no children; these nodes are called *leaf* nodes. For example, in Figure 20.1, node F is one of 7 leaf nodes. A leaf is sometimes described as *barren*, since it has no children. If two nodes have the same parent, we call them *siblings*. For example, in Figure 20.1, nodes B, C, and D are siblings. We don't generally extend languages about family relationships much farther than this.

If we start at a node, and follow branches upward in a tree, we find the node's *ancestors*. For example, in Figure 20.1, node G has ancestors G, C, and A. If we start at a node and follow branches downward, we find *descendants*. For example, in Figure 20.1, node G has descendants G, K, L, M, and P. In computer science, a node is usually considered part of its own set of ancestors, and descendants. This doesn't make sense if you think of genealogy, but it makes a kind of mathematical sense: a node is a leaf in its ancestor tree, and a node is the root of its descendant tree.

This last point is worth repeating. Every node can be considered the root of the tree containing

its descendants. As a result, we often call a node a *sub-tree*, even though that is a slight abuse of terminology. A node is not actually a sub-tree, but it is the root of a sub-tree. The tree-ness comes from looking at the way the nodes are organized collectively; the node-ness comes from being connected to other nodes. For example, in Figure 20.1, node G is a node in the tree, but also the root of a sub-tree containing its descendants G, K, L, M, and P.

In computer science, we talk about empty things like empty lists. An empty tree is simply a tree with no nodes. The diagram of an empty tree looks like a drawing of nothing; if you look carefully enough, you will find empty trees drawn everywhere. A tree with exactly one node in it has a root and a leaf, the same node.

## 20.3 Technical definitions

The previous section established intuitions, but it's time to start being precise. The first thing is a precise, and recursive, definition for the term *tree*.

**Definition 20.1** A tree is defined as follows:

1. A structure with exactly zero nodes is an *empty tree*.
2. A structure with exactly one node is a *tree*.
3. If  $t_1, \dots, t_k$  are non-empty trees, then the structure,  $s$ , whose children are the roots of  $t_1, \dots, t_k$  is a *tree*.

The definition is recursive, and this foreshadows the fact that most algorithms that make use of trees will be recursive. In fact, trees are the first real reason we need to learn recursion; most of the examples of recursion we have seen could have been accomplished with a loop. Not so with trees! Most of the computation we will do on trees will not only use recursion, but would be exceedingly difficult to do without recursion. The definition not only defines trees precisely, but it gives the basic organization of every algorithm for trees: the base cases, and the recursive case.

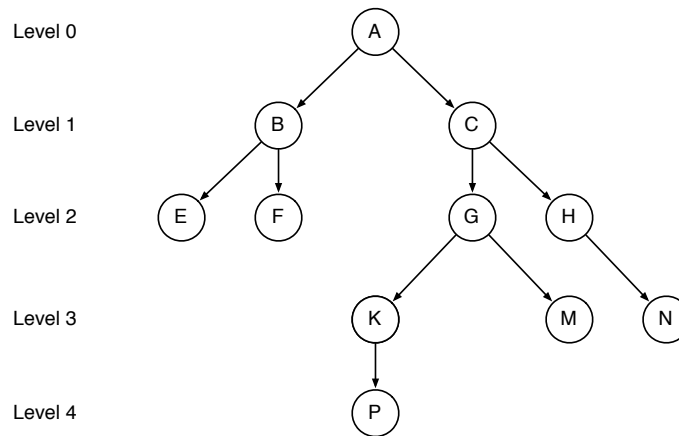
A *path* is a sequence of nodes connected by branches, with no repeated branches allowed. For example, in Figure 20.1, the path between node F and node J goes through nodes F, B, A, D, and J. The *length* of the path is the number of branches in the path. For example, in Figure 20.1, the path between node F and node J has 4 branches. The number of nodes in the path is always one more than the length. There is exactly one path from any node to any other node in the tree. In particular, there is one path from the root to any other node; this path always moves down the tree.

The term *level* describes where a node is in terms of the length of the path from the root to the node. The root node is at level 0, since the path from the root to the root has zero branches. The root's immediate children (if there are any), are at level 1. We usually draw tree diagrams so that all the nodes at the same level are in a row in the diagram. For example, in Figure 20.1, node F is at level 2, and node P is at level 4.

A tree may have many nodes at several levels, but there is always at least one node at the maximum level. For example, in Figure 20.1, node P is at the maximum level 4. The *height* of a tree is always one more than the maximum level of any node in the tree. Another way to put it is that level counts branches, and height counts nodes; The height of a tree is the number of nodes on the longest path from the root to any leaf node in the tree. The height of an empty tree is defined to be zero. For example, in Figure 20.1, the tree's height is 5.

The *ancestors* of a node are nodes on the path from the node to the root. In computer science, we usually say that the ancestors of a node form a *set*, which gives us a mathematical language to

Figure 20.2: A binary tree with 11 nodes, 5 leaf nodes, and height 5.



talk about elements and membership, and other useful concepts. So, according to the definition, a node is an element in its own ancestor set. A node  $u$  is the *descendant* of a node  $v$  if  $v$  is an ancestor of  $u$ . So a node  $s$  is an element in its own descendant set. It may seem strange to include a node as one of its own ancestors, but it simplifies the definitions a lot.

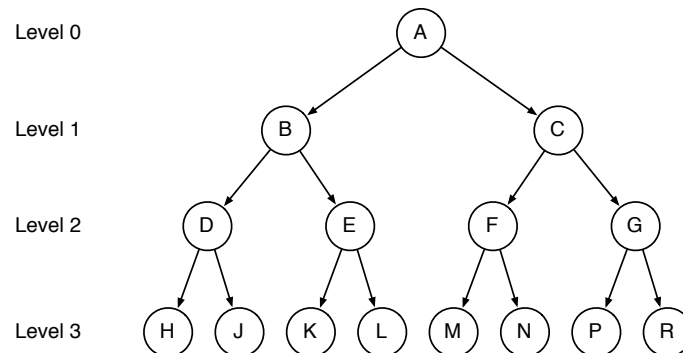
A *sub-tree* is a tree in every sense of the above definition, but it may be part of a larger tree. If  $u$  is a node in a tree, the sub-tree rooted at  $u$  consists of  $u$  and all of its descendants. Sometimes we abuse terminology a little and call  $u$  a subtree, even though we know it's a node; the phrase “subtree rooted at” always adds syllables to a sentence but does not always increase clarity.

## 20.4 Binary trees

A *binary tree* is a special kind of tree with no node having more than 2 children. In part 3 of Definition 20.1, a tree  $s$  can have no more than  $k$  children, and for a binary tree,  $k = 2$ . For example, Figure 20.2 shows a binary tree. Since there are at most 2 children, we can label them *left* and *right*. Sometimes we say *left sub-tree*, or *left node*. In a diagram, we do put the left node to the left of the right node. Either or both of the subtrees could be empty.

A *complete binary tree* is a binary tree that has exactly two children for every node, except for nodes at the maximum level, where the nodes are barren. For example, Figure 20.3 shows a complete binary tree. An empty tree is not complete, because it has no nodes, and therefore no barren nodes. A binary tree with exactly one node (both a leaf and a root) is complete. A complete binary tree has  $2^l$  nodes at level  $l$ . Go ahead, count them: in Figure 20.3, there are 4 nodes at level 2. One of the most important aspects of a tree is the relationship between the number of nodes in a tree, and the tree's height. In total, a complete binary tree whose height is  $h$  contains  $2^h - 1$  nodes. In Figure 20.3, the height is 4, and there are 15 nodes total. This fact implies that a complete binary tree with  $n$  nodes has height  $\log_2(n + 1)$ . These facts are going to be useful for us later.

Figure 20.3: A complete binary tree with 15 nodes, 8 leaf nodes, and height 3.



## 20.5 Binary Treenode ADT

In Chapter 16 we introduced the Node ADT, and we've used it a lot. Perhaps you were skeptical of the utility of nodes, and perhaps you were dubious of the claim that working with nodes was good practice. But we shall see that our work with nodes gives us good background to introduce the ADT that will allow us to create trees.

A **treenode** is a simple object with three attributes:

**data** A data value

**left** A **treenode** (or the value `None`)

**right** A **treenode** (or the value `None`)

The only operation provided for the **treenode** ADTs is the initialization:

- `treenode(data, left=None, right=None)` creates a **treenode** to contain the data value and the given left and right values. If left and right are not given, the value `None` is used by default.

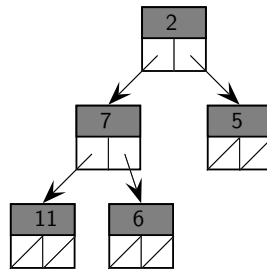
In this implementation, we will make the attributes public. This decision will simplify a lot of the code we write.

```

class treenode(object):

    def __init__(self, data, left=None, right=None):
        """
        Create a new treenode for the given data.
        Pre-conditions:
            data: Any data value to be stored in the treenode
            left: Another treenode (or None, by default)
            right: Another treenode (or None, by default)
        """
        self.data = data
        self.left = left
        self.right = right
  
```

Figure 20.4: A binary tree constructed using the treenode ADT.



A *treenode* is very similar to our nodes from Chapter 16. A node has just one possible reference to another node. A *treenode* has two possible links to other *treenodes*.

Hopefully you can already predict some of the consequences of the similarity between nodes and *treenodes*. The first consequence is the fact that a node-chain is a very special kind of tree. We can imagine chain-like trees, where every *treenode* has at most one child; it has a root, and exactly one leaf. We call this kind of tree *degenerate* because it is not very tree-like. The height of a degenerate tree with  $n$  nodes is exactly  $n$ . Compare this to a complete binary tree with  $n$  nodes, whose height is  $\log_2(n+1)$ .

## 20.6 Trees from treenodes

Just as we did in Chapter 16, we can use the ADT to construct hierarchically organized data structure. We will call this a *primitive binary tree*, though this is not an official computer science term; it is analogous to the term *node-chain*. As we saw in Chapter 10, where we used node-chains to implement Queues and Stacks, we will use primitive binary trees to implement higher-level concepts in later chapters. First, we'll just get the hang of the primitive operations.

Consider the tree in Figure 20.4. There are many ways we could build this tree using the *treenode* ADT, and here's just one of them:

```

import treenode as TN

root = TN.treenode(2)
a = TN.treenode(7)
b = TN.treenode(5)
c = TN.treenode(11)
d = TN.treenode(6)

root.left = a
root.right = b
a.left = c
a.right = d

```

In the above code, five nodes are created, and at first each node is its own one-node tree. Because each tree is assigned to a Python variable, we can construct the tree by connecting the *treenodes* by

name. This is just one way to construct the tree in the figure. Here is another:

```
root = TN.treenode(2,
    TN.treenode(7,
        TN.treenode(11),
        TN.treenode(6)),
    TN.treenode(5))
```

In this construction, we're building the tree without naming the subtrees at all. If you read the code carefully, you'll see that the level of the tree is suggested by the indentation.

Perhaps you are thinking that building trees using the treenode ADT is tedious and possibly error prone. You'd be right! We'll only use the ADT to take care of the basics. Just as we implemented Queues and Stacks using the Node ADT, we'll use the treenode ADT to implement other kinds of hierarchical data structures.





### Sequences and Tree Traversals

- Breadth-order sequence
- Pre-order sequence
- In-order sequence
- Post-order sequence

### Tree Traversals

- Pre-order traversal
- In-order traversal
- Post-order traversal
- Breadth-order traversal

### Algorithms based on Tree Traversals

- Calculating the height of a binary tree
- Counting the number of nodes in a tree
- Searching for a value in the tree
- Generalizations about functions on trees

## 21 — Algorithms on Binary Trees

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe pre-order, in-order, and post-order traversals, and demonstrate their behaviour in example binary trees.
- Employ a binary tree as a solution to a software design problem.

### 21.1 Sequences and Tree Traversals

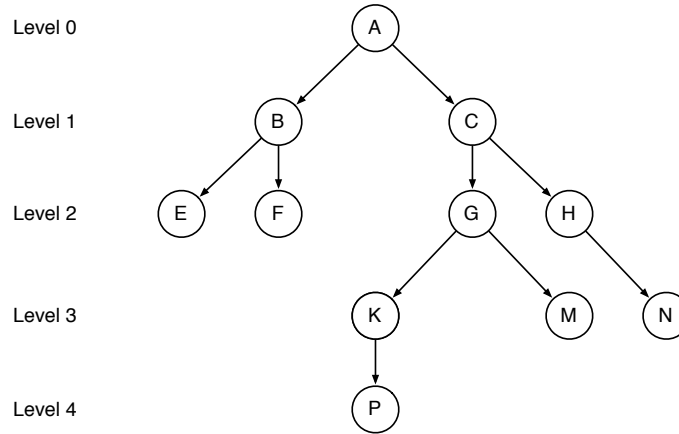
We emphasized in Chapter 10 that linear data structures like lists, stacks, and queues, impose a linear sequence on the data they store. Unlike linear data structures, a tree does not impose a unique or natural ordering on the data it stores. While we could argue that top-to-bottom, left-to-right is a good sequence for data stored in a tree, it is not the only possibility. There are, in fact, four distinct sequences for the data stored in trees. The first is quite a bit different from the remaining three.

We care about sequences, because we have to write programs for binary trees. Our Python programs are sequential, in the sense that they tell the computer to do one thing at a time. When faced with a problem that requires the use of data stored in binary trees, we have to decide how to handle the fact that a treenode may have 2 children. Fortunately, there are 4 options which will form the basis of many programs that deal with data stored in binary trees. First we'll define the sequences, and the property that characterizes them. After that, we'll show four template functions that demonstrate the sequences. These functions are known as "traversals." Following that, a few examples of programs that use one of the traversals as the basis of the algorithm.

#### 21.1.1 Breadth-order sequence

We already have the notion of treenode levels, as being the distance of a treenode from the root. Informally, the *breadth-order* sequence works top-to-bottom, left-to-right, as mentioned just previ-

Figure 21.1: A binary tree. (a) breadth-order sequence: A,B,C,E,F,G,H,K,M,N,P; (b) Pre-order sequence: A,B,E,F,C,G,K,P,M,H,N; (c) In-order sequence: E,B,F,A,P,K,G,M,C,H,N; (d) Post-order sequence: E,F,B,P,K,M,G,N,H,C,A;



ously. Formally, we say that the breadth-order sequence arranges the treenodes so that all the nodes at level  $i$  appear in the sequence after the nodes at level  $i - 1$  and before the nodes at level  $i + 1$ .

For example, given the binary tree in Figure 21.1, the breadth-order sequence of nodes is as follows: A,B,C,E,F,G,H,K,M,N,P.

### 21.1.2 Pre-order sequence

The pre-order sequence organizes the nodes so that the data at the root appears before the data at its subtrees. For example, given the binary tree in Figure 21.1, the pre-order sequence of nodes is as follows: A,B,E,F,C,G,K,P,M,H,N. To make sense of this organization, choose any node in the tree, say the node with data value C. It is in the subtree of A, so it appears after A in the pre-order sequence. But it is the root of its own subtree, so in the pre-order sequence, it appears before every node below it in the tree. It is before G and H, as well as P and M, for example.

### 21.1.3 In-order sequence

The in-order sequence organizes the nodes so that the data at the root appears after the data at its left subtree, but before the data in its right subtree. For example, given the binary tree in Figure 21.1, the in-order sequence of nodes is as follows: E,B,F,A,P,K,G,M,C,H,N. To make sense of this organization, choose any node in the tree, say the node with data value C. It is in the right subtree of A, so it appears after A in the in-order sequence. But it is the root of its own subtree, so in the in-order sequence, it appears after every node below it in its left subtree, e.g., C appears after K. In the in-order sequence it appears before anything in its right subtree, e.g., N.

### 21.1.4 Post-order sequence

The post-order sequence organizes the nodes so that the data at the root appears after the data at its subtrees. For example, given the binary tree in Figure 21.1, the post-order sequence of nodes is as

follows: E,F,B,P,K,M,G,N,H,C,A. To make sense of this organization, choose any node in the tree, say the node with data value C. It is in the subtree of A, so it appears before A in the post-order sequence. But it is the root of its own subtree, so in the post-order sequence, it appears after every node below it in the tree. It is after G and H, as well as P and M, for example.

Take care not to draw the wrong conclusion about pre-order and post-order sequences. The post-order sequence still obeys the left-to-right convention, and this implies that the post-order is not simply the reverse of the pre-order sequence.

## 21.2 Tree Traversals

In this section, we'll give simple programs that demonstrate the traversals. Keep in mind that the programs we are showing are not inherently useful for what they do. In fact, the programs we'll see simply display the tree data to the console. Clearly there are more interesting things to do with data than display. However, these programs demonstrate the application of the sequences discussed in the previous section. We'll use the phrase "process the data" to indicate in an abstract way that some calculations or computations should be done, and those will depend on the application. We'll simply display the data, for now. When you need an algorithm to accomplish something in a binary tree, chances are very high that you'll base the algorithm on one of these traversals, and you won't be displaying data, you'll be doing something more interesting.

We'll start with pre-, in-, and post-order traversals. These are simple recursive algorithms, and quite easy to implement. We'll introduce the breadth-order traversal last. While it is the easiest one to describe in words (top-to-bottom, left-to-right), it is the hardest one to write a program for. This fact tells us something about human language, and the value of diagrams.

### 21.2.1 Pre-order traversal

The pre-order sequence organizes the nodes so that the data at the root appears before the data at its subtrees. Therefore, the pre-order traversal processes the data at the root before processing any data at its subtrees. Taking this idea, we can write a recursive algorithm for a pre-order traversal:

1. Process the root of the subtree.
2. Recursively process the left-subtree in pre-order sequence
3. Recursively process the right-subtree in pre-order sequence

In Python, we'll be a little more careful, by including a base case, as follows:

```
1 def pre_order(tnode):
2     """
3     Display the nodes of a tree in pre-order.
4     :param tnode: a primitive tree
5     :return: nothing
6     """
7     if tnode is None:
8         return
9     else:
10        print(tnode.data, end=" ")
11        pre_order(tnode.left)
12        pre_order(tnode.right)
```

Here, the root is processed at line 8, which happens before the two recursive calls.

On tests and exams and job-interviews, you might be called upon to demonstrate your knowledge of the pre-order sequence without the aid of a computer. In this case, it pays to have a strategy. Use this one for guaranteed success:

1. Write down the root.
2. Draw a box for left and right subtrees.
3. Fill in each box recursively: left, then right.

### 21.2.2 In-order traversal

The in-order sequence organizes the nodes so that the data at the root appears after the data at its left subtree, but before the data in its right subtree. Therefore, the in-order traversal processes the data at the root after it processes the left subtree, but before processing the data at its right subtree. Taking this idea, we can write a recursive algorithm for an in-order traversal:

1. Recursively process the left-subtree in in-order sequence
2. Process the root of the subtree.
3. Recursively process the right-subtree in in-order sequence

In Python, we'll be a little more careful, by including a base case, as follows:

```
1 def in_order(tnode):
2     """
3     Display the nodes of a tree in pre-order.
4     :param tnode: a primitive tree
5     :return: nothing
6     """
7     if tnode is None:
8         return
9     else:
10         in_order(tnode.left)
11         print(tnode.data, end=" ")
12         in_order(tnode.right)
```

Here, the root is processed at line 9, which happens after finishing the left subtree, and before starting the right subtree.

On tests and exams and job-interviews, you might be called upon to demonstrate your knowledge of the in-order sequence without the aid of a computer. In this case, it pays to have a strategy. Use this one for guaranteed success:

1. Draw a box for left and right subtrees.
2. Write down the root between the two boxes.
3. Fill in each box recursively: left, then right.

### 21.2.3 Post-order traversal

The post-order sequence organizes the nodes so that the data at the root appears after the data at its subtrees. Therefore, the post-order traversal processes the data at the root after processing any data at its subtrees. Taking this idea, we can write a recursive algorithm for a post-order traversal:

1. Recursively process the left-subtree in post-order sequence
2. Recursively process the right-subtree in post-order sequence
3. Process the root of the subtree.

In Python, we'll be a little more careful, by including a base case, as follows:

```

1 def post_order(tnode):
2     """
3     Display the nodes of a tree in pre-order.
4     :param tnode: a primitive tree
5     :return: nothing
6     """
7     if tnode is None:
8         return
9     else:
10        post_order(tnode.left)
11        post_order(tnode.right)
12        print(tnode.data, end=" ")

```

Here, the root is processed at line 10, which happens after the two recursive calls.

On tests and exams and job-interviews, you might be called upon to demonstrate your knowledge of the post-order sequence without the aid of a computer. In this case, it pays to have a strategy. Use this one for guaranteed success:

1. Draw a box for left and right subtrees.
2. Write down the root after the boxes
3. Fill in each box recursively: left, then right.

#### 21.2.4 Breadth-order traversal

The breadth-order sequence arranges the treenodes so that all the nodes at level  $i$  appear in the sequence after the nodes at level  $i - 1$  and before the nodes at level  $i + 1$ . Therefore, the breadth-order traversal processes all the data (from left to right) at level  $i$  before it processes any data at level  $i + 1$ . Unlike the previous traversals, the breadth-order traversal is not recursive.

The key to the breadth-order traversal is to understand that if a treenode in the tree is at level  $i$ , its children (if there are any) are at level  $i + 1$ . We have to schedule the processing of a node's level  $i + 1$  children after all the other nodes at level  $i$  are processed. Fortunately, we have a nice tool to do this kind of scheduling: a FIFO queue. To see how this would work, let's assume that all the treenodes at level  $i$  are already in the queue. If we look at any level  $i$  node already in the queue, we can enqueue its children, and by the FIFO protocol, the children are added after everything currently in the queue. In other words, the children at level  $i$  are guaranteed to appear after all the level  $i$  treenodes.

For example, consider the tree in Figure 21.1. Suppose that the treenodes B, C are in a FIFO queue. If we enqueue B's children, the treenodes E and F will be in the queue after C.

The breadth-order first traversal uses a FIFO queue to store treenodes, starting with the root. The algorithm dequeues a treenode, processes it, and then enqueues its children, ignoring empty trees. The complete algorithm is given below:

```

1 def breadth_order(tnode):
2     """

```

```

3      Display the nodes of a tree in breadth-order.
4      :param tnode: a primitive tree
5      :return: nothing
6      """
7      explore = Q.Queue()
8      Queue.enqueue(explore, tnode)
9
10     while explore.size() > 0:
11         current = explore.dequeue()
12         print(current.get_data(), end=" ")
13         left = current.get_left()
14         if left is not None:
15             explore.enqueue(left)
16         right = current.get_right()
17         if right is not None:
18             explore.enqueue(right)

```

## 21.3 Algorithms based on Tree Traversals

The functions given in the previous section are perhaps useful for displaying the contents of a binary tree, during debugging, for example, but not for much else. However, traversals are commonly the foundation of other algorithms that do more interesting tasks. In this section, we'll see some examples. You'll notice a few very important common features of these algorithms.

The recursive definition of a tree, given in the previous chapter, provides guidance in the design of algorithms for trees. They'll almost always be recursive, and even if recursion is not needed, recursion will always be a good choice. Second, the base case(s) will frequently be either an empty tree or a tree with no children (a leaf). Third, the recursive cases can be designed by assuming that the function returns the correct result for subtrees; we'll only have to design the calculations that combine the results from the subtrees.

### 21.3.1 Calculating the height of a binary tree

The height of a tree is the number of nodes on the longest path from the root to any leaf node in the tree. When we see a tree in a diagram, we can almost always identify the deepest leaf nodes, and so the height of a tree is not very difficult to calculate. On the other hand, when a computer is asked to find the height of a tree, starting at the root, it's not at all obvious where the deepest leaf node is. For that reason, a function to find the height of a tree has to try all possibilities.

The key to designing a recursive function for trees is to assume that the function will work on a *treenode*'s subtrees, and then focus on designing the combination step. Let's assume that our tree *t* has two subtrees, the left *lt* and the right, *rt*. We will assume that our function will correctly calculate the height of *lt* and *rt*, leaving us the job of using this information to calculate the height of *t*. Either the two subtrees have equal height, or one of them will have a greater height. In either case, the height of *t* involves one more *treenode* on the path than one either subtree. In Python, we have the following function:

```

1  def height(tnode):
2      """

```



```

3     Purpose:
4         Determine the height of a tree.
5         The height is defined as the length of the longest path
6         from the root to any leaf.
7     Pre-conditions:
8         :param tnode: a treenode
9     Post-conditions:
10        none
11    Return
12        :return: the height of the tree as an integer
13    """
14    if tnode is None:
15        return 0
16    else:
17        lh = height(tnode.left)
18        rh = height(tnode.right)
19        return 1 + max(lh, rh)

```

The recursive calls for the left and right subtrees precede the calculation of the maximum, so we can view this as being based on a post-order traversal.

### 21.3.2 Counting the number of nodes in a tree

To count the number of nodes in a tree, we will use the same recursive approach as above. The base case, an empty tree, has 0 nodes in it, by definition. For the recursive case, we observe that we can count the number of nodes in the left subtree recursively, and the number of nodes in the right subtree recursively. The total number of nodes in the tree is one more than the sum. In Python:

```

1     while explore.size() > 0:
2         current = explore.dequeue()
3         print(current.get_data(), end=" ")
4         left = current.get_left()
5         if left is not None:
6             explore.enqueue(left)
7         right = current.get_right()
8         if right is not None:
9             explore.enqueue(right)
10
11
12 if __name__ == '__main__':
13     # Create the tree we've been using in class
14     root = tn.treenode(2)
15     a = tn.treenode(7)
16     b = tn.treenode(5)

```

As written, the count() function most closely resembles a post-order traversal, since the additions cannot be performed until the second recursive call returns. But it resembles all the recursive traversals to a significant degree.

### 21.3.3 Searching for a value in the tree

We use trees to organize data hierarchically. At this point, you may well wonder at the utility of using trees, since we've not seen much that's special about trees so far. It will come, in the next chapter! But for now, let's assume that, like an unordered list, a data value could be anywhere in the tree. The question is how to go about answering whether or not a given value is in the tree or not.

There are four cases. First, we can check if the value is stored in the root of the given tree. Second, it could be in the left subtree; or it could be in the right subtree, which is the third case. For both of these cases, we will use the power of recursion. Finally, the value we're looking for may not be in the tree at all. This last case is a bit subtle. We can only say for sure that the value is not in the tree if we look at the whole tree, and have no success finding it. Of course, if the tree happens to be an empty tree, it will contain nothing, and in particular it will not contain the value we are looking for. In Python:

```

1      :param tnode: a primitive tree
2      :return: nothing
3      """
4      if tnode is None:
5          return
6      else:
7          post_order(tnode.left)
8          post_order(tnode.right)
9          print(tnode.data, end=" ")
10
11
12 def breadth_order(tnode):
13     """
14     Display the nodes of a tree in breadth-order.
15     :param tnode: a primitive tree
16     :return: nothing
17     """
18     explore = Q.Queue()
```

Note the use of the Boolean operator `or` in the return statement. It might be tempting to write something like the following:

```

def member(tnode, val):
    if tnode == None:
        return False
    elif tnode.data == val:
        return True
    elif member(tnode.left, val):
        return True
    elif member(tnode.right, val):
        return True
    else:
        return False
```

This is a correct function, but it says in eight lines of Python what the previous version says using `or` and 4 lines. Note that most programming languages, Python included, define the `or` operation so that it does not perform unneeded computations. Logically, the expression `True or x` evaluates to `True`, no matter what `x` is. For this reason, in an expression `y or x`, if `y` is `True`, the expression `x` would only be evaluated if `y` were `False`.

#### 21.3.4 Generalizations about functions on trees

We can make a few generalizations about the examples we've seen. First, almost all functions that have to work with trees are recursive, and more or less similar to the recursive traversals. It is conceivable to write non-recursive functions for trees, but to do so, you'd have to write code that basically simulates the call stack, to hold the data that would be pushed to the call stack in a recursive function anyway. The very idea of avoiding recursion is absurd to the extreme. Second, the structure of recursive functions on trees closely mirrors the formal definition of a tree. This is not an accident. The formal definition of a mathematical object, like a tree, identifies very precisely the salient and valuable features<sup>1</sup> of the object. The strategy for designing a recursive tree function on tree  $t$  is to assume that the function will return the correct answer for the subtrees,  $lt$  and  $rt$ , and to focus on the combination step; taking the information returned by the recursive calls, and using it to formulate the correct answer for  $t$ . This is not an accidental strategy. It's based on a fundamental property called mathematical induction.

---

<sup>1</sup>And you thought mathematics was about numbers? It's *not*. It's about making concepts precise enough to reason about rigorously. We learn math, and use math, when we want to reason effectively.



## Introduction

## Binary Search Trees

## Basic operations

Checking if the value is in the BST

Adding a value to the BST

Removing a value from the BST

## The importance of Binary Search Trees

# 22 — Binary Search Trees

## Learning Objectives

After studying this chapter, a student should be able to:

- Describe what a binary search tree is, in terms of its basic organizational principle.
- Describe basic operations of binary search tree, and demonstrate the behaviour in example binary search trees.
- Explain how tree balance affects the efficiency of various binary search tree operations.

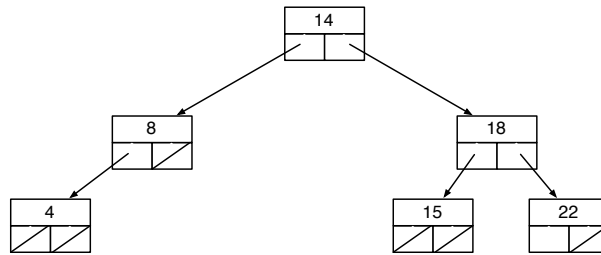
## 22.1 Introduction

Many applications and algorithms depend on being able to store data in collections, like lists and dictionaries, so that the data can be looked up efficiently. It is so common that modern languages like Python provide tools to assist programmers. We've already mentioned lists and dictionaries, but it is also important to point out the Python operator `in` which can be used to check if a data value appears in a collection. For example, given a Python list `alist`, the expression `x in alist` performs a linear search, returning `True` if the value appears in the list. The time complexity of linear search of a list with  $N$  elements is  $O(N)$ .

As we have seen in the previous course, keeping the data organized can make a big difference. For example, if the data in the list is sorted, we can apply binary search to find a given value. The worst case time complexity of binary search of a sorted list with  $N$  elements is  $O(\log N)$ . It is very important to understand how much of an improvement binary search is over linear search. It's not just a little better. It's outrageously better.

To get the benefit of binary search in lists, we have to start with a sorted list, and it's important to understand the costs of sorting. The best sorting methods, for example, the Quick Sort algorithm, have average case time complexity of  $O(N \log N)$ . It turns out that no general purpose sorting algorithm can be expected to do better than  $O(N \log N)$  on average.

Figure 22.1: A binary search tree. Values on the left are smaller. Values on the right are larger.



So if we have a list of data that never changes, sorting it once at the beginning, and using binary search for repeated look-up is a very efficient strategy. On the other hand, if we have an application where values are added and deleted from the list frequently, keeping it in sorted order can be very expensive. It is possible to keep sorting the data after every addition or deletion, paying the cost of sorting, i.e., repeatedly doing  $O(N \log N)$  steps. It is also possible to ensure that the data is always sorted by inserting new data in the appropriate place, to maintain sorted order; this amounts to linear search, repeated for every new data value, at a cost of  $O(N)$  every time. The question is whether we can do better than either of these two options, and the answer is yes.

## 22.2 Binary Search Trees

A *binary search tree* (BST) is a special kind of binary tree that allows us to keep data organized. The binary search tree is a collection, like a list. We'll be able to add data, delete data, and search for data in a binary search tree far more efficiently than using a Python list. The tree structure of the data allows us to avoid having to use a sorting algorithm to keep the data organized.

A binary search tree is just a binary tree, with up to two children per node, and data stored as usual. What makes a binary search tree special is the way the data in the tree is organized.

**Definition 22.1** A binary tree satisfies the *binary search tree property* if at every node  $v$  in the tree:

1. All the data stored in the left subtree of  $v$  is smaller than the data at  $v$ .
2. All the data stored in the right subtree of  $v$  is larger than the data at  $v$ .

If a binary tree satisfies this property, it is called a binary search tree.

An example of a binary search tree is shown in Figure 22.1. The definition above implicitly assumes that a binary search tree does not store multiple copies of the same data value, though the definition could be adjusted to allow multiple copies. Also, the terms *smaller* and *larger* could mean before and after according to a given rule, e.g. alphabetical order. A binary search tree can be used to store numerical data, text data, or any other kind of data that has a natural ordering to the values. It is interesting to note that an in-order traversal of a binary search tree processes the data in sorted order.

Just a word of caution: It is common for students to assume that every binary tree is a binary search tree. This is not true! There are lots of different ways to organize data in a binary tree, and the BST property is only one possible way. Another organizational principle is called the *heap property*.

## 22.3 Basic operations

There are three important operations on binary search trees that we have to discuss:

1. `member_prim(tnode, value)`: checks if the value is in the binary search tree rooted at `tnode`.
2. `insert_prim(tnode, value)`: adds the value to the binary search tree rooted at `tnode`.
3. `delete_prim(tnode, value)`: removes the value from the binary search tree rooted at `tnode`.

We will discuss these in order.

### 22.3.1 Checking if the value is in the BST

If we were searching for a value in an **arbitrary** binary tree (one that does not have the BST property), we would have to check three different places to find a particular value.

1. First, check the root. If the value you want is stored as data in the node, you have successfully found what you were looking for.
2. Second, check if the value is on the left side of the tree.
3. If the item was not in the left side of the tree, check the right side of the tree.

This analysis lends itself to a simple recursive function.

When we are searching for a value in a binary search tree, we can be a little more intelligent about our search. Given a tree *t* with the BST property and a data value:

1. If *t* is *empty*, return False
2. If *t* *stores the same data value*, return True
3. If the data value is *smaller than the root*, look *left* recursively
4. If the data value is *larger than the root*, look *right* recursively

In Python, the function is as follows:

```
def member_prim(tnode, value):
    """
    Purpose:
        Check if value is stored in the binary search tree.
    Preconditions:
        :param tnode: a binary search tree
        :param value: a value
    Postconditions:
        none
    :return: True if value is in the tree
    """
    if tnode is None:
        return False
    else:
        cval = tnode.data
        if cval == value:
            # found the value
            return True
```



```
elif value < cval:
    # use the BST property
    return member_prim(tnode.left, value)
else:
    return member_prim(tnode.right, value)
```

It's important to draw attention to how the function handles a search that comes up empty. Essentially, the search keeps stepping left or right until it steps right off the tree, i.e., the subtree we stepped to is empty. While a human might be able to keep a context in mind, especially if a tree diagram is available, the function does not make any attempt to keep track of which nodes were looked at.

Search in a binary search tree only follows one of the possible two branches at any node. If the search is lucky, the value we're looking for is near the top of the tree. But if we're not lucky, the value is stored in a leaf node, or, the value is not in the tree at all. As a result, the worst case time complexity for `member_prim()` happens when the search ends up exploring a path to the deepest part of the tree.

As with all recursive functions, we analyze the computational complexity of the function by counting the number of recursive calls. The number of recursive calls needed to get to the bottom of the tree depends on the height of the tree. In a complete binary tree with height  $h$ , the number of nodes is  $N = 2^h - 1$ . Solving for  $h$ , we can say that the height of a complete binary tree with  $N$  nodes is  $h = \log(N + 1)$ . Thus, in the worst case, the computational complexity for search in a balanced binary search tree is  $O(\log N)$ .

Binary trees don't have to be balanced, though. Informally, in a balanced tree, most nodes have 2 children. It is very common for trees to have at least a few nodes with only one child. It is possible to build a tree in which each node has no more than 1 child. We call such trees *degenerate* because it is not very tree-like. In a degenerate tree, every node except the leaf node has one child. Therefore the number of nodes in a degenerate tree of height  $N$  is exactly  $N$ .

The worst case time complexity for search in general binary search trees is  $O(N)$ , simply because in the worst case, a binary search tree might be degenerate.

### 22.3.2 Adding a value to the BST

When we talk about adding a value to a binary search tree, we will assume that the value is not already in the tree. If the value happens to be in the tree already, we will not add a second copy. We could modify the rules to allow second copies, but that adds detail without enhancing clarity, so we'll leave it alone.

The key algorithmic insight that we use to put a value into a binary search tree is to put the value exactly where `member_prim()` would find it if it were in the tree already. In other words, our algorithm would start by searching the tree for the value. Since the value is (presumably) not in the tree, eventually, the search will wind up looking at the value `None`, only to realize that the value is not in the tree. The new value that we're adding to the tree will always be a leaf node, as either a left or right child of a node that is currently in the tree. For example, to add the value 10 to the tree in Figure 22.1, we'd search for the value 10 by stepping left from 14, and right from 8. The new node containing the value 10 should be the right subtree of the node containing 8.

Remember that we made note of this situation when we discussed the `member_prim()` function. We observed that search does not make any attempt to keep track of how it arrived at a `None`, i.e., an empty tree. This `None` that we're looking at is exactly the `None` we want to replace, by adding a new

leaf. But this new leaf cannot connect itself to its parent node; the parent node has to connect itself to the leaf, by storing the reference to the new leaf as a left or right subtree.

It is possible to write some tricky Python code to handle the special situation that arises when a node has to attach a new leaf node, but there is a simpler solution. Instead of having special code that only attaches a leaf, we'll always make sure that every node re-attaches its child node to the appropriate place. In most cases, this is simply an assignment that was unnecessary, but in the case of adding a new leaf, the new connection will be made with a very small amount of fuss.

Given a tree  $t$  with the BST property and a data value:

1. If the tree is *empty*, return a new tree with the value stored
2. If the root *stores the same data value*, return the root
3. If the data value is *smaller than the root*, update the *left* subtree, recursively
4. If the data value is *larger than the root*, update the *right* subtree recursively

In Python, the function is as follows:

```

1 def insert_prim(tnode, value):
2     """
3     Insert a new value into the binary tree.
4     Preconditions:
5         :param tnode: a binary search tree
6         :param value: a value
7     Postconditions:
8         If the value is not already in the tree, it is added
9     Return
10        :return: flag, tree
11        Flag is True if insertion succeeded;
12            tree contains the new value
13        Flag is False if the value is already in the tree,
14            tree unchanged
15    """
16
17    if tnode is None:
18        return True, TN.TreeNode(value)
19    else:
20        cval = tnode.data
21        if cval == value:
22            return False, tnode
23        elif value < cval:
24            flag, subt = insert_prim(tnode.left, value)
25            if flag:
26                tnode.left = subt
27            return flag, tnode
28        else:
29            flag, subt = insert_prim(tnode.right, value)
30            if flag:
31                tnode.right = subt

```

```
32         return flag, tnode
```

There are two base cases. The first base case, lines 17-18, handles the creation of a new `treenode`, which happens as soon as the search reaches an empty tree. Pay special attention to the fact that a tuple is being returned: the value `True` to indicate a successful insertion, and the newly created leaf node. The other base case, lines 21-22, handles the situation that the value is already in the tree. Again a tuple is returned: this time the value `False` to indicate that no value was inserted, and the tree unchanged.

The two recursive cases similar. If the new value should be inserted on the left (lines 24-27), the recursive call starts on the left subtree (line 24). But notice how the resulting tuple is used. If the insertion was successful (line 25), the node sets the left subtree (line 26). Here `subt` could be the new leaf node, or it might be just the old subtree with a new leaf attached somewhere far below `tnode`. In either case, the left subtree is attached. The case for the right subtree is essentially the same (lines 29-32).

Since the `insert_prim()` function essentially does a search, the computational complexity again depends on the height of the tree. If the tree is balanced, we have  $O(\log N)$ , as before. If the tree is degenerate, we have  $O(N)$ , again, as before, for exactly the same reasons.

### 22.3.3 Removing a value from the BST

Like the previous two examples, the function to remove a value from the tree first has to find the value in the tree. Unlike insertion, where we always insert a new `treenode` as a leaf in the tree, there are four different situations for deletion:

1. The node to be deleted has no children. For example, node 4 in Figure 22.1. In this case we can just remove the leaf, and the resulting tree will remain a binary search tree. The way to remove the leaf is to replace the reference to the leaf with a `None`.
2. The node to be deleted has only a left subtree. For example, node 8 in Figure 22.1. In this case we can just connect the parent of node 8 to node 8's left child, and the resulting tree will remain a binary search tree. This connection can be implemented by replacing the reference to node 8 with the reference to node 8's left child.
3. The node to be deleted has only a right subtree. This is similar to case 2.
4. The node to be deleted has two subtrees. For example, node 14 in Figure 22.1. This is the only case where the right thing to do is perhaps not obvious. We discuss it below in detail.

These 4 cases have to be handled separately.

#### Case 4: the node to be deleted has a left and right child

Suppose we want to delete the node containing 14 in Figure 22.1. We'll need to make one tree by connecting the two subtrees, somehow. There are a bunch of different ways we could do this, but we have to make sure that, whatever we do, the resulting tree will satisfy the binary search tree property.

If we look carefully at the tree, and think about the binary search tree property, we can make an important observation: every value in 14's left subtree is smaller than any value in 14's right subtree. This observation suggests a strategy: We can connect 14's right subtree, the tree rooted at 18, to the right of the biggest node in 14's left subtree. In this example tree, the value 8 is the biggest value in 14's left subtree, since it has no right child. But in general, 14's left subtree might contain values to the right of 8, which are bigger than 8, but less than 14. This strategy will connect the whole tree starting at 18 as 8's right subtree.

In general, we have to find the largest value in a subtree to manage case 4. Fortunately, the binary search tree property tells us exactly where this value is: keep stepping right until you find a node that has no right subtree.

To summarize, given a tree with the BST property and a data value to delete, we can find the node to be deleted as follows:

- If the tree is *empty*, return False
- If the root *stores the same data value*, *reconnect* the root's children (see Algorithm reconnect)
- If the data value is *smaller than the root*, *delete* the value from the *left* subtree, recursively
- If the data value is *larger than the root*, *delete* the value from the *right* subtree recursively

Once found, the node to be deleted is at the root of a subtree using the reconnect algorithm, as follows. Given a tree with the BST property, whose root is to be deleted:

- If the root has *no children*, return None (the empty tree)
- If the root has *a left child, but no right child*, return the left child.
- If the root has *a right child, but no left child*, return the right child.
- If the root *has 2 children*:
  - Start at the left child, and step down, always going right, to find left's largest node.
  - Attach the root's right child as the right subtree of left's largest node
  - Return left.

The deletion operation has two phases: first, finding the node to be deleted; second, reconnecting the tree to maintain the binary search tree property. To find the node to be deleted, we have to descend into the tree, and in the worst case, we have to go to the deepest node in the tree. The first phase is essentially the same as the search algorithm. In the second phase, cases 1-3 are very simple, requiring a couple of assignments only, and therefore,  $O(1)$ . Case 4 warrants special consideration, since the reconnection requires stepping down the node's left subtree. But observe that the total amount of stepping, in the worst case, reaches the deepest part of the tree. In other words, the worst case depends on the height of the tree, as with search, and with insert. So, when the tree is balanced, the worst-case time complexity is  $O(\log N)$ , or if the tree is degenerate, the worst case time complexity is  $O(N)$ .

## 22.4 The importance of Binary Search Trees

If we can guarantee that the binary search trees are always balanced, every operation we've considered in this chapter has worst case time complexity of  $O(\log N)$ , where  $N$  is the number of data values stored in the tree. In comparison, the cost of keeping a sorted list sorted while adding and deleting values is at best  $O(N)$ . The improvement of  $O(\log N)$  of  $O(N)$  is astonishingly good. It may seem trivial if you misunderstand the nature of the logarithm function. When your costs can be described by  $O(\log N)$ , it's the practical equivalent of saying *no normal, practically realistic increase in the problem size will have any noticeable effect on the runtime*. This is an astounding statement. Keep reading it until you are astounded. We'll wait.

Keep in mind that we have said nothing about keeping the trees balanced. None of our algorithms made any attempt to do that. One can easily imagine a sequence of insertions and deletions that could cause binary search trees to become degenerate, and thereby lose the logarithmic advantage. In fact, there are more complicated algorithms that rearrange the tree at every insertion or deletion to guarantee balanced search trees at all times, but those algorithms are not studied in first year

---

computer science. They are a second year topic!

## 23 — The Table ADT

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe what a Table is, in terms of its possible applications.
- Explain how a variation of the treenode data structure can be used to implement the Table ADT.
- Explain how the binary search tree can be used to make the Table ADT implementation efficient.

### 23.1 Introduction

In the previous chapter we saw how to use binary tree data structure, and the binary search tree property to organize data so that tasks such as data look-up, insertion, and deletion are very efficient. In this chapter we will complete the picture by defining an ADT called a Table, which will build on binary search trees, but will allow arbitrary data to be stored.

The Table ADT will store keyed data. A key is a value that is unique to the data being stored, like a Student number. The key tells us which data we're looking for, but the data may be much more than the key, for example, a complete Student record. The keys have to be comparable, like numbers or strings.

The Table ADT will provide yet another example of an ADT, but more importantly, will be a direct application of binary search trees. Furthermore, by studying the implementation of the Table ADT, we will get an insight into one of Python's key data structures. A dictionary is just Python's own implementation of our Table ADT.

The primary operations for the Table ADT are as follows:

- `is_empty()`: Query if table is empty
- `size()`: Query size of table

- `insert()`: Insert key, value into table
- `delete()`: Delete key, value from table
- `retrieve()`: Retrieve the data for a given key

Our Table ADT will be implemented using the object-oriented variation of the `TreeNode` class we saw in Chapter 22. Our variation, called a `KVTreeNode`, will store two different pieces of information at a node: a key value, which we will use to organize the tree structure, and the data value, which will be the payload. For example, in a student information system, the key would be the Student number, and the data value would be the complete student record.

The `KVTreeNode` is defined as follows:

```
class KVTreeNode(object):
    def __init__(self, key, value, left=None, right=None):
        """
        Create a new KVTreeNode for the given data.
        Pre-conditions:
            key:      A key used to identify the node
            value:    Any data value to be stored in the KVTreeNode
            left:     Another KVTreeNode (or None, by default)
            right:    Another KVTreeNode (or None, by default)
        """
        self.value = value
        self.left = left
        self.right = right
        self.key = key
```

We've made all the attributes public, for convenience.

## 23.2 Table ADT implementation

The implementation of the Table ADT will use the `KVTreeNode` discussed above, as well as slight variations of the binary search tree algorithms discussed in Chapter 22.

We'll make it a class. The `__init__()` method, given below, demonstrates that a Table is an object with two attributes, `root` and `size`:

```
class Table(object):
    def __init__(self):
        self.__root = None
        self.__size = 0
```

The `size` field allows us to keep track of the number of nodes in the Table, so that we don't have to count them when someone asks for the size of the Table.

The `retrieve` method takes a key, and returns the data value associated with the key in the Table, if it exists. The operation uses an internal function, which is very similar to the `member_prim()` function from Chapter 22. The only significant difference is that the key value is used to locate the appropriate treenode (e.g., line 17), but the data value is returned (line 19).

```
1     def retrieve(self, key):
2         """
```



```

3         Return the value associated with the given key.
4         Preconditions:
5             :param key: a key
6         Postconditions:
7             none
8         Return
9             :return: True, value if the key appears in the table
10                False, None otherwise
11     """
12
13     def retrieve_prim(tnode):
14         if tnode is None:
15             return False, None
16         else:
17             ckey = tnode.key
18             if ckey == key:
19                 return True, tnode.value
20             elif key < ckey:
21                 return retrieve_prim(tnode.left)
22             else:
23                 return retrieve_prim(tnode.right)
24
25     return retrieve_prim(self.__root)

```

The other two operations, namely `insert()` and `delete()` are left as exercises.

### 23.3 The efficiency of the Table ADT implementation

The Table ADT implementation is simply a binary search tree, where the keys are used to keep the data values organized. As a result, the operations `retrieve`, `insert`, and `delete` operations have the same computational complexity as the simpler variations studied in Chapter 22. If the trees are balanced, the worst case time complexity for all three operations is  $O(\log N)$ . As a result, Tables can get extremely large, and we wouldn't expect to notice any significant runtime costs when accessing data there. Degenerate trees are still a problem, but not a problem we have to solve in a first year course!



## 24 — Object Oriented Programming

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe the difference between mutable and immutable objects.
- Explain the concept of inheritance, using a simple example.
- Explain the concept of function over-riding.
- Implement a method that over-rides `__str__()` for a simple class.

### 24.1 Mutable vs Immutable

In the previous chapter, we presented some of the basics of object oriented programming. In this chapter we will be exploring some of the implications of the paradigm. It is not the purpose of this course to teach everything there is to know about objects and classes. There are other courses that cover these matters very thoroughly. But there are a few concepts which will be helpful for your work in this course.

Consider the following class definition:

```
class Rectangle(object):  
    def __init__(self, width, length, x, y):  
        self.length = length  
        self.width = width  
        self.x = x  
        self.y = y
```

The initialization method creates and initializes four public attributes. Because they are public, any application script can access these attributes using dot-notation, as follows:

```
1 small = Rectangle(5, 5, 0, 0)      # 5x5 at (0,0)
```

```

2 big = Rectangle(10, 10, 0, 0)      # 10x10 at (0,0)
3
4 # move them
5 small.x += 10
6 small.y += -5
7
8 big.x = 15
9 big.y = 0

```

The application script creates two instances of the Rectangle class, and then, by accessing the public attributes directly, moves the rectangles to different locations. The fact that the attributes are public means that they can be changed by any script. This is convenient, but leaves the object open to certain kinds of errors. For example, lines 5-6 move the small rectangle *by* a certain amount, whereas lines 8-9 move the bigger rectangle *to* a specific place. The difference between += and = is pretty small, and could be prone to error. One way to improve the situation is to make the attributes private, and provide setters and getters, as follows:

```

class Rectangle(object):
    def __init__(self, width, length, x, y):
        self.__length = length
        self.__width = width
        self.__x = x
        self.__y = y

    def move_by(self, by_x, by_y):
        self.__x += by_x
        self.__y += by_y

    def move_to(self, to_x, to_y):
        self.__x = to_x
        self.__y = to_y

```

To keep the example brief, we've omitted the doc-strings. Now, if you want to move the rectangle, you need to use one of the methods, because the attributes are private, and cannot be modified by an application script.

```

small = Rectangle(5, 5, 0, 0)      # 5x5 at (0,0)
big = Rectangle(10, 10, 0, 0)     # 10x10 at (0,0)

# move them
small.move_by(10, -5)

big.move_to(15, 0)

```

This has the advantage of making the intention very clear, as well as being less prone to errors.

The two methods added to the definition change the object's attributes. Thus we call the object mutable. A *mutable* object has at least one public attribute that can be changed, or at least one

method that changes the object's attributes. An *immutable* object has no public attributes, and no methods at all that change the object's attributes.

For example, consider the following class definition:

```
class Currency(object):
    def __init__(self, dollars, cents):
        self.__dollars = dollars
        self.__cents = cents
```

This definition is short, but there are no public attributes, and no methods to change the value of any attribute. Objects of the Currency class represent some amount of money, stored in two private attributes. Because they are private, they cannot be accessed outside the class definition.

A class like this may seem to have limited value, because once the objects are created, they cannot be changed:

```
movie_price = Currency(10, 50)
pop_corn_price = Currency(7, 95)
drink_price = Currency(5, 95)
```

We will design a few methods to make the Currency class somewhat more useful, but we will insist that no one is able to make changes to the value stored in an existing Currency object. In other words, we want Currency objects to be immutable.

Suppose we want to take two Currency objects and add them together to get a combined value; e.g., maybe we want to add the cost of popcorn and a soft drink. One way to do this is to write a method called `add()` in the Currency class as follows:

```
def add(self, other):
    """
    Add the other Currency object to this one.
    :param other: A Currency object
    :return: (none)
    """
    cents = (100 * (self.__dollars + other.__dollars)
             + (self.__cents + other.__cents))
    self.__dollars = cents//100
    self.__cents = cents % 100
```

This example shows that objects can be passed to other methods. In this case, both `self` and `other` are Currency objects, and the `add()` method is part of the Currency class, so all the private attributes are accessible here. However, this method violates our design decision that a currency object will not change its value. The method also makes a fairly arbitrary choice in deciding that `self` should be the one to change, and not `other`. Why not the other way around? There is no good reason to choose either, and some programmer is going to have a bug in their code because they forget which one changes. It's a bad design.

Alternatively, we could design `add()` to return a new Currency object, whose value is the sum of the two currency objects, as follows:

```
def add(self, other):
    """
```

```

    Add two Currency values together, producing a new
    Currency value
    :param other: A Currency object
    :return: A Currency object
    """
    cents = 100 * (self.__dollars + other.__dollars) \
        + (self.__cents + other.__cents)
    result = Currency(cents//100, cents % 100)
    return result

```

A new Currency object is being created to contain the sum, and it is being returned as the answer. It's not doing anything more complicated than the previous version. We could use it in a script as follows:

```

movie_price = Currency(10, 50)
pop_corn_price = Currency(7, 95)
drink_price = Currency(5, 95)

sub_total = Currency(0, 0)
sub_total = sub_total.add(movie_price)
sub_total = sub_total.add(pop_corn_price)
sub_total = sub_total.add(drink_price)

```

While it looks like the variable `sub_total` refers to an object that is changing, a careful understanding of the `add()` method reveals that every time `add()` is called, a new Currency object is created, and a reference to a new Currency object replaces the previous reference to an older Currency object. You could imagine a number of other similar methods, one for subtraction, and one for multiplication, that also return new Currency objects. As long as no method actually changes any attributes, the object can be considered immutable. There is no magic for immutable objects; just a different pattern for how they can be used.

All of the immutable objects we know about are like our Currency class: strings, numbers, and tuples have no methods to change the data stored inside; they can only be used in expressions that give rise to new objects. For example, we can't instruct Python to change the number 3. When we write `3 + 1` we are not changing 3 or 1; we are creating a new number with the value 4. We can use the value 4 in an expression, or we can save it by using an assignment statement.

## 24.2 Inheritance: Giving several classes some common properties

We started the chapter with a simple example of the Rectangle class. A rectangle instance has a position  $(x, y)$ , and a height  $h$  and a width  $w$ . But if you know anything about geometry, you know that there are many other shapes we could want to define: squares, and circles, for example. If we want to use squares and circles in some application, they will also need position and size information. In fact, any shape we can think of will need a position; they have this in common. All shapes have a size, too, but we usually describe circles and squares differently. A circle has a radius, and a square has a side-length.

It would be trivial to create two new classes, one called Circle, and the other called Square, which is more or less copy/paste from the Rectangle class that we saw earlier.

```
class Rectangle(object):
    def __init__(self, width, length, x, y):
        self._length = length
        self._width = width
        self._x = x
        self._y = y

    def move_by(self, by_x, by_y):
        self._x += by_x
        self._y += by_y

    def move_to(self, to_x, to_y):
        self._x = to_x
        self._y = to_y

class Square(object):
    def __init__(self, side, x, y):
        self._side = side
        self._x = x
        self._y = y

    def move_by(self, by_x, by_y):
        self._x += by_x
        self._y += by_y

    def move_to(self, to_x, to_y):
        self._x = to_x
        self._y = to_y

class Circle(object):
    def __init__(self, radius, x, y):
        self._radius = radius
        self._x = x
        self._y = y

    def move_by(self, by_x, by_y):
        self._x += by_x
        self._y += by_y

    def move_to(self, to_x, to_y):
        self._x = to_x
        self._y = to_y
```

This is not a good implementation, because the methods are all very similar. Consider how we would design the Hexagon class: by copy/pasting the code from Rectangle. The problem with copy/paste



is that it duplicates any undiscovered bug in the code. If there's any reason to add or change the behaviour of a duplicated method, the modification has to be copy/pasted as well. Copy/pasting code is simply a bad implementation tool, and software designers have been aware of this problem for decades.

Python (and every other object-oriented language) provides tools that allow programmers to express the idea that some things are common to a group of classes, while other things will be different. The key is to give a name to the common properties, and to give them their own class. For example, we stated above that all shapes have a position. The Shape class below captures the idea of position:

```
class Shape(object):
    def __init__(self, x, y):
        self._x = x
        self._y = y

    def move_by(self, by_x, by_y):
        self._x += by_x
        self._y += by_y

    def move_to(self, to_x, to_y):
        self._x = to_x
        self._y = to_y
```

A Shape knows about position, but not about size. Its methods `move_to()` and `move_by()` give us the tools to move shapes around.

Python (and every other object-oriented language) also lets us indicate that a Rectangle is a kind of a Shape:

```
1 class Rectangle(Shape):
2     def __init__(self, width, length, x, y):
3         Shape.__init__(self, x, y)
4         self.length = length
5         self.width = width
```

On the first line of the class definition for Rectangle, we see the name Shape being used instead of `object`. This is how we tell Python that a Rectangle is a kind of Shape. Specifically, it means that a Rectangle object has all the properties of a Shape as defined above: it has the two Shape attributes, `self.x` and `self.y`, and the two Shape methods `move_to()` and `move_by()`. The `__init__()` method for Rectangle has 5 parameters, and on line 3, it calls the `__init__()` method for Shape, allowing Shape to initialize the Shape attributes for the object. This is unlike most of the method calls we have seen because it mentions the class name, Shape explicitly, to avoid ambiguity.

The same technique can be used to say that a Square is a kind of Rectangle, specifically, a rectangle where the length and width are equal:

```
1 class Square(Rectangle):
2     def __init__(self, side, x, y):
3         Rectangle.__init__(self, side, side, x, y)
```

Square's `__init__()` method does nothing else but instruct Rectangle to initialize its attributes. Rectangle initializes `self.width` and `self.length` to the same value. and Rectangle instructs Shape to initialize `self.x` and `self.y`, as before.

We say that Square is *a kind of* Rectangle, or that Square is a *sub-class* of Rectangle. We say that Rectangle is the *super-class* for Square, or that Rectangle is the *parent class* for Square. Likewise, Rectangle is a sub-class of Shape, and Shape is the parent class of Rectangle. We should point out that this extends to the class `object`: Shape is a sub-class of `object` and `object` is the parent class of Shape. This language is not accidentally similar to the terminology we use for trees. Our example has only shown a sequence of parent-child relationships from `object` to Square, but Currency is sub-class of `object` as well, making it a sibling, so to speak, of Shape. All classes that we have seen, including lists, tuples, dictionaries, strings, numbers and Booleans, are part of this tree of classes, and the root of that tree is the `object` class. It's not a binary tree; the `object` class has many sub-classes. The tree of classes is called a class hierarchy.

Each class inherits all the methods of the classes in its ancestry (up the hierarchy to the root). When each class calls the `__init__()` method of its parent class, it inherits all the attributes of the classes in its ancestry. As a result, Square and Rectangle have `move_to()` and `move_by()` methods, and attributes `self.x` and `self.y`.

We can define the Circle class to be a sub-class of Shape as well.

```
1 class Circle(Shape):
2     def __init__(self, radius, x, y):
3         Shape.__init__(self, x, y)
4         self.radius = radius
```

It inherits attributes and methods from Shape, and methods from `object`.

We can give each of our classes a method to calculate area. For Square and Rectangle, the method is the same; for Circle, we need to use a different calculation. In the definition for class Rectangle, we'd add the method:

```
1     # in Rectangle
2     def area(self):
3         return self.length * self.width
```

and in Circle, we'd add a method of the same name, but with a different calculation:

```
1     # in Circle
2     def area(self):
3         return Math.pi * self.radius ** 2
```

This is how we can give different properties to classes that have the same parent class. It's very interesting to note that we don't have to give Square an area method at all. Square will inherit Rectangle's area method, and that's exactly correct for Square, because length and width are the same.

When we call an object method, for example `sq.area()`, Python looks for the `area()` method in the object's class. If the method is defined there, Python uses it. But if the method is not defined by the class, Python checks the object's ancestor classes, one by one, until it finds one (good), or until it reaches the root of the class hierarchy, and cannot find a method with that name (a runtime error). The first one Python finds is the one that gets called.

### 24.3 Over-riding a method

As it stands, our Currency objects may be correctly summing up, but we have no way to see what their value is. We can't simply use `print()` to print out an object, and we can't access the private attributes outside of the class definition. But we could write a method to construct a string representing the value of a currency object:

```
def to_string(self):  
    return '$'+str(self.__dollars)+'.'+str(self.__cents)
```

The function doesn't display anything. It simply creates a string that could be displayed, as follows:

```
print('Before taxes:', sub_total.to_string())
```

This is a very simple solution to a common problem.

You might wonder why we couldn't simply use Python's `str()` function to display a Currency object. It turns out that the answer is that we can, but the result is not what we want. By default, the `str()` function will create a string reflecting a Currency object's identity (i.e., its address on the heap), rather than its value (a dollar amount). For example, `str(sub_total)` would result in a string that looks like '<Currency object at 0x102d2c4e0>'. The `str()` function currently does not know enough about the Currency object to do any better.

We can do better by changing the name of our `to_string()` method, giving it the name `__str__()` instead. Everything else looks the same:

```
def __str__(self):  
    return '$'+str(self.__dollars)+'.'+str(self.__cents)
```

Now, when we call Python's `str()` function, Python knows to use Currency's `__str__()` method:

```
print('Before taxes:', str(sub_total))
```

The `object` class has a definition for `__str__()`; this method returns a bland and unhelpful return string like '<object at 0x102d2c4e0>'. It is the best `object` can do without help from the programmer. But any class in the hierarchy can define its own version of `__str__()`, which can create a useful string based on the attributes of the class instances. That's what we did above with Currency. It has the same name as the method in `object`, but it has a different behaviour. We say that the `__str__()` method in the Currency class *over-rides* the `__str__()` method in the `object` class.

Python's `str()` function looks something like this:

```
def str(obj):  
    return obj.__str__()
```

In other words, when you call `str(total)`, Python translates that into call to the `__str__()` method in `total`'s class. If `total`'s class has a `__str__()` method, that's the one that gets used. If not, Python checks the parent class for a `__str__()` method, and up the hierarchy as described in the previous section. Since `object` does have a `__str__()` method, we always get something, even if it is not too informative.

The `object` class has dozens of methods like `__str__()` that make programming with objects and classes easier. For example, we could over-ride the `__add__()` method, by changing the name

of `Currency`'s `add()` to `__add__()`. With this new name in place, we can use normal arithmetic syntax on `Currency` objects:

```
if __name__ == '__main__':  
  
    movie_price = Currency(10, 50)  
    pop_corn_price = Currency(7, 95)  
    drink_price = Currency(5, 95)  
  
    sub_total = Currency(0, 0)  
    sub_total = sub_total + movie_price  
    sub_total = sub_total + pop_corn_price  
    sub_total = sub_total + drink_price
```

This works because Python translates expressions using `+` into method calls using `__add__()`, much the way it translates the function `str()` into a method call for `__str__()`.

## 24.4 Summary

We have learned that data can be encapsulated in objects along with the methods that operate on the data. We have seen that objects can be mutable, if some method allows the data to be modified. We can create immutable objects by ensuring that none of the object's methods is allowed to change the data; operations on immutable objects create new class instances containing the result of the operations. We discussed class hierarchy in terms of sub-classes and super-classes, and how a class inherits methods and attributes from its parent class. We have seen that some common Python functions like `str()` can be applied to instances of classes we define by over-riding `object` methods.



## 25 — Three Kinds of Tasks

### Learning Objectives

After studying this chapter, a student should be able to:

- To distinguish between search tasks, decision tasks, and optimization tasks.
- To give examples of search tasks, decision tasks, and optimization tasks.

### 25.1 Introduction

In computer science, and software development, we face many kinds of problems. The fundamental problem, from which all the others seem to flow, is that we want our computers to perform a specific task. In other words, the fundamental problem in computer science is to develop algorithms. An algorithm is a sequence of instructions that accomplish a stated task. We've seen a fair number of algorithms so far. Some are fairly mundane, like calculating an average value for some collection of data, or finding the  $n$ th Fibonacci number. Others are a bit more interesting, like simulating an MM1 queue, or building a Huffman tree.

Before we start our study of algorithms, it is useful to provide a loose classification of the kinds of tasks that we typically need algorithms to do. These are broad categorizations with somewhat vague descriptions, because a lot of different examples fit into them.

The first is the class of *search tasks*. A search task is any task in which the objective is the discovery or construction of a thing (object or value) that satisfies a given set of requirements. Some example search tasks:

- Calculating the square root of a given number. The value we need has the requirement that if it is squared (i.e., multiplied by itself), the result is equal to the given number.
- Determining the solution of a set of equations involving a set of mathematical variables. The thing we need as a result is a set of numerical values, one for each variable, which makes the given equations true.

- Planning an itinerary for a trip from Saskatoon to Hong Kong. The thing we want to have at the end is a plan of travel. This plan must be feasible, in the sense that each step of the itinerary is a real travel service (e.g., a real flight on a real airline), and each step in the sequence completes before the next step is supposed to start.

The thing about search tasks is that the answer can be checked, or *certified*, just by looking at the answer and the requirements.

The second kind of task is a *decision task*. A decision task is any task in which the objective is to determine if a set of requirements can be satisfied. A decision task is usually phrased as a yes/no or true/false question. Some examples of decision tasks:

- Does a given number have a real-valued square root?
- Does a set of equations involving a set of mathematical variables have a solution?
- Is it possible to travel from Saskatoon to Hong Kong?
- Is there a time this week that everyone on the project can meet for an hour?

Decision tasks are closely related to search tasks. Sometimes, the answer to a decision task is to perform the related search task, and answer affirmatively if the search task has a certified answer. Sometimes, the decision task is a lot easier; for example, by checking the sign of a number, we can determine if it has real roots without having to find the roots first.

The third kind of task is an *optimization task*. An optimization task is a task in which the objective is to find the best, according to a given set of criteria, out of all the possible things that satisfy a set of requirements. Some examples of optimization tasks:

- Find the largest (or smallest) number in a collection (list or tree).
- Find an itinerary for travel from Saskatoon to Hong Kong that has the lowest dollar cost.
- Find a code for text data that results in the highest compression rate of all possible codes.
- Find a one-hour meeting time this week at which the largest number of project members can attend.

An optimization task is a special kind of search task, in that we're trying to find a thing (object or value) that satisfies some requirements. An optimization task is distinct from a search task in that the answer to an optimization task cannot be checked by looking at the answer and the requirements; it also requires checking that no other possible answer is better. For example, the square root task (a search task), we can determine if a value is correct by checking if squaring it gives us the right answer. In the task of finding the maximum value from a list (an optimization task), we can't tell if it's the largest value just by looking at it; we have to know something about the other values in the list as well.

## 25.2 Example Problems

In this section we'll describe several example problems, describing their requirements, and giving variations according to the classification scheme discussed in the previous section. Algorithms to solve these problems will be discussed in the next chapter.

As you read through the examples that follow, keep in mind that these are tasks that can be stated fairly simply, and only make use of concepts that we have already studied. As you learn more about computer science, you will encounter a much larger variety!

Also keep in mind the idea of abstraction: the tasks below involve lists and trees and quantities that may not seem to have any hook into a real, practical problem. It would be a mistake for you to



conclude that the tasks have no practical value, simply because they have been stated abstractly. It is because a task can be stated abstractly that a solution to the abstract problem can be applied to any number of real, practical problems. This is a skill that all computer scientists develop: reducing a real problem to an abstract task that represents the essential aspects. Once the trappings of a practical application have been removed, and the essential aspects are revealed, it is very common to discover that the task is well-understood, with well-known algorithms for solving it. Alternatively, it might be similar to a known task, so reading and studying the similar task may suggest a solution. Finally, if the task is not already well-studied, the effort to understand it, and solve it is considered a contribution to science.

### 25.2.1 Subset Sum

Given a list of numbers,  $L$ , and a numeric target value  $T$ , find a list of numbers  $M$ , taken from  $L$ , whose sum is exactly  $T$ . For example, let  $L = [1, 3, 5, 7]$ , and let  $T = 8$ . One solution is  $M = [1, 7]$ .

The subset sum problem is a search problem, because we can check any postulated answer  $M$  by adding the numbers up and comparing it to  $T$ .

### 25.2.2 Maximum Slice

A slice of a list is any contiguous sequence of numbers in the list, starting at index  $a$  and including all numbers at index  $a$  up to but not including index  $b$ . This is analogous to Python's notion of a slice.

Given a list of numbers,  $L$ , find the slice from  $a$  to  $b$  that has the largest sum of all possible slices of  $L$ . For example, let  $L = [1, -2, 3, 4, -5]$ . The slice from index 2 up to but not including index 4 has a sum equal to 7, which is the largest of all possible slices. In the maximum slice problem, we want the maximum sum, and typically, the range  $a$  to  $b$  is not part of the answer.

The maximum slice problem is an optimization problem. We can verify that a given slice has a given sum, we cannot know it is the maximum sum just by looking at the slice.

### 25.2.3 Making Change

Given  $D$ , an integer in the range 0 through 99, and a list  $L$  of coin values, find a list of integers  $C$ , indicating how many of each coin value are needed to have the value of  $D$  exactly.

For example, let  $D = 37$ , and let  $L = [1, 5, 10, 25]$ , i.e., the coins are pennies (no longer in circulation in Canada), nickel, dimes, and quarters. A solution is  $C = [2, 2, 0, 1]$ , meaning: 2 pennies, 2 nickels, 0 dimes, and 1 quarter. A quick calculation proves that the solution has monetary value equal to 37.

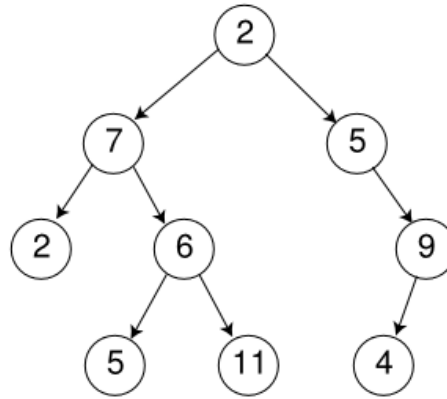
As stated, the problem is a search task. Given a list of the numbers of coins, and their values, we can check if the value of those coins is  $D$ . We can turn it into an optimization task by asking for the smallest number of coins equal to the value  $D$ . In our example above, the solution with the smallest number of coins replaces the 2 nickels with 1 dime:  $C = [2, 0, 1, 1]$ .

### 25.2.4 Maximum Tree Path

Given a binary tree  $T$  with integer data values, and a target value  $V$ , find a path from the root of  $T$  to a leaf node on which the sum of the data values is at least  $V$ . For example, given the tree in Figure 25.1, and  $V = 10$ , the path that takes the left branch every time has data values that add up to 11.

As stated, the problem is a search problem: it's easy to verify that the path sum satisfies the requirement. A variant of the task can be stated as a decision problem: is there a path whose sum is

Figure 25.1: A binary tree with integer values.



greater than  $V$ ? A related optimization task would be to ask for the path in the tree whose sum is the largest over all paths from root to leaf. In this variant, we are given the tree  $T$ , but not a target value.

Another search variant of the Maximum Tree Path problem is to find a path from root to leaf whose sum is not greater than  $V$ .

### 25.2.5 Leap Line

This task is inspired by side-scroller video games, but is highly simplified. Consider a game avatar controlled by a player, who can make the avatar to run and jump in one direction along a track. The track has rewards (e.g., coins), and penalties (e.g., mushrooms) placed uniformly, and the goal is to collect the highest total score by running through the coins to collect them, and jumping over mushrooms to avoid penalties. The avatar cannot jump past the end of the track.

We can create an abstract representation of this game by using integer values in a list  $L$ : positive values for coins, and negative values for mushrooms. The avatar's path through the game is decided by having a choice: either moving from index  $i$  to index  $i + 1$  (a *step*), and adding the value  $L[i + 1]$  to the player's score, or moving from index  $i$  to index  $i + 2$  (a *jump* of size 2), collecting only the value  $L[i + 2]$ , avoiding the item at index  $i + 1$ . The Leap Line task is to determine the maximum possible score for the given list of integers  $L$ .

The task can be generalized a little more, if we change the distance that the avatar can jump. In general, a step takes the avatar from index  $i$  to index  $i + 1$ , and a jump takes the avatar from index  $i$  to index  $i + k$ , where  $k$  is a positive integer.

## Introduction

### Algorithm styles

Brute Force Algorithms

Backtracking Algorithms

Divide and Conquer Algorithms

Greedy Algorithms

Dynamic Programming Algorithms

### How to solve it

## 26 — Algorithms

### Learning Objectives

After studying this chapter, a student should be able to:

- For each of the strategies (brute-force, greedy, divide-and-conquer, recursive backtracking, and dynamic programming), identify a practical example to which it would apply.
- Use a greedy approach to solve an appropriate problem and determine if the greedy rule chosen leads to an optimal solution.
- Use a divide-and-conquer algorithm to solve an appropriate problem.
- Use recursive backtracking to solve a problem such as navigating a maze.

### 26.1 Introduction

In this chapter we'll describe some basic algorithm styles. By studying these styles, you will have a terminology to help you describe algorithms, and also you will have a foot-hold on some general approaches that will help you develop algorithms. These algorithm styles apply to all the task categories presented in the previous chapter.

Designing algorithms is not the same as designing software. We do not insist on robustness, adaptability, or reusability when designing algorithms, even though these are essential for good software. Algorithms may be interesting for their own sake; we might design an algorithm to prove that something can be done, or to provide the basis of some crucial component of software. An algorithm may eventually become part of an application, but the work on designing the algorithm would focus mostly on how to accomplish the task, and not at all on how to turn it into quality software. Once the algorithm has been designed, we can turn it into quality software using the concepts we've studied in previous chapters.

The study of algorithms does not require the use of computers at all, though undoubtedly they

can be quite useful. Some of the most fundamental work on algorithms was done by computer scientists before computers were practical or common. The study of algorithms is the pure science aspect of computer science.

## 26.2 Algorithm styles

### 26.2.1 Brute Force Algorithms

A *brute force* algorithm solves a search task or an optimization task by trying every possibility. To design a brute force algorithm, you have to figure out what the possibilities are, and how to work through them systematically. The easiest brute force algorithm that we've seen is linear search through a list: stepping through the possibilities is as simple as writing a for-loop.

Other tasks from the previous chapter are not as easy. For the Making Change problem, the possibilities are the set of all possible combinations of any number of coins. For the Subset Sum task, the possibilities are every possible sub-list of the list  $L$ . For the Maximum Tree Path, the possibilities are every path in the tree from root to leaf.

A brute force algorithm does not try to explore the possibilities in a clever way (hence *brute*, implying a deficit of cleverness). However, a brute force algorithm will stop as soon as one of the possibilities satisfies the task's requirements.

The only *force* that a brute force algorithm can apply is the speed of the computer. If the number of possibilities is very large, the time complexity of a brute force algorithm might be so high as to render it useless for anything but small examples. A brute force algorithm is usually the simplest to design for any given problem, and writing one can be useful for deepening your understanding of the task you are trying to accomplish.

### 26.2.2 Backtracking Algorithms

A *backtracking* algorithm solves a search task or optimization task by organizing the exploration of possibilities using recursion. Typically, a backtracking algorithm explores the space of possibilities by trying all the choices recursively. As a result, a backtracking algorithm may be a brute force algorithm (but not always). For example, a brute force algorithm for the Maximum Tree Path search task (find a path from root to leaf in  $T$  whose sum is at least  $V$ ), would explore the tree using recursion, because trees are inherently recursive.

Backtracking algorithms are not always brute force. Depending on the task, a backtracking algorithm might be able to avoid a very large number of possibilities by checking whether deeper recursion can possibly lead to a solution or not. If it can be determined that a particular choice can not lead to a solution, that choice can be eliminated before it is tried, and more time can be spent exploring other more feasible choices.

### 26.2.3 Divide and Conquer Algorithms

A *divide and conquer* algorithm solves a task by dividing the task into independent sub-tasks. The independent tasks are usually solved recursively, and the solution to the original task is constructed from the solutions to the sub-tasks. For example, the Quick Sort algorithm is a divide and conquer algorithm: it separates a list into sub-lists, sorts them, and puts the pieces together. A divide and conquer approach is valuable when the computational complexity of solving the sub-tasks is much smaller than the cost of solving the whole problem.

Divide and conquer algorithms are frequently designed as recursive functions. Normal recursive functions that break a task of size  $N$  into one task of size  $N - 1$  and another of size 1 are also divide and conquer algorithms. Typically, though, the subtasks for a divide and conquer algorithm are substantially smaller than the given task, say two problems of size  $N/2$ .

### 26.2.4 Greedy Algorithms

A *greedy* algorithm solves a search task by avoiding the exploration of all the possibilities. To avoid looking at all the possibilities, a greedy algorithm tries to make good choices about the task based on the information that the algorithm has available at the time. The word “tries” in the previous sentence is important, because sometimes, the information that the algorithm has is too limited to allow the algorithm to make the right or perfect choice. But sometimes, the algorithm has all the information it needs to make the right choice.

The binary search algorithm is a greedy algorithm. It decides where to look for a given element of a list based on the value at the midpoint of the list. If the list is sorted, the limited information guarantees that binary search will return the correct answer. But if the list is not sorted, the decision to search half of the list only will almost certainly not return the correct answer.

While binary search is familiar, it is not typical of greedy algorithms. More typically, greedy algorithms construct a solution to a search task or an optimization task piece by piece. For example, we can solve the Make Change task by deciding sequentially how many quarters to use, followed by dimes, nickels and pennies. It’s greedy because a decision about the answer is made without exploring the possibilities. This particular greedy algorithm makes every choice correctly, and so it can solve the optimization version of this task without trying all possibilities.

A greedy algorithm is also used when exploring all possibilities is completely infeasible, and an answer, even if not the right or best answer, is better than no answer at all. Many algorithms in Artificial Intelligence and Machine Learning are greedy algorithms of this sort.

### 26.2.5 Dynamic Programming Algorithms

A *dynamic programming* algorithm is a divide and conquer algorithm combined with the technique of memoization. This kind of algorithm is highly valuable if the divide and conquer algorithm would require the re-solving of the same sub-tasks repeatedly, and the memo allows the algorithm to save time by looking up the answer instead of re-computing it.

Consider for example, the Leap Line problem. It is fairly easy to implement a backtracking algorithm that tries every possible combination of stepping and jumping. If we simply implement a memo (i.e., a dictionary that remembers what calculations have already been made), we can avoid a lot of repeated calculations.

The simplest dynamic programming algorithms can be described as being divide and conquer algorithms supplemented by a memo. Many dynamic programming algorithms are not written using recursion, but build up a data structure, similar to our memos, using loops instead of recursion. This is usually done when recursion would incur too high a cost in terms of its use of the system stack, but frequently results in algorithms that are much more difficult to design and understand afterward.

The memoization incurs a memory cost, in exchange for saving runtime costs. This trade-off usually leads to a very effective solution, but one must always analyze these costs to be sure.

### 26.3 How to solve it

If you are faced with a task that you don't have an algorithm to solve, you might wonder where you should start. This section tries to give some advice, in the form of a sequence of things to try.

1. Try to identify the essential aspects of the problem, independent of the application. You're no doubt trying to solve a problem with a practical application, but it usually pays to strip away the details, and look at the essential aspects. At this point, you may recognize your task as being related to a task you already know how to solve.
2. Determine if your task is a search task, decision task, or optimization task. This will help you focus on what your algorithm needs to do.
3. Design a brute force algorithm for your task. This will help you understand what a solution is, and can be useful for small examples.
4. Try a divide and conquer approach.
5. Design a dynamic programming algorithm, by simply adding a memo to your divide and conquer algorithm, if it does a lot of repeated calculations.
6. Design a greedy algorithm that makes simple choices without trying all possibilities. This might be good enough.
7. Design a backtracking algorithm. It might be brute force, or you might be able to save some work by cutting off fruitless possibilities before they are fully explored.



#### Introduction

#### Subset Sum

Brute Force  
Divide-and-conquer

#### Maximum Slice

Brute Force  
Divide and conquer

#### Make Change

Greedy

#### Maximum Tree Path

Brute force  
Greedy

#### Leap Line

Brute Force  
Dynamic Programming  
Greedy

## 27 — Algorithm Examples

### Learning Objectives

After studying this chapter, a student should be able to:

- To identify examples of different algorithm styles.
- To practice abstraction and generalization for specific tasks.
- To become aware of the kinds of mathematical tools that computer scientists use to express abstractions.

### 27.1 Introduction

The best way to improve your algorithm design skills is to study other algorithms for tasks you haven't seen before. Learning to extract the essential aspects of a task will help you recognize similarities. It's also important to learn about the kinds of mathematical tools and notation that computer scientists use to express abstractions: mathematical logic, set theory, graphs, etc. This chapter will present some algorithms for the tasks first described in Chapter 25.

### 27.2 Subset Sum

Given a list of numbers,  $L$ , and a numeric target value  $T$ , find a list of numbers  $M$ , taken from  $L$ , whose sum is exactly  $T$ . For example, let  $L = [1, 3, 5, 7]$ , and let  $T = 8$ . One solution is  $M = [1, 7]$ .

The subset sum problem is a search problem, because we can check any postulated answer  $M$  by adding the numbers up and comparing it to  $T$ .

#### 27.2.1 Brute Force

For the subset sum problem, the space of possibilities is the space of all possible sublists of the given list. A brute force algorithm has to create every possible sublist of the of the given list, and to



try them all, until one of the sublists has the target sum. This is most easily done using recursion. The internal function, `trysum()`, shown below, has a recursive case (lines 17-24) that makes two recursive calls. First (line 19), it tries to use the first element in the given list. If an answer can be found using that element, the answer is returned (lines 20-21). If not, the function is called again (line 24), this time not using the first element.

```

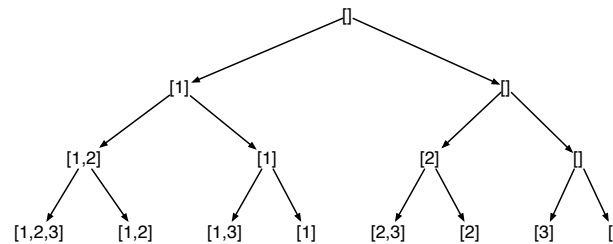
1  def subsetsum_v1(alist, target):
2      """
3      Purpose:
4          Given a list of positive integers, and a target,
5          Return a sublist of integers whose sum is target.
6      Preconditions:
7          alist: a list of positive integers
8          target: a positive integer
9      Return:
10         A tuple True, sublist if sum(sublist) == target
11         Or False, None otherwise
12     """
13
14     def trysum(al, subset):
15         """
16         Recursively try every subset of values from al.
17         """
18         if len(al) == 0 and sum(subset) == target:
19             return True, subset
20         elif len(al) == 0:
21             return False, None
22         else:
23             first = al[0]
24             rest = al[1:]
25             flag, answer = trysum(rest, subset + [first])
26             if flag:
27                 return flag, answer
28             else:
29                 # try without the first value
30                 return trysum(rest, subset)
31
32     return trysum(alist, [])

```

This technique of trying all possible subsets by including or excluding an element is fundamental to problems involving subsets or sublists. You should imagine the collection of all possible subsets, organized like a binary tree. The root of this imaginary tree is empty, since no elements have been added. The left branch is created by putting the first element into the sublist; the right branch is created by excluding the first element. All possible subsets are constructed by either including or excluding elements from the original list. Figure 27.1 shows a small example of a tree of all the subsets of the list `[1, 2, 3]`.

The worst case time complexity of the brute-force algorithm is  $O(2^N)$ , because there are  $2^N$

Figure 27.1: A tree of all the subsets of the list  $[1, 2, 3]$ , that would be created by the function `trysum()` for the Subset Sum task.



sublists of a list of size  $N$ .

### 27.2.2 Divide-and-conquer

A divide and conquer algorithm can be designed for the subset sum task. The idea is pretty simple: A recursive function named `binary()` (lines 8-24) tries to find a subset with the appropriate sum using only elements from the first half of the list (line 16). If an answer cannot be found in the first half of the list, the algorithm will try using only the second half of the list (line 20). If neither of those work, the brute-force algorithm as described above will be used on the whole list (line 24). Notice the base case of the `binary()` function. For lists of size 5 or smaller, a brute force algorithm, `trysum()` is performed.

The code is shown below. To keep the listing as short as possible, the listing does not include the `trysum()` internal function, which is called on line 24. It's identical to the function of the same name in the brute force backtracking algorithm from the previous section.

```

1  def subsetsum_v3(alist, target):
2      """
3      Purpose:
4          Given a list, display all subsets.
5      Preconditions:
6          alist: a list
7      Post-conditions:
8          displays all subsets to the console,
9          one subset per line. Could be a lot!
10     Return:
11         None
12     """
13     # this does divide and conquer, in the form of
14     # dividing the list in half
15     def binary(alist):
16         if len(alist) < 5:
17             # for small lists, just use brute force!
18             return trysum(alist, [])
19         else:
20             mid = len(alist)//2

```

```

21
22         # try to find a solution in a list half the size
23         flag, ans = binary(alist[:mid])
24         if flag: return True, ans
25
26         # if not, try the other half
27         flag, ans = binary(alist[mid + 1:])
28         if flag: return True, ans
29
30         # if not, try the whole list
31         return trysum(alist, [])
32
33     return binary(alist)

```

In the worst case, trying to find a solution in lists of half the size may result in no solution, and so the divide-and-conquer algorithm resorts to using brute force on the whole list. As a result, the worst case time complexity is the same as for the brute force algorithm above. But the algorithm tends to be very fast in a limited range of experiments. This suggests that an average case time complexity analysis may reveal that the algorithm belongs in a lower order complexity class. The average case analysis is beyond the scope of the course!

## 27.3 Maximum Slice

Given a list of numbers,  $L$ , find the slice from  $a$  to  $b$  that has the largest sum of all possible slices of  $L$ . For example, let  $L = [1, -2, 3, 4, -5]$ . The slice from index 2 up to but not including index 4 has a sum equal to 7, which is the largest of all possible slices. In the maximum slice problem, we want the maximum sum, and typically, the range  $a$  to  $b$  is not part of the answer.

The maximum slice problem is an optimization problem.

### 27.3.1 Brute Force

For the maximum slice problem, the set of all possibilities is the set of all slices. A slice can start at any element in the list. For each starting element, the slice can end at any element after the starting element. In other words, there are  $N$  slices starting at index 0,  $N - 1$  slices starting at index 1, etc. In total, there are  $O(N^2)$  possible slices of a list of size  $N$ .

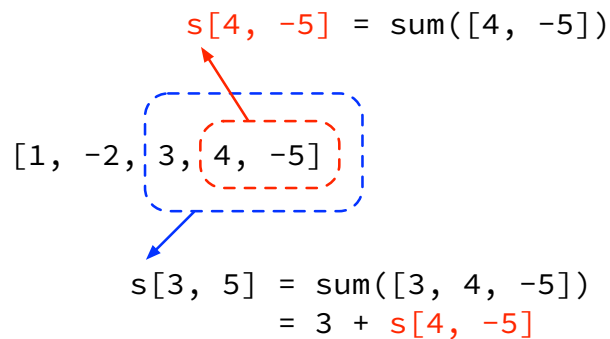
In the brute force algorithm shown below, a couple of for-loops are used to control the start and end of the slices (lines 9-13). Once the start and end of a slice is determined, the sum can be calculated using Python's `sum()` function (line 11).

```

1 def maxslice_brute_force_v0(alist):
2     """
3     Purpose:
4         Find the maximum sum of all slices of alist.
5     Preconditions:
6         alist: a list of numbers
7     Post-conditions:
8         None
9     Return:

```

Figure 27.2: An example of the use of a dictionary to store slice sums in the Maximum Slice task. Calculating the sum of the slice  $[3, 4, -5]$  can be done by storing sums of a smaller slice  $[4, -5]$  in a dictionary, and using them when needed.



```

10     a number, the maximum slice sum
11     """
12     maxsum = alist[0]
13     for i in range(len(alist)):
14         for j in range(i+1, len(alist)):
15             slice = sum(alist[i:j + 1])
16             if slice > maxsum:
17                 maxsum = slice
18     return maxsum

```

Because this is an optimization problem, the brute force algorithm has to step through every possible slice, keeping track of the largest sum it has seen so far (lines 12-13).

The worst case time complexity of this implementation is  $O(N^3)$ , because the sum of each slice (line 11) requires  $O(N)$  steps, and there are  $O(N^2)$  slices. This result should give us a little bit of a nudge. The the number of possible slices is  $O(N^2)$ , but the algorithm has worst case time complexity  $O(N^3)$ , which is a factor of  $N$  higher. We should be wondering if we could eliminate this factor of  $N$  somehow.

It turns out that we can improve this algorithm by employing an idea from dynamic programming. Consider as a concrete example what happens after computing the sum of the slice from index 2 to index 4. Some time later, when the starting index is increased, the brute force algorithm has to calculate the sum of the slice from index 3 to index 4. Notice that part of this slice was previously summed when calculating the sum from 2 to 4. In other words, the brute force algorithm given above recomputes a lot of sums on line 11. We can avoid recomputing any sums by storing the partial sums in a dictionary. See Figure 27.2. In the code below, a dictionary named `s` is used to store slice sums from index  $i$  to index  $j$ .

```

1 def maxslice_brute_force_v1(alist):
2     """
3     Purpose:
4     Find the maximum sum of all slices of alist.

```

```

5     Preconditions:
6         alist: a list of numbers
7     Post-conditions:
8         None
9     Return:
10        a number, the maximum slice sum
11    """
12    # using brute force: look at all possible slices
13    # but store all the partial sums in a dictionary
14    # where s[j] stores the value sum(alist[i,j+1])
15
16    maxsum = alist[0]
17    for i in range(len(alist)):
18        s = {}
19        s[i] = alist[i]
20        if s[i] > maxsum:
21            maxsum = s[i]
22        for j in range(i+1, len(alist)):
23            s[j] = s[j-1] + alist[j]
24            if s[j] > maxsum:
25                maxsum = s[j]
26    return maxsum

```

The new algorithm computes slice sums by using the dictionary *s*. We start on line 19 with a slice of size 1, namely the element at *alist[i]*, and storing it in the dictionary *s*. Line 23 creates a new sum for a slice that's larger by exactly one element, by looking up the previous sum in the dictionary, and adding exactly one more element from the list. By keeping sums in a dictionary, we can limit the cost of computing the sum of any slice to  $O(1)$ , making the time complexity of the improved brute force algorithm  $O(N^2)$ . This algorithm uses dynamic programming to save time, but it is still a brute force algorithm because it checks all the possible slices.

### 27.3.2 Divide and conquer

To apply divide and conquer to a task, it's a good idea to consider what you could do if you split your task into two roughly equal subtasks. For the Maximum Slice task, we have to imagine what could happen to the maximum slice (even though we don't know where it is). It's possible that the maximum slice is contained entirely in the first half of the list, in which case, a recursive function call with the first half of the list will find it. Similarly, it might be in the second half of the list, so a recursive function call with the second half of the list will find it. It's also possible that the maximum slice starts in the first half of the list, and extends into the second half.

The problem is that we don't know where the maximum slice is! So we will try all three possibilities. Our divide and conquer will find the maximum slice for the first half of the list, the last half of the list, and the maximum slice that extends through the middle, and it will return whichever one of these three is better.

In the code below, the function `maxslice_rec()` (lines 36-51) recursively tries to find the maximum slice in both halves (lines 46, 47), and through the middle (lines 48-49).

```
1 def maxslice_DC(alist):
2     """
3     Purpose:
4         Find the maximum sum of all slices of alist.
5     Preconditions:
6         alist: a list of numbers
7     Post-conditions:
8         None
9     Return:
10        a number, the maximum slice sum
11    """
12
13    # internal function
14    def max_tail(left, right):
15        """
16        Calculate the maximum slice that ends at right
17        (from any point starting at left or later)
18        """
19        s = alist[right]
20        maxsum = s
21        # calculate the sums from right to left (backwards)
22        for i in range(right - 1, left - 1, -1):
23            s = s + alist[i]
24            if s > maxsum:
25                maxsum = s
26        return maxsum
27
28    # internal function
29    def max_head(left, right):
30        """
31        Calculate the maximum slice that starts at left
32        (to any point up to and including right)
33        """
34        s = alist[left]
35        maxsum = s
36        for i in range(left + 1, right + 1):
37            s = s + alist[i]
38            if s > maxsum:
39                maxsum = s
40        return maxsum
41
42    # internal function
43    def maxslice_rec(left, right):
44        """
45        Recursively find maximum slice between left and right.
46        """
```

```

47     # using divide and conquer
48     if left == right:
49         return alist[left]
50     else:
51         # divide, and solve
52         mid = (right + left) // 2
53         max_left = maxslice_rec(left, mid)
54         max_right = maxslice_rec(mid + 1, right)
55         max_cross = (max_tail(left, mid)
56                     + max_head(mid + 1, right))
57         # conquer
58         return max(max_left, max_right, max_cross)
59
60     # body of maxslice_DC
61     return maxslice_rec(0, len(alist) - 1)

```

The interesting part of the algorithm is how the maximum slice is computed when it passes through the middle of the list. We've defined two internal functions to manage this. The function `max_tail()` (lines 9-21) uses a dictionary and starts from the middle and steps towards the left side of the list, storing slice sums in the dictionary. The function `max_head()` (lines 23-34) uses a dictionary and starts from the middle and steps towards the right side of the list, storing slice sums in the dictionary. A careful reading of these two functions tells us that both functions have worst case time complexity of  $O(N)$ . In other words, focussing our search on a slice that passes through the middle of the list means we're ignore very many slices; the combination of `max_tail()` and `max_head()` only looks at a very small subset of all the possible slices.

The worst case time complexity of the divide and conquer algorithm can be shown to be  $O(N \log N)$ . The course has not presented the mathematics needed to reach this conclusion, but we can wave our hands a bit. Consider an imaginary binary tree (again), starting with the whole list as the root, and each branch taking the left half or the right half of the list. We can split the list in half at most  $O(\log N)$  times, so the tree has at  $O(\log N)$  levels. One important thing to notice is that every element in the original list appears exactly once in every level of the tree. The other thing to notice is that `max_tail()` and `max_head()` are called on every level in the tree. Since these functions have time complexity  $O(N)$ , we can see that the total amount of work on each level is also  $O(N)$ . Considering there are  $O(\log N)$  levels in the imaginary tree, we get the result  $O(N \log N)$ . This kind of analysis is a bit beyond what we're hoping for in this course, and the hand-waving above is not intended to be absolutely convincing. The tools you need for this kind of analysis will be obtained from further computer science courses.

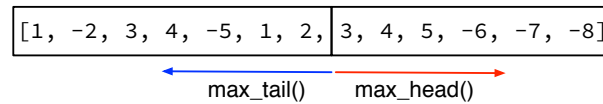
## 27.4 Make Change

Given  $D$ , an integer in the range 0 through 99, and a list  $L$  of coin values, find a list of integers  $C$ , indicating how many of each coin value are needed to have the value of  $D$  exactly.

For example, let  $D = 37$ , and let  $L = [1, 5, 10, 25]$ , i.e., the coins are pennies (no longer in circulation in Canada), nickel, dimes, and quarters. A solution is  $C = [2, 2, 0, 1]$ , meaning: 2 pennies, 2 nickels, 0 dimes, and 1 quarter. A quick calculation proves that the solution has monetary value equal to 37.



Figure 27.3: An example showing how `max_tail()` and `max_head()` work together to calculate the maximum slice that passes through the middle of a list. This technique is part of the divide and conquer algorithm for Maximum Slice.



As stated, the problem is a search task. Given a list of the numbers of coins, and their values, we can check if it the value of those coins is  $D$ . We can turn it into an optimization task by asking for the smallest number of coins equal to the value  $D$ . In our example above, the solution with the smallest number of coins replaces the 2 nickels with 1 dime:  $C = [2, 0, 1, 1]$ .

### 27.4.1 Greedy

A greedy algorithm is used to avoid systematically considering every possibility. It does so by making choices about which possibility to consider using the information it has available. A greedy algorithm for the Make Change task simply calculates how many coins of each value that it needs. In our example, we assume the use of quarters (25), dimes (10), nickels (5), and pennies (1).

The greedy algorithm considers these coins in order of decreasing value, i.e., quarters, dimes, nickels, pennies. We can calculate how many quarters we need by dividing the amount by the value of the coin, using integer division (line 13). Taking the remainder of division by the value of a coin (line 15) tells us how much remains to be accounted for by smaller coins.

```

1  def change_v2(cents):
2      """
3      Purpose:
4          Make change for the given cents value.
5          Assumes coin values 25c, 10c, 5c, 1c
6      Pre-conditions:
7          :param cents: an integer
8      Return:
9          a list of counts for the coins used.
10     """
11     coins = [25, 10, 5, 1]
12     coin_index = 0
13     counts = [0] * len(coins)
14     remaining = cents
15     while remaining > 0:
16         counts[coin_index] = remaining // coins[coin_index]
17         remaining = remaining % coins[coin_index]
18         coin_index += 1
19     return counts

```

It's interesting to consider the time complexity of this algorithm. The loop on lines 12-16 repeats the calculation once for each coin, and there are 4 coins. So the time complexity is  $O(1)$ . This

explains why even children can master this algorithm very easily. It's fast, and doesn't require a lot of memory!

You might consider how many combinations of coins there are that have a total value of  $D$ . A brute force algorithm to explore all of these possible combinations would be quite a bit more expensive in terms of time!

## 27.5 Maximum Tree Path

Given a binary tree  $T$  with integer data values, and a target value  $V$ , find a path from the root of  $T$  to a leaf node on which the sum of the data values is at least  $V$ . For example, given the tree in Figure 25.1, and  $V = 10$ , the path that takes the left branch every time has data values that add up to 11. As stated, the problem is a search problem: it's easy to verify that the path sum satisfies the requirement.

### 27.5.1 Brute force

To find the maximum sum of data values from root to leaf in a given tree, we have to reach each leaf, and check the sum on the path from the root. This is most easily done recursively, by determining the maximum sum from the current node through one of its children. This idea is implemented in the `max_path()` function below.

```

1  def max_path(atree):
2      """
3      Find the maximum sum of all paths from root to leaf
4      in atree.
5      :param atree: A TreeNode
6      :return: the maximum sum
7      """
8      if atree is None:
9          return 0
10     else:
11         left = max_path(atree.left)
12         right = max_path(atree.right)
13         return max(left, right) + atree.data

```

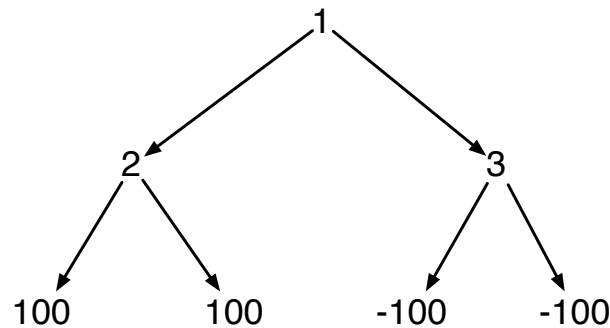
The computational complexity of this algorithm depends on the size of the tree. Each node in the tree is added to a sum exactly once, so if there are  $N$  nodes in the tree, the algorithm has  $O(N)$  time complexity. The shape of the tree doesn't matter, nor does the arrangement of the values.

### 27.5.2 Greedy

A greedy algorithm for this task serves mainly as an example of what kinds of decisions could be made to avoid systematically working through every possibility. As we noted above, the brute force algorithm for this task is  $O(N)$ , which is completely feasible. We usually would not feel obligated to explore better algorithms to avoid the  $O(N)$  cost of the brute force algorithm.

The greedy algorithm chooses a single path from root to leaf by checking the values stored at its children. The greedy choice, which illustrates the reason why it's called "greedy", is to choose the path that leads to the child with the largest value.

Figure 27.4: A binary tree with high values hidden below a node with a relatively low value. A greedy algorithm for Maximum Tree Path will not return the correct result.



```

1 def greedy_path(atree):
2     """
3     Find the maximum sum of all paths from root to leaf
4     in atree.
5     :param atree: A TreeNode
6     :return: the maximum sum
7     """
8     if atree.is_leaf():
9         return atree.data
10    elif atree.left.data > atree.right.data:
11        # the left branch seems promising
12        return atree.data + greedy_path(atree.left)
13    else:
14        # the right branch seems promising
15        return atree.data + greedy_path(atree.right)

```

The greedy algorithm can easily make a choice that is not correct. For example, the binary tree in Figure 27.4 has high value leaf nodes below the node with value 2, and low value leaf nodes below the node with value 3. The greedy algorithm will choose the right branch from the root, because 3 is bigger than 2.

This example shows that a greedy choice does not always lead to the correct answer. The information available at the time of the decision was not sufficient to guarantee the right answer.

The computational complexity of the greedy search is the length of a path from root to leaf, i.e., the height of the tree. Depending on the way the tree branches, the height of the binary tree could be as low as  $O(\log N)$ , or as high as  $O(N)$ , where  $N$  is the number of nodes in the tree. So in the worst case, the computational complexity is  $O(N)$ .

It seems important to emphasize that a greedy algorithm like this can be useful when it is infeasible to explore every possibility. The real value of this example is demonstrating how the information available to an algorithm could be used, and that it might not be enough to obtain a correct or best solution.

For some problems, such as the Make Change problem above, a greedy algorithm can produce the right answer. For this problem, it might not. This demonstrates a richness of variety in the kinds of problems we want to solve with algorithms. It also demonstrates another very important kind of trade-off: when time costs are high, we may have to compromise the quality of the answer our algorithms return. In other words, sometimes *good enough* is good enough.

## 27.6 Leap Line

We are given a list  $L$  of positive and negative integers. A path through  $L$  is decided by having a choice at every point in  $L$ : either moving from index  $i$  to index  $i + 1$  (a *step*), and adding the value  $L[i + 1]$  to the score, or moving from index  $i$  to index  $i + 2$  (a jump of size 2), collecting only the value at  $L[i + 2]$ , avoiding the item at index  $i + 1$ . The Leap Line task is to determine the maximum possible score for the given list of integers  $L$ .

### 27.6.1 Brute Force

A brute force algorithm for this task requires systematic examination of every possible path through the list. The algorithm below accomplishes this using recursion. At every index in the list, there are 2 possibilities: take a step from index  $i$  to index  $i + 1$ , or take a jump from index  $i$  to index  $i + 2$ . The best path from index  $i$  is the better of the two choices. The base case for the recursion considers the options available when the avatar is near the end of the list.

The internal function `jump_or_step()` systematically checks every path using recursion.

```

1 def maximumScoreFrom_v0(track):
2     """
3     Purpose:
4         Calculate the maximum score that can be obtained from
5         stepping and jumping along the given track
6     Pre-conditions:
7         :param track: a list of integers
8     Return:
9         the maximum score
10    """
11
12    def jump_or_step(loc):
13        """
14        Calculate the maximum score that can be obtained from
15        starting at the given location
16        """
17        step_loc = loc + 1
18        jump_loc = loc + 2
19
20        if step_loc == len(track):
21            return track[loc]
22        elif jump_loc >= len(track):
23            return track[loc] + jump_or_step(step_loc)
24        else:

```

```

25         return track[loc] + max(jump_or_step(step_loc),
26                                 jump_or_step(jump_loc))
27
28     return jump_or_step(0)

```

As with all brute force algorithms, the worst case time complexity depends on the number of possibilities. As with any recursive function that has 2 recursive calls in the recursive case, the computations can be modelled by an imaginary binary tree. Most of the nodes in the imaginary tree have 2 children, so the number of nodes in the tree is  $O(2^h)$  where  $h$  is the height of the tree. When the choice is between stepping (+1) or jumping (+2) through a list of size  $N$ , there will be leaf nodes from height  $N/2$  down to  $N$ . So  $h = O(N)$ , and thus the time complexity is  $O(2^N)$ .

### 27.6.2 Dynamic Programming

It's really instructive to see how adding a memo can help an algorithm like this. In this case, the memo keeps track of the value of the best path starting from index  $i$ . Because there is only one list, about half of all paths will step on index  $i$ , and the other half will jump over  $i$ . In any case, the brute force algorithm is expensive because it recomputes the best path starting from index  $i$  every time a possible path reaches index  $i$ . The memo avoids this repeated computation by storing the value of the best path starting from  $i$  the first time index  $i$  is encountered, and using a lookup for the value every other time.

The following code is essentially the brute force algorithm with an added memo. See line 10 for the creation of the dictionary, lines 19-20 to check to see if the value from  $i$  has already been computed once before, and 28-29 for storing the value in the memo the first time.

```

1  def maximumScoreFrom_v1(track):
2      """
3      Purpose:
4          Calculate the maximum score that can be obtained from
5          stepping and jumping along the given track
6      Pre-conditions:
7          :param track: a list of integers
8      Return:
9          the maximum score
10     """
11     def jump_or_step(loc):
12         """
13         Calculate the maximum score that can be obtained from
14         stepping and jumping along the given track, starting
15         at the given location
16         """
17         # check if the best score is already known
18         if loc in memo:
19             return memo[loc]
20
21         step_loc = loc + 1
22         jump_loc = loc + 2

```

```

23
24     if step_loc == len(track):
25         return track[loc]
26     elif jump_loc >= len(track):
27         result = track[loc] + jump_or_step(step_loc)
28         memo[loc] = result
29         return result
30     else:
31         result = track[loc] + max(jump_or_step(step_loc),
32                                   jump_or_step(jump_loc))
33         memo[loc] = result
34         return result
35
36     # body of maximumScoreFrom_v1()
37     # using memoization
38     # memo[loc] stores the best score starting from loc.
39     memo = {}
40     return jump_or_step(0)

```

With a memo, it can be shown that the dynamic programming algorithm requires  $O(N)$  time. Showing that result by algorithm analysis is beyond the expectations in this course, but we can run experiments to support this claim.

### 27.6.3 Greedy

It turns out that under certain conditions, a greedy solution can obtain the correct result for the Leap Line problem. The algorithm shown below chooses to step to  $i + 1$  or leap to  $i + 2$  depending on whether the value at  $i + 1$  is positive (step) or negative (leap). See lines 22-27.

```

1  def maximumScoreFrom_v2(track):
2      """
3      Purpose:
4          Calculate the maximum score that can be obtained from
5          stepping and jumping along the given track
6      Pre-conditions:
7          :param track: a list of integers
8      Return:
9          the maximum score
10     """
11
12     def jump_or_step(loc):
13         """
14         Calculate the maximum score that can be obtained from
15         stepping and leaping along the given track, starting
16         at the given location
17         """
18         step_loc = loc + 1

```

```
19     jump_loc = loc + 2
20
21     if step_loc == len(track):
22         return track[loc]
23     elif jump_loc >= len(track):
24         return track[loc] + leap_or_step(step_loc)
25     elif track[step_loc] < 0:
26         # greedy choice: avoid negatives
27         return track[loc] + leap_or_step(jump_loc)
28     else:
29         # greedy choice: collect 1
30         return track[loc] + leap_or_step(step_loc)
31
32     return jump_or_step(0)
```

The greedy algorithm looks only at one possibility, and since the list has  $N$  elements, the greedy algorithm has worst case time complexity of  $O(N)$ .

The greedy algorithm will obtain the correct maximum score if the list contains only 2 values, e.g., -1 and 1, and the leap is limited to  $i + 2$ . If the values in the list have a more variety, or if the leaps are longer, a greedy solution cannot find the correct maximum score. This claim will be stated without any justification, but the tools you need to do the analysis are part of future computer science courses.





## 28 — Case Study: Huffman codes

### Learning Objectives

After studying this chapter, a student should be able to:

- Describe the difference between fixed-length codes and variable-length codes.
- Describe the Huffman tree algorithm for designing variable length codes for a given message.
- Outline and explain the time complexity for the various stages of encoding and decoding using Huffman codes.
- Explain why a pair of queues makes the most efficient implementation of the Huffman tree algorithm.

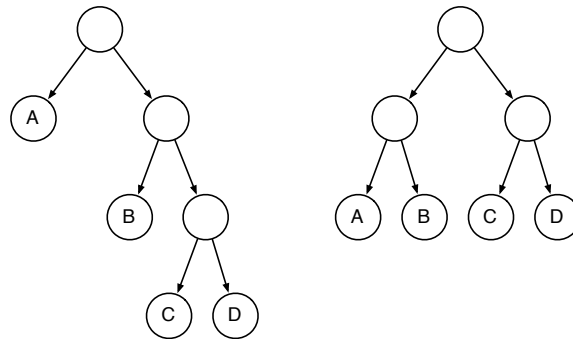
### 28.1 Introduction

In this section, we'll study an interesting application of trees, namely the construction of a code to support file compression for text files. This is an interesting case study for a number of reasons. File compression is a real application. The use of trees to support file compression is far from obvious, but easy enough to present at an introductory level. The algorithms we will study in this chapter reveal a depth of analysis and design, capturing the spirit of the field of computer science, going beyond the notion of *learning to program*.

Computers use binary digits ("bits") to represent all data, but particularly for the purposes of this chapter, text data. Numbers, letters, and symbols like + are represented inside the computer using sequences of bits. The representation amounts to being a code.

*Fixed-length* codes, like ASCII, use bit-sequences that are all the same length. ASCII codes were designed to use 7 bits to represent the typical characters and symbols you can type on a keyboard. A modern, standardized set of codes, known as Unicode, allows for 8 bit codes (e.g., UTF-8) or 16 bit codes (e.g., UTF-16). Standardized codes like these allow computer applications to exchange

Figure 28.1: Two binary trees.



information, and having a fixed-length means that computers, and computer programs can manipulate the data relatively easily.

The practical motivation for studying Huffman codes is the fact that in human languages, some characters are used more frequently than others. For example, the letter 'e' appears far more frequently in English text than the letter 'w'. Using this fact motivates the concept of a variable-length code, in which more frequently used characters have shorter codes. Variable length codes can result in more efficient transmission of text data from one computer to another. We can get the most out of this idea by designing a code that is specific to the message being transmitted: characters or symbols that do not appear in the text need not be coded.

It turns out that there is a relationship between codes and binary trees. A given code (fixed-length or variable-length), can be represented by a binary tree in which the characters or symbols to be encoded are leaf nodes (importantly, leaf nodes have no children). Since there is a unique path from the root to any leaf node, we can describe the path from root to leaf using '0' if the path takes a left branch, or '1' if the path takes a right branch. A fixed-length code has all the symbols (leaf nodes) at the same level in the binary tree. A variable-length code may have symbols (leaf nodes) at different levels of the tree.

For example, consider Figure 28.1, which has two trees. On the left, is a tree where the leaf nodes are at different levels, and on the right is a tree where the leaf nodes are at the same level. The tree on the right corresponds to a fixed-length code. Using the rule that taking a left branch corresponds to a '0', and a right branch corresponds to a '1', we can see the following codes:

00	'A'
01	'B'
10	'C'
11	'D'

Using the fixed-length code, the message 'ABACABAD' is encoded by the sequence '0001001000010011'.

Using the tree on the left corresponds to a variable-length code. Using the rule that taking a left branch corresponds to a '0', and a right branch corresponds to a '1', we can see the following codes:

0	'A'
10	'B'
110	'C'
111	'D'

Using the fixed-length code, the message 'ABACABAD' is encoded by the sequence '01001100100111'.

The variable-length code would be more efficient at representing messages with a high number of 'A' in them, and relatively few 'D'. But if the message had roughly equal numbers of letters 'A' through 'D', the fixed length code would be more efficient. A simple exercise in counting reveals that in the example above, the variable-length code needs fewer binary digits than the fixed-length code. The difference is small here, because the message was small; a larger message will demonstrate more impressive savings. As we will see, the algorithm we will study in this section will find the most efficient code given a particular message. For some messages, it will produce a fixed-length code, precisely because all the symbols with roughly the same frequency in the message.

To make the discussion very concrete, we will assume that we are given a message that we want to encode. The messages we use in our discussion will be relatively short, but in general, they could be as long as you like: essays, Wikipedia articles, DNA sequences, 100 years of temperature data. . .

Huffman codes are designed by constructing a binary tree based on the symbol frequencies in the message, so that common symbols are nearer the root than uncommon symbols. The process to build a Huffman code has two steps:

1. Determine the frequency of the symbols in the message.
2. Build a binary tree using frequency of the symbols.

We'll take a closer look at the second step a bit later. Once the tree is built, the message can be encoded, either from the tree directly, or from a table created for symbol-code pairs. What ever form is chosen, the data structure that is used to encode the message is called a *codec*.

In order to decode a message that was encoded using a Huffman code, the transmission needs to include in the message the information needed to decode the message. In other words, the encoded message has to include the codec in the transmission. When an encoded message is received, the decoding program has to read the codec, and then apply the codec to the encoded message to recover the original message.

## 28.2 Using frequency information to design a code

The first step in the encoding process is to analyze the message. The analysis has to record the symbols in the message, and count how often they are used. In Python, we could use a dictionary whose keys are symbols, and whose values are counts.

```
# analyze the frequency of characters in the message
frequencies = {}
for char in message:
    if char in frequencies:
        # seen char before
        frequencies[char] += 1
    else:
        # first time seeing char
```

```
frequencies[char] = 1
```

For example, if the message to be encoded is 'ABACABAD', we can count 4 'A', 2 'B', and 1 each of 'C' and 'D'.

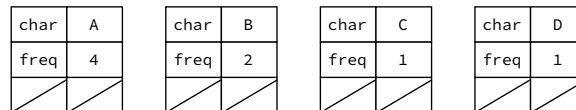
The second step is a little more interesting. We want to build a binary tree that has symbols stored in the leaf nodes, and we want to organize the tree so that uncommon symbols are deeper in the tree than more common symbols.

The Huffman algorithm for building a tree with the desired organization is simple and clever. It builds the binary tree up in stages. Initially, every symbol in the message is stored, along with its frequency, in a leaf node, i.e., a binary tree with no children, and all these leaf nodes are stored in a list.

The tree-building process continues as follows: As long as there are at least 2 trees in the list, the two trees with the lowest frequencies are selected and removed from the list, and merged. The merge process creates a new tree whose children are the two selected trees, and whose frequency is the sum of the frequencies of its children. The new tree is placed in the list with the other trees. Note that if two trees are removed from the list, and an new tree is placed in the list, the list gets smaller by one element.

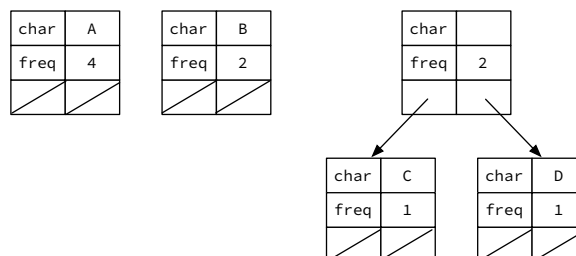
The process of selection, removal, merging, and replacing continues until there is only one tree left in the list, which is the tree we'll use to derive the code for the message. This tree is called a Huffman tree.

Using our example message 'ABACABAD', the four symbols in our message would be stored as follows:

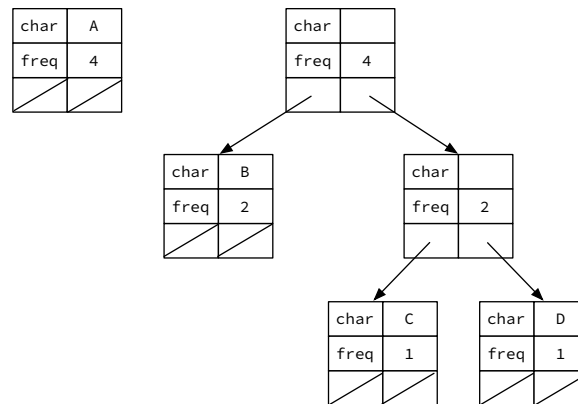


It's important to note that the frequencies are stored in the treenodes as well. Also, these Huffman trees would be stored in a list, which is not explicitly shown in the diagram.

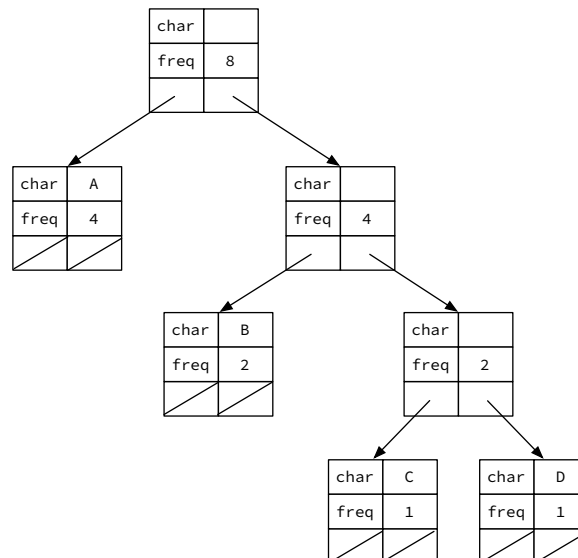
In our example, the two lowest frequency trees are the trees for 'C' and 'D'. These are removed from the list, and merged, giving them a common parent node with no character, but a combined frequency of 2:



The tree-building process continues by choosing the two lowest frequency trees in the list, namely the tree for 'B', and the tree combining 'C' and 'D'. These trees are merged, giving them a common parent node with no character but a combined frequency of 4:



Since there are only two trees left, these are combined in the same way, giving them a common parent node with no character but a combined frequency of 8:



This tree is similar in structure to one of the trees in Figure 28.1, which should not be a big surprise.

There are two places in the Huffman tree building process that involve arbitrary choices. If two or more trees have the same lowest frequency, it does not matter which one (or two) are chosen; the choice is arbitrary. As long as the new tree node has the selected nodes as children, it does not matter which is the left child, and which is the right child; the choice is arbitrary. As a result of arbitrary choices, the process will result in different Huffman trees depending on the choices made arbitrarily. In this situation, we say that Huffman trees are not unique. By this we mean that any of the trees that might have resulted are as good for the purpose of producing codes as any of the others.

## 28.3 Encoding and decoding

Once we have a Huffman tree we have several options for how to use it. Every symbol in the message appears in the tree as a leaf. A symbol's Huffman code is simply the description of the path from the

root of the Huffman tree to the leaf, using a '0' when a left branch is taken, and a '1' when a right branch is taken. Since every leaf has a unique path from root to leaf, every symbol has a unique code.

The simplest way to use a Huffman tree is to traverse the tree building up a code, i.e., a string of '0' and '1' sequences. The sequence is extended by a '0' every time the traversal descends into a left subtree, and by a '1' every time the traversal descends into a right subtree. When the traversal reaches a leaf node, the code, and the symbol can be entered into a dictionary, with the symbol as the key, and the code as the value.

Once this dictionary is constructed, the encoding process is a matter of stepping through the message, and using the dictionary to look up each symbol's code.

```
encoded = ''
for char in message:
    encoded += code[char]
```

Decoding an encoded message is a little bit trickier. We will assume that a code-symbol dictionary has been constructed; this is simply the inverse of the symbol-code dictionary we built from the Huffman tree. In a real application, the decoding dictionary will need to be constructed from the information transmitted along with the encoded message.

The reason decoding is a bit trickier is that Huffman codes are variable length. We don't know in advance which bit sequences belong to which symbols.

The process of decoding an encoded message consists of stepping through the encoded message one character at a time, building up a short sequence of binary digits until the sequence is recognized by the decoding dictionary as a code. Once a code is recognized, the corresponding symbol can be added to the decoded message.

```
message = ''
first = 0
last = first + 1
while first < len(encoded):
    tryit = encoded[first:last]
    if tryit in codes:
        message += codes[tryit]
        first = last
        last = first + 1
    else:
        last += 1
```

## 28.4 Efficiency considerations

The Huffman tree algorithms were described at a fairly high level, temporarily ignoring practical considerations. Careful implementation of these ideas will result in a fairly efficient program. Let us suppose that the message has  $k$  unique symbols, and has length  $N$ ; further, we'll assume the encoded message has length  $M$ . Counting the frequency of the symbols should require  $O(N)$  steps; we have to look at each symbol once to count it. Given a codec for the message, encoding the message requires  $O(N)$  codes to be joined together. Given a codec, decoding the message requires looking at each binary digit in the encoded message, i.e.,  $O(M)$ .



The one place where we have to work a bit harder is in the merge process. The process has to build  $k$  leaf nodes, and performs  $k - 1$  merges (because the list of trees gets smaller by one every time). The process requires us to find the two smallest trees in the list. If  $k$  is small enough, it doesn't matter how we find the two smallest trees.

But considering the costs here, even if  $k$  is not too big, leads to an interesting study. There are two obvious ways to find the smallest tree in a list of trees. First, we could perform a linear search. Repeating the linear search  $k$  times will result in a time complexity of  $O(k^2)$ . The second option is to keep the list sorted, in order of increasing frequency. With the list organized, finding the two smallest trees is easy: we can remove the first two trees in the list in  $O(1)$  time. However, we have to keep the list sorted. That could involve a linear search, putting the newly merged tree into the sorted list at the correct position to keep the list organized by increasing frequency, which requires  $O(k)$  steps. Or we could simply re-sort the list, requiring  $O(k \log k)$  steps.

There is an interesting property of the tree-merging process that we can use to our advantage. Consider what happens when we merge two trees: the result is a single tree whose frequency is the sum of the two selected trees. It turns out that at any stage of the merging process, any tree we construct by merging two other trees will have a greater frequency than any previously merged tree. Consider the sequence of merges in the example above: our merged trees had frequencies 2, 4, and 8, in that order. Because this is guaranteed to be true, we do not have to sort the trees we build, we can store them in a FIFO queue.

We can implement an efficient merge process by keeping two queues: the sorted leaf nodes, and the queue of merged trees. The smallest tree will appear at the front of one of these two queues, and the appropriate tree can be dequeued in  $O(1)$  time. The result is always a merged tree, so it should be enqueued (on the back, FIFO, remember) to the queue of merged trees, again in  $O(1)$  time. The process of merging trees will still require  $O(k)$  merges, but each merge is  $O(1)$  time! The dominating cost of this algorithm is the cost to sort the initial list of treenodes, i.e.  $O(k \log k)$  time.

A good implementation would hide the two queues behind an ADT where the dequeue operation returned the smallest tree, and the enqueue operation puts trees only on the merged list.



# **Appendices**



## A — Abstract Data Types

### A.1 The full ADT Registry implementation

The node-based queue is discussed in Chapter 8. Here we present the implementation without remark.

```
# Synopsis:
#   Defines the Registry ADT
#   A registry records a Boolean value for a range of integers.
#   It can be used to record observations or events.

class Registry(object):

    def __init__(self, size, value):
        """
        Purpose:
            Initialize a registry of a give size filled with the
            given value
        Pre-conditions:
            size: number of elements in the registry
            value: the default initial value for all elements
        """
        self.__reg = [value for i in range(size)]

    def set(self, i):
        """
        Purpose:
            Set the registry element at i to True
        Pre-conditions:
```

```

        i: an index, in the correct range for reg
    Post-conditions:
        the registry value at i is set to True
    Return:
        (none)
    """
    self.__reg[i] = True

def reset(self, i):
    """
    Purpose:
        Set the registry element at i to False
    Pre-conditions:
        i: an index, in the correct range for reg
    Post-conditions:
        the registry value at i is set to False
    Return:
        (none)
    """
    self.__reg[i] = False

def is_registered(self, i):
    """
    Purpose:
        Returns the value stored at registry element i
    Pre-conditions:
        i: an index, in the correct range for reg
    Post-conditions:
        (none)
    Return:
        the value stored at element i
    """
    return self.__reg[i]

```

## A.2 The full ADT Statistics implementation

The node-based queue is discussed in Chapter 8. Here we present the implementation without remark.

```

# Synopsis:
#   Defines the Statistics ADT
#   Calculate mean and variance and other statistical summaries
#   of numeric data.

```

```

# Implementation
# Do the calculations without storing all the data!
# Use a Python dictionary as a record to store three quantities:
#   _count':      the number of data values added
#   _avg':        the running average of the values added
#   _sumsqdiff':  the sum of the square differences between the
#                 values added and the mean so far
# These values can be modified every time a new data value is
# added, so that the mean and variance can be calculated quickly
# as needed. This approach means that we do not need to store
# the data values themselves, which could save a lot of space.

class Statistics(object):
    def __init__(self):
        """
        Purpose:
            Initialize a Statistics object instance.
        """
        self._count = 0      # how many data values seen so far
        self._avg = 0        # the running average so far
        self._sumsqdiff = 0  # the sum of the square differences

    def add(self, value):
        """
        Purpose:
            Use the given value in the calculation of mean and
            variance.
        Pre-Conditions:
            :param value: the value to be added
        Post-Conditions:
            none
        Return:
            :return none
        """
        self._count += 1
        k = self._count      # convenience
        diff = value - self._avg  # convenience
        self._avg += diff / k
        self._sumsqdiff += ((k - 1) / k) * (diff ** 2)

    def mean(self):
        """
        Purpose:
            Return the average of all the values seen so far.

```



```
    Post-conditions:
        (none)
    Return:
        The mean of the data seen so far.
        Note: if no data has been seen, 0 is returned.
        This is clearly false.
    """
    return self._avg

def var(self):
    """
    Purpose:
        Return the variance of all the values seen so far.
        (variance is the average of the squared difference
        between each value and the average of all values)
    Post-conditions:
        (none)
    Return:
        The variance of the data seen so far.
        Note: if 0 or 1 data values have been seen, 0 is
        returned. This is clearly false.
    """
    return self._sumsqdiff / self._count

def sampvar(self):
    """
    Purpose:
        Return the sample variance of all the values seen
        so far.
    Post-conditions:
        (none)
    Return:
        The sample variance of the data seen so far.
        Note: if 0 or 1 data values have been seen, 0 is
        returned. This is clearly false.
    """
    return self._sumsqdiff / (self._count - 1)
```

## B — Applications of Stacks and Queues

### B.1 Queueing simulation

The queueing simulation, based on the concept of a M/M/1 queue, is discussed in Chapter 11. Here we present the implementation without remark.

```
# Synopsis:
#   Application: M/M/1 Queueing simulation
#   (based on example by Sedgewick, Wayne, Dondero)
#
# Simulate the arrival and service of customers using a queue.
# Calculate the average wait time for customers, given
# assumptions about the rate of arrival and service.

import Queue as Q
import Statistics as S
import random as rand

# the assumptions
arrival_rate = 20/60 # the number of customer arrivals per minute
service_rate = 20/60 # the number of customers served per minute
sim_length = 8*60 # the end of the simulation in minutes

def sample_time(x):
    """
    Return a random sample time until an event.
    In the long run, the events will have a rate of x
    events per unit time.
```

```
Preconditions:
    x: the desired arrival rate, per unit time
Post-conditions:
    none
Return:
    a random sample time that obeys the rate x
'''
return rand.expovariate(x)

# customers wait in a service queue
service_queue = Q.Queue()

# we're interested in the average time in the queue
waiting = S.Statistics()

# keep track of when a new customer will arrive
nextArrival = sample_time(arrival_rate)

# keep track of when a service event is complete
nextService = nextArrival + sample_time(service_rate)

# simulate the arrival-service process
while nextService < sim_length:

    # handle all arrivals that occur before service is complete
    while nextArrival < nextService:
        service_queue.enqueue(nextArrival)
        nextArrival = nextArrival + sample_time(arrival_rate)

    # service time is over; serve first customer in the queue
    this_arrival = service_queue.dequeue()

    # how long was the customer waiting?
    waited = nextService - this_arrival
    waiting.add(waited)

    # determine when service to the next customer will end
    if service_queue.is_empty():
        nextService = nextArrival + sample_time(service_rate)
    else:
        nextService = nextService + sample_time(service_rate)

print('Average wait time:', waiting.mean())
```

## B.2 Bracket checking

The bracket-checking program is discussed in Chapter 11. Here we present the implementation without remark.

```
# Synopsis:
#   Application: Bracket Matching Algorithm
#
# Check a string for matching brackets
# Uses the Queue and Stack ADTs

import Queue as Q
import Stack as S

example = ' ( ( ( ( ) ) ) ) '

# create the initial empty containers
chars    = Q.Queue()
brackets = S.Stack()
unmatched_close = False

# put all the characters in the Queue
for c in example:
    chars.enqueue(c)

# brackets match if and only if every '(' has
# a corresponding ')'
while not chars.is_empty() and not unmatched_close:
    c = chars.dequeue()
    if c == '(':
        brackets.push(c)
    elif c == ')' and not brackets.is_empty():
        brackets.pop()
    elif c == ')' and brackets.is_empty():
        unmatched_close = True
    else:
        pass

# check how the analysis turned out
if unmatched_close:
    print("Found a ')' with no matching '('")
elif not brackets.is_empty():
    print("At least one '(' without a matching ')")
else:
```

```
print('Brackets matched')
```

### B.3 Post-fix Arithmetic evaluator

The post-fix arithmetic evaluator is discussed in Chapter 11. Here we present the implementation without remark.

```
# Synopsis:
#   Post-fix arithmetic evaluation

# An application of the Stack ADT
# and the Queue ADT

import TQueue as Q
import TStack as S

def evaluate(expr):
    """
    Evaluate a postfix expression.
    Pre-conditions:
        expr: a string containing a postfix expression
    Post-Conditions:
        none
    Return:
        the value of the expression
    """

    # create the initial empty data structures
    expression = Q.Queue()
    evaluation = S.Stack()

    expr_list = expr.split()
    # put all the items in the Queue
    for c in expr_list:
        expression.enqueue(c)

    # evaluate the expression
    while not expression.is_empty():
        c = expression.dequeue()

        if c == '*':
            v1 = evaluation.pop()
            v2 = evaluation.pop()
            evaluation.push(v1*v2)
        elif c == '/':
```

```
        v1 = evaluation.pop()
        v2 = evaluation.pop()
        evaluation.push(v2/v1)
    elif c == '+':
        v1 = evaluation.pop()
        v2 = evaluation.pop()
        evaluation.push(v1+v2)
    elif c == '-':
        v1 = evaluation.pop()
        v2 = evaluation.pop()
        evaluation.push(v2-v1)
    else:
        evaluation.push(float(c))

    return evaluation.pop()

example = "3 4 + 5 *"
print('( '+example+')', '=', evaluate(example))
```





## C — List-based Stacks and Queues

### C.1 The full list-based Queue implementation

The list-based queue is discussed in Chapter 12. Here we present the implementation without remark.

```
# Synopsis:
#   Implements the Queue ADT
#
# A queue (also called a FIFO queue) is a compound data
# structure in which the data values are ordered according to
# the FIFO (first-in first-out) protocol.

class Queue(object):

    def __init__(self):
        """
        Purpose
            creates an empty queue
        Return
            an empty queue
        """
        self.data = list()

    def is_empty(self):
        """
        Purpose
            checks if the given queue has no data in it
        Return:
            True if the queue has no data, or false otherwise
        """
```

```
        return len(self.data) == 0

def size(self):
    """
    Purpose
        returns the number of data values in the given queue
    Return:
        The number of data values in the queue
    """
    return len(self.data)

def enqueue(self, value):
    """
    Purpose
        adds the given data value to the given queue
    Pre-conditions:
        queue: a queue created by create()
        value: data to be added
    Post-condition:
        the value is added to the queue
    Return:
        (none)
    """
    self.data.append(value)

def dequeue(self):
    """
    Purpose
        removes and returns a data value from the given queue
    Pre-conditions:
        queue: a queue created by create()
    Post-condition:
        the first value is removed from the queue
    Return:
        the first value in the queue
    """
    return self.data.pop(0)

def peek(self):
    """
    Purpose
        returns the value from the front of given queue
        without removing it
    """
```

```

Pre-conditions:
    queue: a queue created by create()
Post-condition:
    None
Return:
    the value at the front of the queue
"""
return self.data[0]

```

## C.2 The full list-based Stack implementation

The list-based stack is discussed in Chapter 12. Here we present the implementation without remark.

```

# Synopsis:
#   Defines the Stack ADT
#
# A stack (also called a pushdown or LIFO stack) is a compound
# data structure in which the data values are ordered according
# to the LIFO (last-in first-out) protocol.
#
# Implementation:
# This implementation uses Python lists directly.

class Stack(object):
    def __init__(self):
        """
        Purpose
            initializes an empty Stack object
        """
        self.data = list()

    def is_empty(self):
        """
        Purpose
            checks if the stack has no data in it
        Return:
            True if the stack has no data, or false otherwise
        """
        return len(self.data) == 0

    def size(self):
        """
        Purpose
            returns the number of data values in the stack

```

```
Return:
    The number of data values in the queue
"""
return len(self.data)

def push(self, value):
    """
    Purpose
        adds the given data value to the stack
    Pre-conditions:
        value: data to be added
    Post-condition:
        the value is added to the stack
    Return:
        (none)
    """
    self.data.append(value)

def pop(self):
    """
    Purpose
        removes and returns a data value from the stack
    Post-condition:
        the top value is removed from the stack
    Return:
        returns the value at the top of the stack
    """
    return self.data.pop()

def peek(self):
    """
    Purpose
        returns the value from the front of stack
        without removing it
    Post-condition:
        None
    Return:
        the value at the front of the stack
    """
    return self.data[-1]
```

## D — Node-based Stacks and Queues

### D.1 The full list-based Queue implementation

The node-based queue is discussed in Chapter 16. Here we present the implementation without remark.

```
# CMPT 145 Course material
# Copyright (c) 2017-2020 Michael C Horsch
# All rights reserved.
#
# This document contains resources for homework assigned to students of
# CMPT 145 and shall not be distributed without permission. Posting this
# file to a public or private website, or providing this file to a person
# not registered in CMPT 145, constitutes Academic Misconduct, according
# to the University of Saskatchewan Policy on Academic Misconduct.
#
# Synopsis:
#   Implements the Queue ADT
#
# A queue (also called a FIFO queue) is a compound data
# structure in which the data values are ordered according to
# the FIFO (first-in first-out) protocol.
#
# Implementation notes:
# This implementation uses the linked node structure.

import node as N

class Queue(object):
```

```
def __init__(self):
    """
    Purpose
        creates an empty queue
    """
    self.__size = 0          # how many elements in the queue
    self.__front = None     # the node chain starts here
    self.__back = None      # the node chain ends here

def size(self):
    """
    Purpose
        returns the number of data values in the queue
    Pre-conditions:
        queue: a queue created by create()
    Return:
        The number of data values in the queue
    """
    return self.__size

def is_empty(self):
    """
    Purpose
        checks if the given queue has no data in it
    Return:
        True if the queue has no data, or false otherwise
    """
    return self.__size == 0

def enqueue(self, value):
    """
    Purpose
        adds the given data value to the queue
    Pre-conditions:
        value: data to be added
    Post-condition:
        the value is added to the queue
    Return:
        (none)
    """
    new_node = N.node(value, None)
```

```
if self.is_empty():
    self.__front = new_node
    self.__back = new_node
else:
    prev_last_node = self.__back
    prev_last_node.set_next(new_node)
    self.__back = new_node

self.__size += 1

def dequeue(self):
    """
    Purpose
        removes and returns a data value from the queue
        Note: the queue cannot be empty!
    Post-condition:
        the first value is removed from the queue
    Return:
        the first value in the queue, or None
    """
    assert not self.is_empty(), 'dequeued an empty queue'

    prev_first_node = self.__front
    result = prev_first_node.get_data()
    self.__front = prev_first_node.get_next()
    self.__size -= 1
    if self.__size == 0:
        self.__back = None
    return result

def peek(self):
    """
    Purpose
        returns the value from the front of queue
        without removing it
        Note: the queue cannot be empty!
    Post-condition:
        None
    Return:
        the value at the front of the queue
    """
    assert not self.is_empty(), 'peeked into an empty queue'
```



```

first_node = self.__front
result = first_node.get_data()
return result

```

## D.2 The full list-based Stack implementation

The node-based stack is discussed in Chapter 16. Here we present the implementation without remark.

```

# CMPT 145 Course material
# Copyright (c) 2017-2020 Michael C Horsch
# All rights reserved.
#
# This document contains resources for homework assigned to students of
# CMPT 145 and shall not be distributed without permission. Posting this
# file to a public or private website, or providing this file to a person
# not registered in CMPT 145, constitutes Academic Misconduct, according
# to the University of Saskatchewan Policy on Academic Misconduct.
#
# Synopsis:
#   Implements the Stack ADT
#
# A stack (also called a pushdown or LIFO stack) is a compound
# data structure in which the data values are ordered according
# to the LIFO (last-in first-out) protocol.
#
# Implementation:
# This implementation uses the linked node structure.

import node as N

class Stack(object):

    def __init__(self):
        """
        Purpose
            creates an empty stack
        """
        self.__size = 0          # how many elements in the stack
        self.__top = None        # the node chain starts here

    def size(self):
        """

```

```
    Purpose
        returns the number of data values in the stack
    Return:
        The number of data values in the stack
    """
    return self.__size

def is_empty(self):
    """
    Purpose
        checks if the stack has no data in it
    Return:
        True if the stack has no data, or false otherwise
    """
    return self.__size == 0

def push(self, value):
    """
    Purpose
        adds the given data value to the stack
    Pre-conditions:
        value: data to be added
    Post-condition:
        the value is added to the stack
    Return:
        (none)
    """
    new_node = N.node(value, self.__top)
    self.__top = new_node
    self.__size += 1

def pop(self):
    """
    Purpose
        Removes and returns a data value from the stack.
        Note: the stack cannot be empty!
    Post-condition:
        the first value is removed from the stack
    Return:
        the first value in the stack, or None
    """
    assert not self.is_empty(), 'popped an empty stack'
```

```
prev_first_node = self.__top
result = prev_first_node.get_data()
self.__top = prev_first_node.get_next()
self.__size -= 1
return result

def peek(self):
    """
    Purpose
        returns the value from the front of given stack
        without removing it
        Note: the stack cannot be empty!
    Post-condition:
        None
    Return:
        the value at the front of the stack
    """
    assert not self.is_empty(), 'peeked into an empty stack'

    first_node = self.__top
    result = first_node.get_data()
    return result
```

## E — The treenode ADT

### E.1 The full Treenode ADT implementation


The node-based queue is discussed in Chapter 20. Here we present the implementation without remark.

```
# CMPT 145 Course material
# Copyright (c) 2017-2020 Michael C Horsch
# All rights reserved.
#
# This document contains resources for homework assigned to students of
# CMPT 145 and shall not be distributed without permission. Posting this
# file to a public or private website, or providing this file to a person
# not registered in CMPT 145, constitutes Academic Misconduct, according
# to the University of Saskatchewan Policy on Academic Misconduct.
#
# Synopsis:
#   Defines the tree node ADT
#
# A treenode is a simple container with three pieces of information
#   data: the contained information
#   left:  a reference to another treenode or None
#   right: a reference to another treenode or None

class treenode(object):

    def __init__(self, data, left=None, right=None):
        """
        Create a new treenode for the given data.
```

```
Pre-conditions:
    data: Any data value to be stored in the treenode
    left: Another treenode (or None, by default)
    right: Another treenode (or None, by default)
"""
self.data = data
self.left = left
self.right = right
```



Maximum Slice  
Subset Sum  
Maximum Tree Path  
Make Change  
Leap Line

## F — Algorithms

A variety of algorithms were discussed in Chapter 27. Here we present the full implementations without remark.

### F.1 Maximum Slice

```
# CMPT 145 - Algorithms
# The Maximum Slice Problem
# Given a list A containing (positive and negative) numbers,
# Find the slice A[a:b] that has the maximum sum

import random as rand
import time as time

# pedagogical example:
def allslices(alist):
    """
    Purpose:
        Display all slices to the console.
    Preconditions:
        alist: a list of numbers
    Post-conditions:
        Outputs all slices to the console.
    Return:
        None
    """
    for a in range(len(alist)):
        for b in range(a, len(alist)):
```

```
        print(alist[a:b+1])

# version 0: naively sums each slice
def maxslice_brute_force_v0(alist):
    """
    Purpose:
        Find the maximum sum of all slices of alist.
    Preconditions:
        alist: a list of numbers
    Post-conditions:
        None
    Return:
        a number, the maximum slice sum
    """
    maxsum = alist[0]
    for i in range(len(alist)):
        for j in range(i+1, len(alist)):
            slice = sum(alist[i:j + 1])
            if slice > maxsum:
                maxsum = slice
    return maxsum

# version 1: sums the slices more cleverly
def maxslice_brute_force_v1(alist):
    """
    Purpose:
        Find the maximum sum of all slices of alist.
    Preconditions:
        alist: a list of numbers
    Post-conditions:
        None
    Return:
        a number, the maximum slice sum
    """
    # using brute force: look at all possible slices
    # but store all the partial sums in a dictionary
    # where s[j] stores the value sum(alist[i,j+1])

    maxsum = alist[0]
    for i in range(len(alist)):
        s = {}
        s[i] = alist[i]
        if s[i] > maxsum:
            maxsum = s[i]
```

```
        for j in range(i+1, len(alist)):
            s[j] = s[j-1] + alist[j]
            if s[j] > maxsum:
                maxsum = s[j]
    return maxsum

# version 2 uses divide and conquer
# Divide the list into 2 halves
# The maximum slice can occur in one of 3 ways:
# 1. Starts and finishes on the left half
# 2. Starts and finishes on the right half
# 3. Starts somewhere in the left half, and
#    continues somewhere to the right half
def maxslice_DC(alist):
    """
    Purpose:
        Find the maximum sum of all slices of alist.
    Preconditions:
        alist: a list of numbers
    Post-conditions:
        None
    Return:
        a number, the maximum slice sum
    """

    # internal function
    def max_tail(left, right):
        """
        Calculate the maximum slice that ends at right
        (from any point starting at left or later)
        """
        s = alist[right]
        maxsum = s
        # calculate the sums from right to left (backwards)
        for i in range(right - 1, left - 1, -1):
            s = s + alist[i]
            if s > maxsum:
                maxsum = s
        return maxsum

    # internal function
    def max_head(left, right):
        """
        Calculate the maximum slice that starts at left
        (to any point up to and including right)
```



```

    """
    s = alist[left]
    maxsum = s
    for i in range(left + 1, right + 1):
        s = s + alist[i]
        if s > maxsum:
            maxsum = s
    return maxsum

# internal function
def maxslice_rec(left, right):
    """
    Recursively find maximum slice between left and right.
    """
    # using divide and conquer
    if left == right:
        return alist[left]
    else:
        # divide, and solve
        mid = (right + left) // 2
        max_left = maxslice_rec(left, mid)
        max_right = maxslice_rec(mid + 1, right)
        max_cross = (max_tail(left, mid)
                     + max_head(mid + 1, right))

        # conquer
        return max(max_left, max_right, max_cross)

# body of maxslice_DC
return maxslice_rec(0, len(alist) - 1)

#148
def allslices(alist):
    for a in range(len(alist)):
        for b in range(a, len(alist)):
            print(alist[a:b+1])

allslices([1,2,3,4])

# put the versions through their paces
if __name__ == '__main__':
    examples = [[1, 2, 3, 4, 5],
                [5, 4, -3, 2, 1],
                [1, 2, -3, 4, 5],

```

```

        [1, -2, 3, 4, -5],
        [1, -2, 3, 4, -5, 1, 2, 3, 4, 5, -6, -7, -8],
        [rand.randint(-100, 100) for i in range(1000)],
        [rand.randint(-100, 100) for i in range(2000)]
    ]

    for ex in examples:
        print('Example: list of length:', len(ex))

        print('Brute Force version 0:')
        start = time.process_time()
        result = maxslice_brute_force_v0(ex)
        end = time.process_time()
        print('Result:', result, 'Time:', (end - start))

        print('Brute Force version 1:')
        start = time.process_time()
        result = maxslice_brute_force_v1(ex)
        end = time.process_time()
        print('Result:', result, 'Time:', (end - start))

        print('Divide and conquer:')
        start = time.process_time()
        result = maxslice_DC(ex)
        end = time.process_time()
        print('Result:', result, 'Time:', (end - start))

    print()

```

## F.2 Subset Sum

```

# CMPT 145 - Algorithms
# The Subset Sum problem
#
# Given a list of integers L, and a target T
# Find a subset of numbers from L whose sum is T

# pedagogical example, working through subsets
def all_subsets(alist):
    """
    Purpose:
        Given a list, display all subsets.
    Preconditions:
        alist: a list
    """

```

```

Post-conditions:
    displays all subsets to the console,
    one subset per line. Could be a lot!
Return:
    None
"""
def allsub(al, sub):
    if len(al) == 0:
        print(sub)
    else:
        first = al[0]
        rest = al[1:]
        allsub(rest, [first]+sub) # with first
        allsub(rest, sub)        # without first

allsub(alist, [])

# Brute force solution
def subsetsum_v1(alist, target):
    """
    Purpose:
        Given a list of positive integers, and a target,
        Return a sublist of integers whose sum is target.
    Preconditions:
        alist: a list of positive integers
        target: a positive integer
    Return:
        A tuple True, sublist if sum(sublist) == target
        Or False, None otherwise
    """

    def trysum(al, subset):
        """
        Recursively try every subset of values from al.
        """
        if len(al) == 0 and sum(subset) == target:
            return True, subset
        elif len(al) == 0:
            return False, None
        else:
            first = al[0]
            rest = al[1:]
            flag, answer = trysum(rest, subset + [first])
            if flag:
                return flag, answer

```

```
        else:
            # try without the first value
            return trysum(rest, subset)

    return trysum(alist, [])

# Backtracking solution
def subsetsum_v2(alist, target):
    """
    Purpose:
        Given a list of positive integers, and a target,
        Return a sublist of integers whose sum is target.
    Preconditions:
        alist: a list of positive integers
        target: a positive integer
    Return:
        A tuple True, sublist if sum(sublist) == target
        Or False, None otherwise
    """

    def trysum(al, subset):
        """
        Recursively try every subset of values from al.
        """
        if sum(subset) == target:
            return True, subset
        elif len(al) == 0 or sum(subset) > target:
            return False, None
        else:
            first = al[0]
            rest = al[1:]
            flag, answer = trysum(rest, subset + [first])
            if flag:
                return flag, answer
            else:
                # try without the first value
                return trysum(rest, subset)

    return trysum(alist, [])

# Backtracking solution
def subsetsum_v22(alist, target):
    """
    Purpose:
```

```

    Given a list of positive integers, and a target,
    Return a sublist of integers whose sum is target.
Preconditions:
    alist: a list of positive integers
    target: a positive integer
Return:
    A tuple True, sublist if sum(sublist) == target
    Or False, None otherwise
"""

def trysum(al, subset, thesum):
    """
    Recursively try every subset of values from alist.
    """
    if thesum == target:
        return True, subset
    elif len(al) == 0 or thesum > target:
        return False, None
    else:
        first = al[0]
        rest = [v for v in al[1:] if v <= (target - thesum)]
        flag, answer = trysum(rest, subset + [first], thesum+first)
        if flag:
            return flag, answer
        else:
            # try without the first value
            return trysum(rest, subset, thesum)

    return trysum(sorted(alist, reverse=True), [], 0)

# Divide and Conquer solution, not very smart
def subsetsum_v3(alist, target):
    """
    Purpose:
        Given a list, display all subsets.
    Preconditions:
        alist: a list
    Post-conditions:
        displays all subsets to the console,
        one subset per line. Could be a lot!
    Return:
        None
    """

```

```
def trysum(alist, subset):
    """
    Recursively construct a subset from alist, checking
    if the sum is equal to the target.
    """
    if sum(subset) == target:
        return True, subset
    elif len(alist) == 0:
        return False, None
    else:
        tryit = alist[0]
        # try to find a subset sum using the first value
        flag, answer = trysum(alist[1:], subset + [tryit])
        if flag:
            return flag, answer
        else:
            # try without the first value
            return trysum(alist[1:], subset)

# this does divide and conquer, in the form of
# dividing the list in half
def binary(alist):
    if len(alist) < 5:
        # for small lists, just use brute force!
        return trysum(alist, [])
    else:
        mid = len(alist)//2

        # try to find a solution in a list half the size
        flag, ans = binary(alist[:mid])
        if flag: return True, ans

        # if not, try the other half
        flag, ans = binary(alist[mid + 1:])
        if flag: return True, ans

        # if not, try the whole list
        return trysum(alist, [])

return binary(alist)
```

```
import random as rand
```

```
import time as time
if __name__ == '__main__':

    example1 = [1, 2, 3, 4, 5]
    example2 = [1, 3, 5, 7]
    target = 8

    for ex in [example1, example2]:
        print('Example: list of length:', len(ex))

        print('Brute Force (Version 1):')
        start = time.process_time()
        subset = subsetsum_v1(ex, target)
        end = time.process_time()
        print('Target:', target,
              'Time:', end - start,
              'Result:', subset)

        print('Backtracking (Version 2):')
        start = time.process_time()
        subset = subsetsum_v2(ex, target)
        end = time.process_time()
        print('Target:', target,
              'Time:', end - start,
              'Result:', subset)

        print('Divide and Conquer (Version 3):')
        start = time.process_time()
        subset = subsetsum_v3(ex, target)
        end = time.process_time()
        print('Target:', target,
              'Time:', end - start,
              'Result:', subset)

    print()

    lo = 0
    hi = 5000
    example_size = 22
    example3 = [rand.randint(lo, hi)
                 for i in range(example_size)]

    number_of_trials = 5
    for i in range(number_of_trials):
        target = rand.randint(hi, hi*100)
```

```
print('Brute Force (Version 1) on list of size',
      example_size, ':')
start = time.process_time()
subset = subsetsum_v1(example3, target)
#subset = None
end = time.process_time()
print('Target:', target,
      'Time:', end - start,
      'Result:', subset)

print('Backtracking (Version 2) on list of size',
      example_size, ':')
start = time.process_time()
subset = subsetsum_v2(example3, target)
end = time.process_time()
print('Target:', target,
      'Time:', end - start,
      'Result:', subset)

print('Backtracking (Version 2.2) on list of size',
      example_size, ':')
start = time.process_time()
subset = subsetsum_v22(example3, target)
end = time.process_time()
print('Target:', target,
      'Time:', end - start,
      'Result:', subset)

print('Divide and Conquer (Version 3) on list of size',
      example_size, ':')
start = time.process_time()
subset = subsetsum_v3(example3, target)
end = time.process_time()
print('Target:', target,
      'Time:', end - start,
      'Result:', subset)

print()
```

### F.3 Maximum Tree Path

```
# CMPT 145 - Algorithms
# The Maximum Path Problem
# Given a Binary Tree (not necessarily a Binary Search Tree)
# containing (positive and negative) numbers,
```



```

# Find the path from root to leaf that has the maximum sum

import random as rand
import time

# a quick and dirty implementation of TreeNodes
class Treenode(object):
    def __init__(self, data, left=None, right=None):
        self.data = data
        self.left = left
        self.right = right

    def is_leaf(self):
        return self.left is None and self.right is None

def random_tree(h, lor=-200, hir=200):
    """
    Construct a complete binary tree of height h with
    randomly chosen data values
    :param h: the height of the tree
    :param lor: an integer
    :param hir: an integer, lor <= hir
    :return: a complete binary Treenode of height h
    """
    if h == 1:
        return Treenode(rand.randint(lor, hir))
    else:
        return Treenode(rand.randint(lor, hir),
                        random_tree(h - 1, lor, hir),
                        random_tree(h - 1, lor, hir))

# Brute force solution: check all possible paths
def max_path(atree):
    """
    Find the maximum sum of all paths from root to leaf
    in atree.
    :param atree: A Treenode
    :return: the maximum sum
    """
    if atree is None:
        return 0
    else:

```

```

        left = max_path(atree.left)
        right = max_path(atree.right)
        return max(left, right) + atree.data

# Greedy solution
# Note: this algorithm may not return the true maximum!
def greedy_path(atree):
    """
    Find the maximum sum of all paths from root to leaf
    in atree.
    :param atree: A TreeNode
    :return: the maximum sum
    """
    if atree.is_leaf():
        return atree.data
    elif atree.left.data > atree.right.data:
        # the left branch seems promising
        return atree.data + greedy_path(atree.left)
    else:
        # the right branch seems promising
        return atree.data + greedy_path(atree.right)

if __name__ == '__main__':
    # build the tree;  $O(2^N)$  time!
    height = 21
    maxval = 100
    print('Building tree...')
    start = time.process_time()
    some_tree = random_tree(height)
    end = time.process_time()
    print('...took:', end - start, 'seconds')

    print('Brute Force Implementation:')
    start = time.process_time()
    result = max_path(some_tree)
    end = time.process_time()
    print("Result:", result, 'Time:', end - start)

    print('Greedy Implementation:')
    start = time.process_time()
    result = greedy_path(some_tree)
    end = time.process_time()
    print("Result:", result, 'Time:', end - start)

```

## F.4 Make Change

```
# CMPT 145 - Algorithms
# The Make Change Problem
#
# Given a number D in the range 0-99
# Find the smallest collection of coins whose value is D

# a greedy solution, adding coins one at a time
def change_v1(cents):
    """
    Purpose:
        Make change for the given cents value.
        Assumes coin values 25c, 10c, 5c, 1c
    Pre-conditions:
        :param cents: an integer
    Return:
        a list of counts for the coins used.
    """
    coins = [25, 10, 5, 1]
    coin_index = 0
    counts = [0, 0, 0, 0]
    remaining = cents
    while remaining > 0:
        if coins[coin_index] <= remaining:
            counts[coin_index] += 1
            remaining -= coins[coin_index]
        else:
            coin_index += 1
    return counts

# a greedy solution, calculating each quantity
# of coins exactly
def change_v2(cents):
    """
    Purpose:
        Make change for the given cents value.
        Assumes coin values 25c, 10c, 5c, 1c
    Pre-conditions:
        :param cents: an integer
    Return:
        a list of counts for the coins used.
    """
```

```
coins = [25, 10, 5, 1]
coin_index = 0
counts = [0] * len(coins)
remaining = cents
while remaining > 0:
    counts[coin_index] = remaining // coins[coin_index]
    remaining = remaining % coins[coin_index]
    coin_index += 1
return counts

# A brute force solution
# Find the smallest of all combinations of coins
def change_v3(cents):
    """
    Make change for the given cents value.
    Assumes coin values 25c, 10c, 5c, 1c
    :param cents: an integer
    :return: a list of counts for the coins used.
    """

    coins = [25, 10, 5, 1]

    def combinations(counts):
        """
        Cycle through every possible combination, looking
        for a combo that adds up to the right value, and
        has the smallest number of coins
        :param counts: a 4-tuple
        :return: a pair (True, list) if list is the best
                combination
        """
        value = sum([counts[i] * v for i, v in enumerate(coins)])
        if value == cents:
            return True, counts
        elif value > cents:
            # don't add more coins to this combination
            # because it's already too big
            return False, None
        else:
            (c0, c1, c2, c3) = counts
            # add 1 to each number of coins separately
            trying = [(c0 + 1, c1, c2, c3), (c0, c1, c2 + 1, c3),
                     (c0, c1 + 1, c2, c3), (c0, c1, c2, c3 + 1)]
            # try to find the best combination, by using
            # the lowest number of coins
```

```

        best_size = 100
        best_counts = None
        for combo in trying:
            flag, res = combinations(combo)
            if flag and sum(res) < best_size:
                best_size = sum(res)
                best_counts = res
        if best_size == 100:
            # nothing worked!
            return False, None
        else:
            return True, best_counts

    flag, result = combinations((0, 0, 0, 0))
    if flag:
        return result
    else:
        return None

if __name__ == '__main__':
    examples = [0, 25, 37, 49, 51, 87, 99, 104]

    for e in examples:
        #print('Version 1:', e, change_v1(e))
        print('Greedy:', e, change_v2(e))
        #print('Version 3:', e, change_v3(e))

    print()

```

## F.5 Leap Line

```

# CMPT 145 - Algorithms
# The Leap Line Problem
# Given a list A containing (positive and negative) numbers,
# Assuming a collection policy
# Find the maximum score that can be collected
#
# Collection policy:
#   Mario must leap over 1 item or step into it
#   Leaping over item prevents collection
#   Stepping into an item ensures collection

import time as time

```



```

    return jump_or_step(0)

# version 1: Using Memoization
def maximumScoreFrom_v1(track):
    """
    Purpose:
        Calculate the maximum score that can be obtained from
        stepping and jumping along the given track
    Pre-conditions:
        :param track: a list of integers
    Return:
        the maximum score
    """
    def jump_or_step(loc):
        """
        Calculate the maximum score that can be obtained from
        stepping and jumping along the given track, starting
        at the given location
        """
        # check if the best score is already known
        if loc in memo:
            return memo[loc]

        step_loc = loc + 1
        jump_loc = loc + 2

        if step_loc == len(track):
            return track[loc]
        elif jump_loc >= len(track):
            result = track[loc] + jump_or_step(step_loc)
            memo[loc] = result
            return result
        else:
            result = track[loc] + max(jump_or_step(step_loc),
                                      jump_or_step(jump_loc))
            memo[loc] = result
            return result

    # body of maximumScoreFrom_v1()
    # using memoization
    # memo[loc] stores the best score starting from loc.
    memo = {}
    return jump_or_step(0)

```

```

# version 2: Greedy
# looks one location ahead, and leaps over -1
def maximumScoreFrom_v2(track):
    """
    Purpose:
        Calculate the maximum score that can be obtained from
        stepping and jumping along the given track
    Pre-conditions:
        :param track: a list of integers
    Return:
        the maximum score
    """

    def jump_or_step(loc):
        """
        Calculate the maximum score that can be obtained from
        stepping and leaping along the given track, starting
        at the given location
        """
        step_loc = loc + 1
        jump_loc = loc + 2

        if step_loc == len(track):
            return track[loc]
        elif jump_loc >= len(track):
            return track[loc] + leap_or_step(step_loc)
        elif track[step_loc] < 0:
            # greedy choice: avoid negatives
            return track[loc] + leap_or_step(jump_loc)
        else:
            # greedy choice: collect 1
            return track[loc] + leap_or_step(step_loc)

    return jump_or_step(0)

examples = [
    ['easy',
     [0, -1, -1, 1, -1, -1, 0]],
    ['medium',
     [0, -1, 1, -1, -1, 1, -1, -1, 1, -1, 1, 0]],
    ['challenging',
     [0, -1, 1, -1, -1, 1, -1, -1, 1, -1, 1, -1, -1, 0]],
    ['hardmode',
     [0, -1, 1, -1, -1, 1, -1, -1, 1, -1, 1, -1, -1, 1,
      1, 1, -1, -1, -1, 1, -1, 1, -1, -1, 1, -1, -1, 1,

```



```

        -1, 1, -1, -1, 1, 1, 1, -1, 0]],
    ['very_hard',
     [0, -1, 1, -1, -1, 1, -1, -1, 1, -1, 1, -1, -1, 1,
      1, 1, -1, 1, -1, -1, 1, 1, 1, 1, 1, -1, 1, -1, -1,
      1, -1, -1, 1, -1, 1, 1, -1, -1, 1, -1, -1, 1, -1,
      1, -1, -1, 0]]
]

for e in examples:
    print('Example: list of length:', len(e[1]))

    print('Brute Force version (v0):', end=" ")
    start = time.process_time()
    result = maximumScoreFrom_v0(e[1])
    end = time.process_time()
    print('Result', '(' + e[0] + '):', result,
          'Time:', (end - start))

    print('Memoized version (v1):', end=" ")
    start = time.process_time()
    result = maximumScoreFrom_v1(e[1])
    end = time.process_time()
    print('Result', '(' + e[0] + '):', result,
          'Time:', (end - start))

    print('Greedy version (v2):', end=" ")
    start = time.process_time()
    result = maximumScoreFrom_v2(e[1])
    end = time.process_time()
    print('Result', '(' + e[0] + '):', result,
          'Time:', (end - start))

print()

```