

10.2. Hash Tables 415

Hash Codes in Python

The standard mechanism for computing hash codes in Python is a built-in function with signature `hash(x)` that returns an integer value that serves as the hash code for object `x`. However, only immutable data types are deemed hashable in Python. This restriction is meant to ensure that a particular object's hash code remains constant during that object's lifespan. This is an important property for an object's use as a key in a hash table. A problem could occur if a key were inserted into the hash table, yet a later search were performed. Among Python's built-in data types, the immutable `int`, `float`, `str`, `tuple`, and `frozenset` classes produce robust hash codes, via the `hash` function, using techniques similar to those discussed earlier in this section. Hash codes for character strings are well crafted based on the algorithm of Dijkstra [1962]. Hash codes for tuples are computed with a similar technique based upon a combination of the hash codes of the individual elements of the tuple. When hashing a `frozenset`, the order of the elements should be irrelevant, and so a natural option is to compute the exclusive-or of the individual hash codes of the elements. For mutable types, such as a `list`, a `TypeError` is raised. Instances of user-defined classes are treated as unhashable by default, with a `TypeError` raised by the `hash` function. However, a function that computes hash codes can be implemented in the form of a special method named

`hash`

within

a class. The returned hash code should reflect the immutable attributes of an instance. It is common to return a hash code composed of three numeric red, green, and blue components might implement the method as:

```
def
```

```
hash
```

```
(self):
```

```
    return hash( ( self.
```

```
red, self.
```

```
green, self.
```

```
blue) ) # hash combined tuple
```

An important rule to obey is that if a class defines equivalence through

```
eq
```

```
,
```

then any implementation of

```
hash
```

must be consistent, in that if `x == y`, then

`hash(x) == hash(y)`. This is important because if two instances are considered

to be equivalent and one is used as a key in a hash table, a search for the second instance should result in the dictionary. If `hash` is not consistent, it ensures that `hash(5)` and `hash(5.0)` are the same.

10.2. Hash Tables 413

collisions for common groups of strings. In particular, "temp01" and "temp10" collide using this function, as do "stop", "tops", "pots", and "spot". A better hash code should somehow take into consideration the positions of the x_i 's. An alternative hash code, which does exactly this, is to choose a nonzero constant, a , and use as a hash code the value $x_0a^{n-1} + x_1a^{n-2} + \dots + x_{n-2}a + x_{n-1}$.

Mathematically speaking, this is simply a polynomial in a that takes the components $(x_0, x_1, \dots, x_{n-1})$ of an object x as its coefficients. This hash code is therefore called a polynomial hash code. By Horner's rule (see Exercise C-3.50), this polynomial can be computed as $x_{n-1} + a(x_{n-2} + a(x_{n-3} + \dots + a(x_2 + a(x_1 + ax_0)) \dots))$.

Intuitively, a polynomial hash code uses multiplication by different powers as a way to spread out the influence of each component across the resulting hash code.

Of course, on a typical computer, evaluating a polynomial will be done using the finite bit representation for a hash code; hence, the value will periodically overflow the bits used for an integer. One should be mindful that such overflows are occurring and choose the constant a so that it has some nonzero, low-order bits, which will serve to preserve some of the information content even as we overflow. We have done some experimental studies that suggest that 33, 37, 39, and 41 are particularly good choices for a when working with character strings that are English words. In fact, in a list of over 50,000 English words formed as the union of the word lists provided in two of the previous chapters, 37, 39, or 41 produced less than 7 collisions in each case!

Cyclic-Shift Hash Codes

A variant of the polynomial hash code replaces multiplication by a with a cyclic shift of a partial sum by a certain number of bits. For example, a 5-bit cyclic shift of the 32-bit value 00111101100101101010100010101000 is achieved by taking the leftmost five bits and placing those on the rightmost side of the representation, resulting in 10110010110101010001010100000111. While this operation has little natural meaning in terms of arithmetic, it accomplishes the goal of varying the bits of the calculation. In Python, a 32-bit integer is

10.6. Exercises 455

C-10.32 Perform experiments on our ChainHashMap and ProbeHashMap classes to measure its efficiency using random key sets and varying limits on the load factor (see Exercise R-10.15).

C-10.33 Our implementation of separate chaining in ChainHashMap conserves memory by representing empty buckets in the table as None, rather than as empty instances of a secondary structure. Because many of these buckets will hold a single item, a better option would be to have the table directly reference the Item instance, and to reserve use of secondary containers for buckets that have two or more items. Modify our implementation to provide this additional optimization.

C-10.34 Computing a hash code can be expensive, especially for lengthy keys. In our hash table implementation, inserting an item, and recompute each item's hash code each time we resize our table. Python's dict class makes an interesting trade-off. The hash code is computed once, when an item is inserted, and the hash code is stored as an extra field of the item component.

C-10.35 Describe how to perform a removal from a hash table that uses linear probing to resolve collisions where we do not use a special marker to represent deleted elements. That is, we must rearrange the contents so that it appears that the removed entry was never there.

C-10.36 The quadratic probing strategy has a clustering problem related to the way it looks for open slots. Namely, it checks buckets $A[(h(k) + i^2) \bmod N]$, for $i = 1, 2, \dots, N-1$.

a. Show that $i^2 \bmod N$ will assume at most $(N+1)/2$ distinct values, for N prime, as i ranges from 1 to $N-1$. As a part of this justification, note that $i^2 \bmod N = (N-i)^2 \bmod N$ for all i .

b. A better strategy is to choose a prime N such that $N \bmod 4 = 3$ and then to check the buckets $A[(h(k) \pm i^2) \bmod N]$ as i ranges from 1 to $(N-1)/2$, alternating between plus and minus. Show that this alternate version is guaranteed to check every bucket in A .

C-10.37 Refactor our ProbeHashMap design so that the sequence of secondary probes for collision resolution can be more easily customized. Demonstrate your new framework by providing several examples.

C-10.38 Design a variation of binary search for performing the multimap operation $\text{nd_all}(k)$ implemented with a sorted search table that includes duplicates, and show that it runs in time $O(s + \log n)$, where n is the number of elements in the dictionary and s is the number of items with given key k .

742 Index

directed, 620, 621, 657

acyclic, 655-657

strongly connected, 623

mixed, 621

reachability, 651-654

shortest paths, 654simple, 622traversal, 638-650

undirected, 620, 621

weighted, 659-696

greedy method, 603, 660, 661Guibas, Leonidas, 535

Gutttag, John, 108, 254, 298

Harmonic number, 131, 180, 728

hash code, 411-415

cyclic-shift, 413-414

polynomial, 413

hash function, 415

hash table, 410-426

clustering, 419

collision, 411

collision resolution, 417-419

double hashing, 419

linear probing, 418

quadratic probing, 419

HashMapBase class, 422-423

header sentinel, 270heap, 370-384

bottom-up construction, 380-384

heap-sort, 384, 388-389

HeapPriorityQueue class, 377-378, 382

heapq module, 384

height of a tree, 309-310, 474

height-balance property, 481, 483

Hell, Pavol, 696Hennessy, John, 719heuristic, 289hierarchy, 82

Hoare, C. A. R., 580

hook, 342, 468, 478Hopcroft, John, 254, 298, 535, 696Hopper, Grace, 36

Horner-s method, 146

HTML, 236-238, 251, 582Huang, Bing-Chao, 580Huffman coding, 601-602I/O complexity, 711idfunction, 29

identifier, 41DLE, 2

immutable type, 7, 11, 415

implied method, 76

import statement, 48

inoperator, 14-15, 75

in-degree, 621in-place algorithm, 389, 396, 559, 702

incidence collection, 630incident, 621

incoming edges, 621

independent, 729, 730

index, 186

negative, 14

zero-indexing, 9, 14

IndexError, 20, 33, 34, 83, 232, 303

induction, 138-139, 162

Preface vii

Contents and Organization

The chapters for this book are organized to provide a pedagogical path that starts with the basics of Python programming and object-oriented design. We then add foundational techniques like algorithm analysis and recursion. In the main portion of the book, we present fundam

1. Python Primer

2. Object-Oriented Programming

3. Algorithm Analysis

4. Recursion

5. Array-Based Sequences

6. Stacks, Queues, and Deques

7. Linked Lists

8. Trees

9. Priority Queues

10. Maps, Hash Tables, and Skip Lists

11. Search Trees

12. Sorting and Selection

13. Text Processing

14. Graph Algorithms

15. Memory Management and B-Trees

A. Character Strings in Python

B. Useful Mathematical Facts

A more detailed table of contents follows this preface, beginning on page xi.

Prerequisites

We assume that the reader is at least vaguely familiar with a high-level programming language, such as C, C++

- Variables and expressions.

- Decision structures (such as if-statements and switch-statements).

- Iteration structures (for loops and while loops).

- Functions (whether stand-alone or object-oriented methods).

For readers who are familiar with these concepts, but not with how they are expressed in Python, we provide a p
give a comprehensive treatment of Python.

Index 745

open function, 29, 31

operand stack, 704

operators, 12-17

arithmetic, 13

bitwise, 14

chaining, 17comparisons, 13logical, 12

overloading, 74

precedence, 17

operator, 12

ordfunction, 29

order statistic, 571

OrderedDict class, 457

Orlin, James, 696

osmodule, 49, 159, 182, 357

out-degree, 621

outgoing edge, 621

over-ow, 506

override, 82, 100

p-norm, 53

packing a sequence, 44

palindrome, 181, 615

parameter, 24-28

actual, 24

default value, 26formal, 24

keyword, 27

positional, 27

parent class, 82

parent node, 301

parenthetic string representation, 339

partial order, 16

partition, 679, 681-684

pass statement, 38, 478

path, 302, 623

compression, 684

directed, 623length, 356, 660simple, 623

pattern matching, 208, 584-593

Boyer-Moore algorithm, 586-589

brute force, 584-585Knuth-Morris-Pratt algorithm,
590-593

Patterson, David, 719Perkovic, Ljubomir, 55

permutation, 150Peters, Tim, 568

Phillips, Dusty, 108polymorphism, 26, 77, 93polynomial function, 119, 146polynomial **hash** code, 413

portability, 58

position, 279-281, 305, 438

positional list, 277-285

abstract data type, 279-281

PositionalList class, 281-285, 287, 628

positional parameter, 27post-x notation, 252, 253, 359

postorder tree traversal, 329

10.5. Sets, Multisets, and Multimaps 449

```
def
  or
  (self, other) : # supports syntax S | T
  ...Return a new set that is the union of two existing sets...
  result = type( self )( ) # create new instance of concrete class
  for e in self :
    result.add(e)
  for e in other:
    result.add(e)
  return result
```

Code Fragment 10.15: An implementation of the MutableSet.

or
method,
which computes the union of two existing sets.

An implementation of the behavior that computes a new set as a union of two others is given in the form of the

or
special method, in Code Fragment 10.15.

An important subtlety in this implementation is the instantiation of the resulting set. Since the MutableSet class is a trait, it must belong to a concrete subclass. When computing the union of two such concrete instances, the result should presumably be an instance of the same class as the operands. The function `type(self)` returns a reference to the actual class of the instance identified as `self`, and the subsequent parentheses in expression `type(self)()` call the default constructor for that class.

In terms of efficiency, we analyze such set operations while letting n denote the size of S and m denote the size of T for an operation such as $S | T$. If the concrete sets are implemented with hashing, the expected running time of the implementation in Code Fragment 10.15 is $O(m+n)$, because it loops over both sets, performing constant-time operations in the form of a containment check and a possible insertion into the result.

Our implementation of the in-place version of a union is given in Code Fragment 10.16, in the form of the

or
special method that supports syntax $S |= T$.

Notice that in this case, we do not create a new set instance, instead we modify and return the existing set, after which the in-place version of the union has expected running time $O(m)$ where m is the size of the second set, because we only have to loop through that second set.

```
def
  or
  (self, other) : # supports syntax S |= T
  ...Modify this set to be the union of itself and another set...
  for e in other:
    self.add(e)
  return self # technical requirement of in-place operator
```

Code Fragment 10.16: An implementation of the MutableSet.

or
method,
which performs an in-place union of one set with another.

50 Chapter 1. Python Primer

next number in a sequence based upon one or more past numbers that it has generated. Indeed, a simple yet popular pseudo-random number generator chooses its next number based solely on the most recently chosen number and some additional parameters using the following

$$\text{next} = (a * \text{current} + b) \% n;$$

where a, b, a and n are appropriately chosen integers. Python uses a more advanced technique known as a Mersenne twister. It turns out that the sequences generated by these techniques can be proven to be statistically uniform, which is usually good enough for most applications requiring random numbers, such as games. For applications, such as computer security settings, where one needs unpredictable random sequences, this kind of randomness comes from outer space.

Since the next number in a pseudo-random generator is determined by the previous number(s), such a generator always needs a place to start, which is called its seed. The sequence of numbers generated for a given seed will always be the same.

One common trick to get a different sequence each time a program is run is to use a seed that will be different for each run. For example, we could use some timed input from a user or the current system time in milliseconds.

Python's random module provides support for pseudo-random number generation by defining a Random class; instances of that class serve as generators with independent state. This allows different aspects of a program to rely on their own pseudo-random number generators. The random module (essentially using a single generator instance for all top-level calls).

Syntax

Description

`seed(hashable)`

Initializes the pseudo-random number generator based upon the `hash` value of the parameter `random()`

Returns a pseudo-random floating-point value in the interval `[0.0,1.0)`.

`randint(a,b)`

Returns a pseudo-random integer in the closed interval `[a,b]`.

`randrange(start, stop, step)`

Returns a pseudo-random integer in the standard Python range indicated by the parameters.

`choice(seq)`

Returns an element of the given sequence chosen pseudo-randomly.

`shuffle(seq)`

Reorders the elements of the given sequence pseudo-randomly.

Table 1.8: Methods supported by instances of the Random class, and as top-level functions of the random module.

2.3. Class Definitions 75

Common Syntax

Special Method Form

$a+b$

$a.$

`add`

(b) ; alternatively $b.$

`radd`

(a)

$a \cdot b$

$a.$

`sub`

(b) ; alternatively $b.$

`rsub`

(a)

a

b

$a.$

`mul`

(b) ; alternatively $b.$

`rmul`

(a)

a/b

$a.$

`truediv`

(b) ; alternatively $b.$

`rtruediv`

(a)

$a//b$

$a.$

`floordiv`

(b) ; alternatively $b.$

`rfloordiv`

(a)

$a \% b$

$a.$

`mod`

(b) ; alternatively $b.$

`rmod`

(a)

a

b

$a.$

`pow`

(b) ; alternatively $b.$

`rpow`

(a)

$a \ll b$

$a.$

`lshift`

another occurrence. The efficiency of the Boyer-Moore algorithm relies on creating a lookup table that quickly determines where a mismatched character occurs elsewhere in the pattern. In particular, we define a function `last(c)` as follows: If `c` is in `P`, `last(c)` is the index of the last (rightmost) occurrence of `c` in `P`. Otherwise, we conventionally define `last(c) = -1`.

If we assume that the alphabet is of fixed, finite size, and that characters can be converted to indices of an array (for example, by using their character code), the `last` function can be easily implemented as a lookup table with worst-case $O(1)$ -time access to the value `last(c)`. However, the table would have length equal to the size of the alphabet (rather than the size of the pattern), and time would be required to initialize the entire table.

We prefer to use a **hash** table to represent the `last` function, with only those characters from the pattern occurring in the structure. The space usage for this approach is proportional to the number of characters in the pattern, and thus $O(m)$. The expected lookup time remains independent of the problem (although the worst-case bound is $O(m)$). Our complete implementation of the Boyer-Moore pattern-matching algorithm is given in Code Fragment 13.2.

```

1 def find
2   boyer
3   moore(T, P):
4     ...Return the lowest index of T at which substring P begins (or else -1)...
5     3n, m = len(T), len(P) # introduce convenient notations
6     4if m == 0 : return 0 # trivial search for empty string
7     5last = {} # build 'last' dictionary
8     6for k in range(m):
9       7last[P[k]] = k # later occurrence overwrites
10    8# align end of pattern at index m-1 of text
11    9i = m - 1 # an index into T
12    10 k = m - 1 # an index into P
13    11while i < n:
14      12 if T[i] == P[k]: # a matching character
15        13 if k == 0 :
16          14 return i # pattern begins at index i of text
17        15 else:
18          16 i -= 1 # examine previous character
19          17 k -= 1 # offset to hTandP
20        18 else:
21          19 j = last.get(T[i], -1) # last(T[i]) is -1 if not found
22          20 i += m - min(k, j + 1) # case analysis for jump step
23          21 k = m - 1 # restart at end of pattern
24    22return -1

```

Code Fragment 13.2: An implementation of the Boyer-Moore algorithm.