

## Chapter 3

# Container Classes

---

### Objectives

- To understand the list ADT as a general container class for manipulating sequential collections.
- To understand how lists are implemented in Python and the implications this has for the efficiency of various list operations.
- To develop intuition about collection algorithms such as selection sort and use Python operator overloading to make new sortable classes.
- To learn about Python dictionaries as an implementation of a general mapping and understand the efficiency of various dictionary operations.

### 3.1 Overview

Program design gets more interesting when we start considering programs that manipulate large data sets. Typically, we need more efficient algorithms to operate on large collections. Oftentimes the key to an efficient algorithm lies in how the data is organized, that is, the so-called *data structures* on which the algorithms operate. Object-oriented programs often use *container classes* to manage collections of objects. An instance of a container class manages a single collection. Objects can be inserted into and retrieved from the container object at run-time. Python includes a number of container classes as built-in types. You are probably familiar with lists and dictionaries, which are the two main container classes in Python.

In this chapter, we review the basics of Python lists and dictionaries and also take a look at how these containers are implemented in Python. Knowing how



Remove object (*list1.remove(obj1)*) Deletes the first occurrence of *obj1* in *list1*.

### 3.3 A Sequential Collection: A Deck of Cards

```
class Deck(object):

    def __init__(self):
        """post: Create a 52-card deck in standard order"""

    def shuffle(self):
        """Shuffle the deck
        post: randomizes the order of cards in self"""

    def deal(self):
        """Deal a single card
        pre: self is not empty
        post: Returns the next card in self, and removes it from self."""
```













### 3.4.2 Comparing Cards

That leaves us with the problem of putting the hand in order. The sorting problem is an important and well-studied one in computer science. We'll take a quick look at it here and revisit it again in later chapters. If we want to put some things in a particular order, the first problem we have to solve is what exactly the ordering should be.

In the case of our bridge program, we want to order our **Card** objects, grouping them first by suit and then ordering by rank within suit. Usually orderings are determined by a relation such as “less than.” For example, if we say that a list of numbers is in increasing order, that means that for any two numbers  $x$  and  $y$  in the list, if  $x < y$  then  $x$  must precede  $y$  in the list. Similarly, we need a way of comparing cards so that we can order them in our **Hand**. In Chapter 2, we saw how Python operator overloading allows us to build new classes that “act like” existing classes. Here, we would like our cards to behave like numbers so that we can compare them using the standard Python operators such as  $<$ ,  $=$ ,  $>$ , and so on.

We can do this by defining methods for these operations in the `Card` class. Here are the definitions of the “hook” functions for these operators.

```
def __eq__(self, other):

    return (self.suit_char == other.suit_char and
            self.rank_num == other.rank_num)

def __lt__(self, other):

    if self.suit_char == other.suit_char:
        return self.rank_num < other.rank_num
    else:
        return self.suit_char < other.suit_char

def __ne__(self, other):

    return not(self == other)

def __le__(self, other):

    return self < other or self == other
```

Notice that we’ve given “primitive” definitions for `__eq__` and `__lt__`; the rest of the necessary operators can easily be defined in terms of these two. We have not bothered to write definitions for `__gt__` and `__ge__` because Python gives us these for free. In an expression such as `x > y`, when the `>` operator is not implemented







### 3.5 Python List Implementation

When we analyzed the selection sort above, we concentrated on the `max` operation, which turned out to be  $\Theta(n)$ , but we ignored the `insert` and `remove` methods for lists. It turns out that both of these methods have the same time complexity as `max`. How do we know that? Just as the choice of using a Python list to implement our collection classes `Deck` and `Hand` determines the relative efficiency of the methods in these classes, the choice of data structures in the implementation of Python lists determines the efficiency of various list operations. Therefore, understanding the true efficiency of various operations requires some understanding of Python's underlying data structures.

### 3.5.1 Array-based Lists

So how can we efficiently store and access a collection of objects in computer memory? Recall that computer memory is simply a sequence of storage locations. Each storage location has a number associated with it (much like an index) called its *address*. A single data item may be stored across a number of contiguous memory locations. To retrieve an item from memory, we need a way to either look up or compute the starting address of the object. If we want to store a collection of objects, we need to have some systematic method for figuring out where each object in the collection is located.

Consider the case when all of the objects in a collection are the same size, that is they all require the same number of bytes to be stored. This would be the case with a homogeneous (all the same type) collection. A simple method for storing the collection would be to allocate a single contiguous area of memory sufficient to hold the entire collection. The objects could then be stored one after the next. For example, suppose an integer value requires 4 bytes (32 bits) of memory to store. A collection of 100 integers could be stored sequentially into 400 bytes of memory. Let's say the collection of integers starts at the memory location with the address 1024. This means the number at index 0 in the list starts at address 1024, index 1 is at 1028, index 2 is at 1032, etc. The location of the  $i$ th item can be computed simply using the formula *address of  $i$ th* = 1024 + 4 \*  $i$ .

What we have just described is a data structure known as an *array*. Arrays are a common data structure used for storing collections, and many programming languages use arrays as a basic container type. Arrays are very memory efficient and support quick random access (meaning we can “jump” directly to the item we want) via the address calculation we just discussed. By themselves, however, they are somewhat restrictive. One issue is the fact that arrays must generally be



Knowing that Python lists are implemented as dynamically resizing arrays, we are now in a position to analyze the run-time efficiency of various list operations.

The situation for arbitrary insertion operations anywhere in the list is a little different. Because the elements of an array are in contiguous memory locations, to insert into the middle of an array we have to first create a “hole” by shifting all of the following items one place to the right. When the insertion is at the very front of the list, the Python interpreter has to move all  $n$  elements currently in the array. So the insertion operation is still  $\Theta(n)$  even if the size doubling trick is used when the array is full.

Python lists are an example of a sequential data structure. There is an inherent ordering of the data. Even in our implementation of the randomly shuffled deck, the items in the underlying list are still indexed by the natural numbers (0, 1, 2, ...), which gives the collection a natural ordering. In fact, one can view lists abstractly as just a kind of *mapping* from indexes to items in the list. That is, each valid index is associated with (maps to) a particular list item.

Python lists are an example of a sequential data structure. There is an inherent ordering of the data. Even in our implementation of the randomly shuffled deck, the items in the underlying list are still indexed by the natural numbers (0, 1, 2, ...), which gives the collection a natural ordering. In fact, one can view lists abstractly as just a kind of *mapping* from indexes to items in the list. That is, each valid index is associated with (maps to) a particular list item.

The idea of mapping is very general and need not be restricted to using numbers as the indexes. If you think about it a bit, you can probably come up with all sorts of useful collections that involve other sorts of mappings. For example, a phone book is a mapping from names to phone numbers. Mappings pop up everywhere in programming, and that is why Python provides an efficient built-in data structure for managing them, namely a dictionary.

### 3.6.1 A Dictionary ADT

You have probably run across Python dictionaries before, but perhaps not given them much thought. A dictionary is a data structure that allows us to associate keys with values, that is, it implements a mapping. Abstractly, we can think of a dictionary as just a set of ordered (**key**, **value**) pairs. Viewed as an ADT, we just need a few operations in order to have a useful container type.

#### Create

post: Returns an empty dictionary.

#### put(**key**,**value**)

post: The value **value** is associated with **key** in the dictionary. (**key**,**value**) is now the one and only pair in the dictionary having the given key.

#### get(**key**)

pre: There is an **X** such that (**key**,**X**) is in the dictionary.

post: Returns **X**.

#### delete(**key**)

pre: There is an **X** such that (**key**,**X**) is in the dictionary.

post: (**key**,**X**) is removed from the dictionary.

There are many programming situations that call for dictionary-like structures. Some programming languages such as Python and Perl provide built-in implementations of this important ADT. Other languages such as C++ and Java provide them as part of a standard collection library.

### 3.6.2 Python Dictionaries

A Python dictionary provides a particular implementation of the dictionary ADT. Let's start with a simple example. Remember in our **Card** example we needed to be able to turn characters representing suits into full suit names. That's a perfect job for a dictionary. We could define a suitable Python dictionary like this:



As you can see, the syntax for a dictionary literal resembles our abstract description of a dictionary being a set of pairs. In Python, the key-value pairs are joined with a colon. In this case, we are saying that the string "c" maps to the string "Clubs", "d" maps to "Diamonds", etc.

```
>>> suits
{'h': "Hearts", "c": "Clubs", "s": "Spades", "d": "Diamonds"}
>>> suits.get("c")
'Clubs'
>>> suits["c"]
'Clubs'
>>> suits["s"]
'Spades'
>>> suits["j"]
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
KeyError: "j"
>>> suits.get("j")
>>> suits.get("x", "Not There")
'Not There'
```

The abstract **put** operation for changing entries in a dictionary or extending it with new entries is implemented via assignment in Python. Again, this makes the syntax for working with dictionaries very similar to that of lists. Here are a few examples:



```
h Hearts
c Clovers
s Shovels
d Diamonds
>>> 'c' in suits
True
>>> 'x' in suits
False
```

### 3.6.3 Dictionary Implementation

As with virtually any ADT, there are numerous ways one could go about implementing dictionaries. The choice of implementation will determine how efficient the various operations will be. One simple representation would be to store the dictionary entries as a list of key-value pairs. A `get` operation would involve some form of lookup on the list to find the pair with the specified key. Other operations could also be performed using simple list manipulation. Unfortunately, this approach will not be very efficient, as some of the operations will require  $\Theta(n)$  effort. (An exact analysis of the situation is left as an exercise.)

Python uses a more efficient data structure called a *hash table*. Hash tables are covered in-depth in section 13.5. Here we just want to give you some intuition so that you can understand the efficiency of various dictionary operations. That will enable you to judge the efficiency of algorithms that use Python dictionaries.

The heart of a hash table is a *hashing function*. A hashing function takes a key as a parameter and performs some simple calculations on it to produce a number. Since all data on the computer is ultimately stored as bits (binary numbers), it's pretty easy to come up with hashing functions. Python actually has a built-in function `hash` that does this. You can try it out interactively.



But what happens when two keys hash to the same spot? This is called a *collision*. Dealing with collisions is an important issue that is covered in section 13.5. For now it suffices to note that there are good techniques for dealing with this problem. Using these techniques and ensuring a table of adequate size yields data structures that, in practice, allow for constant time operations. Python dictionaries are very efficient and can easily handle thousands, even millions of entries, provided you have enough memory available. The Python interpreter itself relies heavily on the use of dictionaries to maintain namespaces, so the dictionary implementation has been highly optimized.

Let’s put our new knowledge of dictionaries to use in a program that combines several Python container classes to build Markov models. A Markov model is a statistical technique for modelling systems that change over time. One application of Markov models is in the area of systems for natural language understanding. For example, a speech recognition system can use predictions about what word is likely to come next in a sentence in order to decide among homonyms such as “their,” “they’re,” and “there.”

The basic idea behind a Markov model of language is that one can make predictions about the next word of an utterance by looking at some small sequence of preceding words. For example, a trigram model looks at the preceding two words to predict the next (third) word in a sequence. More or fewer words could be used as a “window” depending on the application. For example, a bigram model would predict the probabilities for the next word based only on the immediately preceding word. Our initial design will be for a trigram model; extending the program to arbitrary length prefixes is left as an exercise.

02/10/2018 - RS000000000000000000000000381359 - Data Structures and Algorithms Using Python and C++



We are now in a position to write the code for this class.

02/10/2018 - RS0000000000000000000000000381359 - Data Structures and Algorithms Using Python and C++





of cup,' sort!' forehead you However, house the went to me unhappy  
up impossible settled We had help the always in see, forgot tree of you  
'for night? because hadn't her ear. all confused sit took the care went  
quite do up, 'How three An Turtle, the was soldiers, solemnly, so went of  
the sharply. to Rabbit 'Tis there last, a with that o'clock and below, he  
Writhing, don't to wig, she into three, But said there,' offended. turning  
This some (she together.'" such be because and what the had to hatters  
better This Mouse new said the pool whiting. with could from bank-the  
mile said I she all! turning when 'Begin By how as head them, little,  
and Latitude he

## 3.7 Chapter Summary

- Container objects are used to manage collections. Items can be added to and removed from containers at run-time.
- The built-in Python list is an example of a container class.
- Lists define a sequential collection where there is a first item and each item (except the last) has a natural successor.
- Lists can be used to store both sorted and unsorted sequences. Selection sort is a  $\Theta(n^2)$  algorithm for sorting a sequence.
- Python lists are implemented using arrays of references. When a list grows too large for the current array, Python automatically allocates a new larger one. This technique allows **append** operations to be done in  $\Theta(1)$  (amortized) time, but operations that insert or delete items in the midst of the list require  $\Theta(n)$  time.





9. Which of the following is not true of Python dictionaries?
  - a) They are implemented as hash tables.
  - b) Values must be immutable.
  - c) Lookup is very efficient.
  - d) All of the above are true.
10. A trigram model of natural language
  - a) uses a prefix of three words to predict the next word.
  - b) uses a prefix of two words to predict the next word.
  - c) is more useful than a Markov model.
  - d) is used to send money overseas.

### Short-Answer Questions

1. Using the **Deck** and **Hand** classes from this chapter, write snippets of code to do each of the following:
  - a) Print out the names of all 52 cards.
  - b) Print out the names of 13 random cards.
  - c) Choose 13 cards at random from a 52-card deck and show the cards in value order (Bridge hand order).
  - d) Deal and display four 13-card hands dealt from a shuffled deck.
2. What is the run-time efficiency ( $\Theta$ ) of the two shuffling algorithms discussed in the chapter (using two lists vs. in place). The discussion suggested that the latter is more efficient. Is this consistent with your  $\Theta$  analysis? Explain.
3. Suppose you are involved in designing a system that must maintain information about a large number of individuals (for example, customer records or health records). Each person will be represented with an object that contains all of their critical information. Your job is to design a container class to hold all of these records. The following operations must be supported:

`add(person)` – adds `person` object to the collection

`remove(name)` – removes the person named `name` from the collection.

`lookup(name)` – returns the record for the person named `name`.

`list_all` – returns a list of all the records in the collection in order by name.

- (a) The objects are stored in a Python list in the order that they are added.
- (b) The objects are stored in a Python list in order by name.
- (c) The objects are stored in a Python dictionary indexed by name.

- (a) an unordered Python list.
- (b) a sorted Python list.
- (c) a Python dictionary. (Note: the elements of the set will be the keys, you can just use `None` or `True` as the value.)

1. Modify the `Deck` class to keep track of the current size of the deck using an instance variable. Does this change the run-time efficiency of the `size` operation? Do a bit of research to answer this question.
2. Look into the functions provided by the Python `random` module to simplify the shuffling code in the `Deck` class.
3. Suppose we want to be able to place cards back into a deck. Modify the `Deck` class to include the operations `addTop`, `addBottom`, and `addRandom` (the last one inserts the card at a random location in the deck).

4. Instead of shuffling a deck of cards, another way to get a random distribution is to deal cards from random locations of an ordered deck. Implement a `Deck` class that uses this approach. Analyze the efficiency of the operations you provide.
5. It can be inconvenient to test programs involving decks of cards if cards are always dealt in random order. One solution is to allow the deck to be “stacked” in a particular order. Design a `Deck` class that allows the contents of the deck to be read from a file.
6. Modify the `sort` method in the `Hand` class so that it sorts the hand “in place.” Hint: look at the in-place shuffling algorithm.
7. Another way to put a hand in order is to place each card into its proper location as it is added to the hand. This algorithm is called an *insertion sort*. Implement a version of `Hand` that uses this method to keep the hand in order.
8. Implement an extended `Deck` class with operations suitable for playing the card game war. You will need to be able to create an empty deck and place cards into it.
9. Write a program to play the following simple solitaire game.  $N$  cards are dealt face up onto the table. If two cards have a matching rank, new cards are dealt face up on top of them. Dealing continues until the deck is empty or no two stacks have matching ranks. The player wins if all the cards are dealt. Run simulations to find the probability of winning with various values of  $N$ .
10. Write a program that deals and evaluates poker hands.
11. Write a program to simulate the game of blackjack.
12. Write a program to deal and evaluate bridge hands to determine if they have an opening bid.
13. Modify the Markov gibberish generator so that it works at the level of characters rather than words. Note: you should not need to modify the class to do this, only how it is used.
14. Extend the Markov gibberish generator to allow the size of the prefix to be determined when the model is created. The constructor will take a parameter specifying the length of the prefix. Experiment with different prefix lengths on texts of various size to see what happens. Combining this with the previous project, you can produce a very versatile and entertaining gibberish generator.

- 02/10/2018 - RS000000000000000000000000381359 - Data Structures and Algorithms Using Python and C++





# Chapter 4

## Linked Structures and Iterators

## Objectives

- ## 4.1 Overview

107



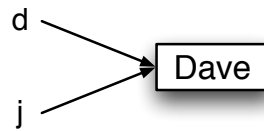


Figure 4.1: Two variables assigned to an object

same object as the name `d`; it does not create a new string object. A good analogy is to think of assignment as placing a sticky note with the name written on it onto the object. At this point, the data object `Dave` has two sticky notes on it: one with the name `d` and one with the name `j`. Figure 4.1 should help clarify what is happening. In diagrams such as this, we use an arrow as an intuitive way to show the “value” of a reference; the computer actually stores a number that is the address of what our arrow is pointing to.

Of course, the Python interpreter can’t use sticky notes, it keeps track of these associations internally using the namespace dictionary. We can actually access that dictionary with the built-in function called `locals()`.

```
>>> print locals()
{'__builtins__': <module '__builtin__' (built-in)>, '__name__': '__main__',
 'j': 'Dave', '__doc__': None, 'd': 'Dave'}
```

In this example, you can see that the local dictionary includes some Python special names (`__builtins__`, `__name__`, and `__doc__`) some of which you may recognize. We’re not really concerned about those here. The point is that our assignment statements added the two names `d` and `j` to the dictionary. Notice, when the dictionary is printed, Python shows us the names as keys and a representation of the actual data objects to which they map as values. Keep in mind that the namespace dictionary actually stores the address of the object (also called a *reference* to the object). Since we usually care about the data, not locations, the Python interpreter automatically shows us a representation of what is stored at the address, not the address itself.

If, out of curiosity, we should want to find the actual address of an object, we can do that. The Python `id` function returns a unique identifier for each object; in most versions of Python, the `id` function returns the memory address where the object is stored.

```
>>> print id(d), id(j)
432128 432128
```

As you can see by the output of the `id` function, after the assignment statement `j = d`, both the names `j` and `d` refer to the same data object. Internally, the Python interpreter keeps track of the fact that there are two references to the string object containing "Dave"; this is referred to as the *reference count* for the object.

Continuing with the example, let's do a couple more assignments.

```
>>> j = 'John'
>>> print id(d), id(j)
432128 432256
>>> d = 'Smith'
>>> print id(d), id(j)
432224 432256
```

When we assign `j = 'John'`, a new string object containing "John" is created. Using our sticky note analogy, we have moved the sticky note `j` to the newly created data object containing the string "John". The output of the `id` function following the statement `j = 'John'` shows that the name `d` still refers to the same object as before, but the name `j` now refers to an object at a different memory location. The reference count for each of the two string objects is now one.

The statement `d = 'Smith'` makes the name `d` refer to a new string object containing "Smith". Note that the address for the string object "Smith" is different from the string object "Dave". Again, the address that the name maps to changes when the name is assigned to a different object. This is an important point to note: *Assignment changes what object a variable refers to, it does not have any effect on the object itself.* In this case, the string "Dave" does not change into the string "Smith", but rather a new string object is created that contains "Smith".

At this point, nothing refers to the string "Dave" so its reference count is now zero. The Python interpreter automatically deallocates the memory for the string object containing "Dave", since there is no longer a way to access it. By deallocating objects that can no longer be accessed (when their reference count changes to zero), the Python interpreter is able to reuse the same memory locations for new objects later on. This process is known as *garbage collection*. Garbage collection adds some overhead to the Python interpreter that slows down execution. The gain is that it relieves the programmer from the burden of having to worry about memory allocation and deallocation, a process that is notoriously knotty and error prone in languages that do not have automatic memory management.

It is also possible for the programmer to explicitly remove the mapping for a given name.

The statement `del d` removes the name `d` from the namespace dictionary so it can no longer be accessed. Attempting to execute the statement `print d` now would cause a `NameError` exception to be raised just as if we had never assigned an object to `d`. Removing that name reduces the reference count for the string `"Smith"` from one to zero so it will now also be garbage collected.

The Python memory model also makes assignment a very efficient operation. An expression in Python always evaluates to a reference to some object. Assigning the result to a name simply requires that the name be added to the namespace dictionary (if it's not already present) along with the four- or eight-byte reference. In a simple assignment like `j = d` the effect is to just copy `d`'s reference over to `j`'s namespace entry.

```
>>> lst1 = [1, 2, 3]
>>> lst2 = lst1
>>> lst2.append(4)
>>> lst1
[1, 2, 3, 4]
```

02/10/2018 - RS0000000000000000000000000381359 - Data Structures and Algorithms Using Python and C++

aliasing crop up only when the shared object happens to be mutable. Things like strings, ints, and floats simply can't change, so aliasing is not an issue for these types.

When we want to avoid the side effects of aliasing, we need to make separate copies of an object so that changes to one copy won't affect the others. Of course a complex object such as a list might itself contain references to other objects, and we have to decide how to handle those references in the copying process. There are two different types of copies known as *shallow copies* and *deep copies*. A shallow copy has its own top-level references, but those references refer to the same objects as the original. A deep copy is a completely separate copy that creates both new references and, where necessary, new data objects at all levels. The Python `copy` module contains useful functions for copying arbitrary Python objects. Here's an interactive example using lists to demonstrate.

```
>>> import copy
>>> b = [1, 2, [3, 4], 6]
>>> c = b
>>> d = copy.copy(b)      # creates a shallow copy
>>> e = copy.deepcopy(b)  # creates a deep copy
>>> print b is c, b == c
True True
>>> print b is d, b == d
False True
>>> print b is e, b == e
False True
```

In this code, `c` is the same list as `b`, `d` is a shallow copy, and `e` is a deep copy. By the way, there are numerous ways to get a shallow copy of a Python list. We could also have used slicing (`d = b[:]`) or list construction (`d = list(b)`) to create a shallow copy.

So what's up with the output? The Python `is` operator tests whether two expressions refer to the exact same object, whereas the Python `==` operator tests to see if two expressions yield equivalent data. That means `a is b` implies `a == b` but not vice versa. In this example, you can see that assignment does not create a new object since `b is c` holds after the initial assignment. However both the shallow copy `d` created by slicing and the deep copy `e` are distinct new objects that contain equivalent data to `b`. While these copies contain equivalent data, their internal structures are not identical. As depicted in Figure 4.2, the shallow copy simply contains a copy of the references at the top level of the list, while the deep copy contains a copy of the mutable parts of the structure at all levels. Notice that the

