

Data structures in Python

Lecture 04

Dr. Colin Rundel

Dictionaryes

Dictionaries

Python `dicts` are a *heterogeneous, ordered**, *mutable* containers of key value pairs. Each entry consists of a *key* (immutable) and a *value* (anything) - they are designed for the efficient lookup of values using a key.

`dicts` are typically constructed using `{}` with `:`

```
1 {'abc': 123, 'def': 456}
```

```
{'abc': 123, 'def': 456}
```

or via `dict()` using a list of key value tuples,

```
1 dict([('abc', 123), ('def', 456)])
```

```
{'abc': 123, 'def': 456}
```

If your keys are strings, then you can pass keyword arguments to `dict()`,

```
1 dict(hello=123, world=456) # can't use `def` here as it is reserved keyword
```

```
{'hello': 123, 'world': 456}
```

Allowed key values

The keys for a `dict` must be *immutable* objects (e.g. integer, double, string, or tuple). Values may be of any type (mutable or immutable).

```
1 {1: "abc", 1.1: (1,1), "one": ["a","n"], (1,1): lambda x: x**2}
```

```
{1: 'abc', 1.1: (1, 1), 'one': ['a', 'n'], (1, 1): <function <lambda> at 0x11cca6e80>}
```

Using a mutable object (e.g. a list) as a key will result in an error,

```
1 {[1]: "bad"}
```

```
TypeError: unhashable type: 'list'
```

when using a tuple, you need to be careful that all elements are also immutable,

```
1 {(1, [2]): "bad"}
```

```
TypeError: unhashable type: 'list'
```

dict “subsetting”

The `[]` operator exists for `dicts` but is used for key-based value look ups,

```
1 x = {1: 'abc', 'y': 'hello', (1,1): 3.14159}
```

```
1 x[1]
```

'abc'

```
1 x['y']
```

'hello'

```
1 x[(1,1)]
```

3.14159

```
1 x[0]
```

KeyError: 0

```
1 x['def']
```

KeyError: 'def'

Value inserts & replacement

Since dictionaries are mutable, it is possible to insert new key value pairs as well as replace the value associated with an existing key.

```
1 x = {1: 'abc', 'y': 'hello', (1,1): 3.14159}
```

New key/value pairs can be inserted,

```
1 x['def'] = -1
2 x
```

```
{1: 'abc', 'y': 'hello', (1, 1): 3.14159, 'def': -1}
```

Existing values can also be replaced,

```
1 x['y'] = 'goodbye'
2 x
```

```
{1: 'abc', 'y': 'goodbye', (1, 1): 3.14159, 'def': -1}
```

Removing keys

```
1 x
```

```
{1: 'abc', 'y': 'goodbye', (1, 1): 3.14159, 'def': -1}
```

The `del` can be used to remove a key and its value,

```
1 del x[(1,1)]  
2 x
```

```
{1: 'abc', 'y': 'goodbye', 'def': -1}
```

All keys and values can be removed via the `clear()` method,

```
1 x.clear()  
2 x
```

```
{}
```

Assigning `None` to a key does not remove it,

```
1 x[1] = None  
2 x
```

```
{1: None}
```

Other common methods

```
1 x = {1: 'abc', 'y': 'hello'}
```

```
1 len(x)
```

2

```
1 list(x)
```

```
[1, 'y']
```

```
1 tuple(x)
```

```
(1, 'y')
```

```
1 1 in x
```

True

```
1 'hello' in x
```

False

```
1 x.keys()
```

```
dict_keys([1, 'y'])
```

```
1 x.values()
```

```
dict_values(['abc', 'hello'])
```

```
1 x.items()
```

```
dict_items([(1, 'abc'), ('y', 'hello')])
```

```
1 x | {(1,1): 3.14159}
```

```
{1: 'abc', 'y': 'hello', (1, 1): 3.14159}
```

```
1 x | {'y': 'goodbye'}
```

```
{1: 'abc', 'y': 'goodbye'}
```

Iterating dictionaries

Dictionaries can be used with for loops (and comprehensions). They will both iterate over the *keys* only. To iterate over the *keys* and *values* use `items()`.

```
1 for z in {1: 'abc', 'y': 'hello'}:  
2     print(z)
```

```
1  
y
```

```
1 [z for z in {1: 'abc', 'y': 'hello'}]
```

```
[1, 'y']
```

```
1 for k,v in {1: 'abc', 'y': 'hello'}.items():  
2     print (k,v)
```

```
1 abc  
y hello
```

```
1 [(k,v) for k,v in {1: 'abc', 'y': 'hello'}.items()]
```

```
[(1, 'abc'), ('y', 'hello')]
```

Exercise 1

Write a function that takes two dictionaries as arguments and merges them into a single dictionary. If there are any duplicate keys, the value from the second dictionary should be used.

```
1 x = {"a": 1, "b": 2, "c": 3}
2 y = {"c": 5, "d": 6, "e": 7}
3
4 def merge(d1, d2):
5     return None
```

Sets

Sets

In Python a `set` is a *heterogeneous, unordered, mutable* container of **unique** immutable elements.

A `set` is constructed using `{}` (without using `:`) or via `set()`,

```
1 {1,2,3,4,1,2}
```

```
{1, 2, 3, 4}
```

```
1 set((1,2,3,4,1,2))
```

```
{1, 2, 3, 4}
```

```
1 set("mississippi")
```

```
{'i', 'p', 'm', 's'}
```

All of the elements must be immutable (and therefore hashable),

```
1 {1,2,[1,2]}
```

```
TypeError: unhashable type: 'list'
```

Subsetting sets

Sets do not use the `[]` operator for element checking or removal,

```
1 x = set(range(5))  
2 x
```

```
{0, 1, 2, 3, 4}
```

```
1 x[4]
```

```
TypeError: 'set' object is not subscriptable
```

```
1 del x[4]
```

```
TypeError: 'set' object doesn't support item deletion
```

Modifying sets

Sets have their own special methods for adding and removing elements,

```
1 x = set(range(5))  
2 x
```

{0, 1, 2, 3, 4}

```
1 x.add(9)  
2 x
```

{0, 1, 2, 3, 4, 9}

```
1 x.remove(9)  
2 x.remove(8)
```

KeyError: 8

```
1 x
```

{0, 1, 2, 3, 4}

```
1 x.discard(0)  
2 x.discard(8)  
3 x
```

{1, 2, 3, 4}

Set operations

```
1 x = set(range(5))  
2 x
```

{0, 1, 2, 3, 4}

```
1 3 in x
```

True

```
1 x.isdisjoint({1,2})
```

False

```
1 x <= set(range(6))
```

True

```
1 x >= set(range(3))
```

True

```
1 5 in x
```

False

```
1 x.isdisjoint({5})
```

True

```
1 x.issubset(range(6))
```

True

```
1 x.issuperset(range(3))
```

True

Set operations (cont)

```
1 x = set(range(5))  
2 x
```

{0, 1, 2, 3, 4}

```
1 x | set(range(10))
```

{0, 1, 2, 3, 4, 5, 6, 7, 8, 9}

```
1 x & set(range(-3,3))
```

{0, 1, 2}

```
1 x - set(range(2,4))
```

{0, 1, 4}

```
1 x ^ set(range(3,9))
```

{0, 1, 2, 5, 6, 7, 8}

```
1 x.union(range(10))
```

{0, 1, 2, 3, 4, 5, 6, 7, 8, 9}

```
1 x.intersection(range(-3,3))
```

{0, 1, 2}

```
1 x.difference(range(2,4))
```

{0, 1, 4}

```
1 x.symmetric_difference(range
```

{0, 1, 2, 5, 6, 7, 8}

More comprehensions

It is possible to use comprehensions with both `sets` and `dicts`,

Set comprehension,

```
1 {x.lower() for x in "The quick brown fox jumped a lazy dog"}
```

```
{'h', 'u', 'k', 'm', 'd', 'y', 'z', ' ', 'l', 'b', 'r', 'w', 'e'}
```

Dict comprehension,

```
1 names = ["Alice", "Bob", "Carol", "Dave"]
2 grades = ["A", "A-", "A-", "B"]
3
4 {name: grade for name, grade in zip(names, grades)}
```

```
{'Alice': 'A', 'Bob': 'A-', 'Carol': 'A-', 'Dave': 'B'}
```

tuple comprehensions

Note that `tuple` comprehensions do not exist,

```
1 # Not a tuple
2 (x**2 for x in range(5))
```

<generator object <genexpr> at 0x11cd99a40>

the above code constructs a generator (more on these later).

If necessary you can use a list comprehension and then cast to a tuple,

```
1 # Is a tuple – via casting a list to tuple
2 tuple([x**2 for x in range(5)])
```

(0, 1, 4, 9, 16)

```
1 tuple(x**2 for x in range(5))
```

(0, 1, 4, 9, 16)

deques (double ended queue)

Dequeues are *heterogeneous, ordered, mutable* collections of elements and behave in much the same way as `lists`. They are designed to be efficient for adding and removing elements from the beginning and end of the collection.

These are not part of the base language and are available as part of the built-in `collections` module. We will discuss libraries / modules next time, for now to get access we will import the `deque` function from `collections`.

```
1 from collections import deque
2 deque("A", 2, True)
```

```
deque(['A', 2, True])
```

growing and shrinking

```
1 x = deque(range(3))
2 x
```

```
deque([0, 1, 2])
```

Values can be added to the beginning and end via `.appendleft()` and `.append()` respectively,

```
1 x.appendleft(-1)
2 x.append(3)
3 x
```

```
deque([-1, 0, 1, 2, 3])
```

values can be removed via `.popleft()` and `.pop()`,

```
1 x.popleft()
```

```
-1
```

```
1 x.pop()
```

```
3
```

```
1 x
```

```
deque([0, 1, 2])
```

maxlen

`deque`s can be constructed with an optional `maxlen` argument which sets their maximum size - if this is exceeded values from the opposite end will be dropped.

```
1 x = deque(range(3), maxlen=4)
2 x
```

```
deque([0, 1, 2], maxlen=4)
```

```
1 x.append(0)
2 x
```

```
deque([0, 1, 2, 0], maxlen=4)
```

```
1 x.append(0)
2 x
```

```
deque([1, 2, 0, 0], maxlen=4)
```

```
1 x.append(0)
2 x
```

```
deque([2, 0, 0, 0], maxlen=4)
```

```
1 x.appendleft(-1)
2 x
```

```
deque([-1, 2, 0, 0], maxlen=4)
```

```
1 x.appendleft(-1)
2 x
```

```
deque([-1, -1, 2, 0], maxlen=4)
```

```
1 x.appendleft(-1)
2 x
```

```
deque([-1, -1, -1, 2], maxlen=4)
```

Some algorithms and data structures basics

Big-O notation

This is a tool that is used to describe the complexity (usually in time but also in memory usage) of an algorithm. The goal is to broadly group algorithms based on how their complexity grows as the size of an input grows.

Consider a mathematical function that exactly captures this relationship (e.g. the number of steps in a given algorithm given an input of size n). The Big-O value for that algorithm will then be the largest term involving n in that function.

Generally algorithms will vary depending on the exact nature of the data and so often we talk about Big-O in terms of expected complexity and worse case complexity, we also often consider amortization for these worst cases.

Common terminology

Complexity	Big-O
Constant	$O(1)$
Logarithmic	$O(\log n)$
Linear	$O(n)$
Quasilinear	$O(n \log n)$
Quadratic	$O(n^2)$
Cubic	$O(n^3)$
Exponential	$O(C^n)$

Note that these terms ignore all smaller order terms as well as the constant in front of the highest order term (this can be significant in practice).

Vector / Array

Linked List

Hash table

Time complexity in Python

Operation	list (array)	dict (& set)	deque
Copy	$O(n)$	$O(n)$	$O(n)$
Append	$O(1)$	—	$O(1)$
Insert	$O(n)$	$O(1)$	$O(n)$
Get item	$O(1)$	$O(1)$	$O(n)$
Set item	$O(1)$	$O(1)$	$O(n)$
Delete item	$O(n)$	$O(1)$	$O(n)$
<code>x in s</code>	$O(n)$	$O(1)$	$O(n)$
<code>pop()</code>	$O(1)$	—	$O(1)$
<code>pop(0)</code>	$O(n)$	—	$O(1)$

Exercise 2

For each of the following scenarios, which is the most appropriate data structure and why?

- A fixed collection of 100 integers.
- A queue (first in first out) of customer records.
- A stack (first in last out) of customer records.
- A count of word occurrences within a document.
- The heights of the bars in a histogram with even binwidths.

Data structures in R

To tie things back to Sta 523 - the following R objects are implemented using the following data structures.

- Atomic vectors - Array of the given type (int, double, etc.)
- Generic vectors (lists) - Array of SEXP's (R object pointers)
- Environments - Hash map with string-based keys
- Pairlists - Linked list