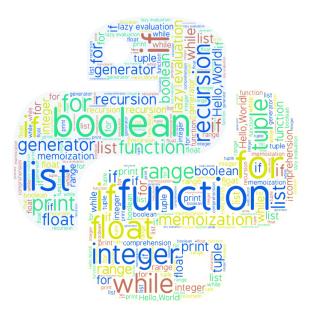
## Python 101

Lec04 Useful data structures and modules

thoum

August 19, 2019

## What have we learned so far?



# Turing Complete

In theory, we *can* do everything<sub>a computer can</sub> with what we have learned so far. That does not imply efficiency.

Python provides datatypes and $MANY$ useful modules that will help us. Today, we will take a peek at some of them.

## Container Datatypes

Things that holds other things (e.g. tuples, lists, etc.).

So far, we have used lists and tuples to store series of values. But there are different containers for different times.

## **Dictionary**

Similar to a list, but instead of *int* as indices, we can use objects that are hashable as indices(a.k.a. key).

Instead of comparing the key with every key it has, dictionary usually  $^{1}$  does it efficiently.

<sup>&</sup>lt;sup>1</sup>Hash Collision

## Hashing

Mapping data of arbitrary size onto data of a fixed size. Mapping words to its first letters, mapping students to student IDs, mapping your socks to the colors of white, grey and black to pair them, they are all hashing.

## Hashing Example

```
a = 1, b = 2, ... z = 26
score of a word = \Sigma value(c_i) \mod 101
```

```
# maps word to number in between 0~100

def score(word):
    s = sum((ord(c)-ord('a')+1 for c in word))
    return s % 101

print(score('hardwork'))
print(score('luck'))
print(score('luck'))
# kind of how dictionary works(ignore collision for now)
lst = [0] * 101

lst[score('newton')] = 'gravity'
lst[score('einstein')] = 'e=mc2'

print(lst[score('einstein')])
print(lst[score('newton')])
```

## Hashing

Hashing by itself is an important topic with wide range of usage(Encryption, Bitcoins...), but will not go into further details.

## **Using Dictionaries**

```
avengers = {} # use curly braces to create a dictionary
# the names are the keys, and the hero identity are values
avengers['brucebanner'] = 'hulk'
avengers['tonystark'] = 'ironman'
avengers['peterparker'] = 'spiderman'
avengers['steverogers'] = 'captainamerica'
print(avengers['tonystark'])
# overwrite key's value
avengers['steverogers'] = 'blueskull'
# no key: error
# print(avengers['thanos'])
```

## **Using Dictionaries**

```
# dictionaries are iterable
print(len(d))
# iterate over keys
for k in d:
    print(k)

# iterate over values
for v in d.values():
    print(v)

# get the key:value pair
for k, v in d.items():
    print(k, v)
```



Unordered collection of *unique* hashable objects.

# set examples

#### removing duplicates

```
lst = [1,1,1,1,5,1,2,3,10,10]
a = set(lst)
print(a)
```

## set examples

#### testing membership

list version: 3.6s, set version: 0.05s

```
import time
import random

1 = list(range(10000))
s = set(1)

start_time = time.time()
for _ in range(50000):
    random.randint(-10000, 10000) in 1
end_time = time.time()
print("search in list took: ", end_time - start_time)

start_time = time.time()
for _ in range(50000):
    random.randint(-10000, 10000) in s
end_time = time.time()
print("search in set took: ", end_time - start_time)
```

## set examples

For other operations like union(), intersection(), difference(), look here. Or

#### dir(set)

to see the methods available.

#### Counter

#### Counts occurrences

```
from collections import Counter
words = ['red', 'blue', 'red', 'green', 'blue', 'blue']
cnt = Counter(words)

print(cnt['green'])
print(cnt['chicken'])
print(cnt.most_common(2))
```

Other operations are on the website.

### deque

Double Ended Queue

Use instead of *list* when inserting, popping happens at the beginning of the list. *e.g.* keeping track of values for moving average?

Insert, pop at the beginning of the list creates an overhead of shifting every other element to the left.

#### deque vs list

#### list:3.920s, deque: 0.015s

```
rom collections import deque
import time
q = deque() # an empty deque
    1 = []
        1.insert(0, i)
    q = deque() # an empty deque
       q.appendleft(i)
  __name__ == '__main__':
   print(timeit.timeit("f()", setup="from __main__ import f",number=2))
    print(timeit.timeit("g()", setup="from __main__ import g",number=2))
```

## deque vs list

# List outperforms deque in random access list:0.111s, deque: 3.597s

```
from collections import deque
import random
import timeit

1 = list(range(1000000))
q = deque(range(1000000))

def f():
    return 1[random.randint(0,999999)]

def g():
    return q[random.randint(0,999999)]

if __name__ == '__main__':
    print(timeit.timeit("f()", setup="from __main__ import f",number=100000))
    print(timeit.timeit("g()", setup="from __main__ import g",number=100000))
```

What to use?

In choosing the right *thing* to use, having a slight idea of complexity would help.



How long does an operation take? How much memory does it use?

## Time Complexity

When input size is N Some operation might take constant time.

e.g. lst[i]

Some operation might take time proportional to N.

e.g. iterating over a list

Some operation might take time proportional to  $N^2$ .

e.g. matrix multiplication

Some operation might take time proportional to Nlog(N).

e.g. sorting

#### Pool

Parallel Computing is an advanced topic, as it requires careful coding to avoid errors like deadlock. (e.g. dining philosophers) But, Python's Pool allows us to do simple parallel computing.(use with care)

#### Pool

```
from multiprocessing import Pool

def f(n):
    print(n)
    return n*n

l = [1,2,3,4,5,6,7,8]

with Pool() as p:
    # the order of execution is not guaranteed
    result = p.map(f, 1)
print(result)
```

#### Pool

The order of execution is not guaranteed, although the end result would be ordered.

```
-/wor../pyt../lec05 <master> > python3 pool.py
[1, 4, 9, 16, 25, 36, 49, 64]
~/wor../pyt../lec05 <master> > python3 pool.py
[1, 4, 9, 16, 25, 36, 49, 64]
~/wor../pyt../lec05 <master> > python3 pool.py
   4, 9, 16, 25, 36, 49, 64]
```

# Pool Explained

Pool(4) means we will use 4 cpus. Pool() defaults to os.cpu\_count().(Mine has 4)

## Pool Explained

with keyword handles initalization and cleanup of certain operations

```
# p is only effective within the 'with'
with Pool() as p:
    # do something

p = Pool()
# do something
p.close()
p.join()
```

#### Pool Performance

```
rom multiprocessing import Pool
import time
def is_prime(n):
    if n & 1 == 0:
    d = 3
    while d * d <= n:
       if n % d == 0:
       d = d + 2
start_time = time.time()
result = [is_prime(n) for n in 1]
end_time = time.time()
print("single took: ", end_time-start_time)
start_time = time.time()
with Pool() as p:
    result = p.map(is_prime, 1)
end_time = time.time()
print("multi took: ", end_time-start_time)
```