

Lecture 8, 1 September 2022

Lists, arrays, dictionaries: implementation details

- What are the salient differences?
- How are they stored?
- What is the impact on performance?

Arrays

- Contiguous block of memory
- Typically size is declared in advance, all values are uniform
- $a[0]$ points to first memory location in the allocated block
- Locate $a[i]$ in memory using index arithmetic
 - Skip i blocks of memory, each block's size determined by value stored in array
- **Random access** -- accessing the value at $a[i]$ does not depend on i
- Useful for procedures like sorting, where we need to swap out of order values $a[i]$ and $a[j]$
 - $a[i], a[j] = a[j], a[i]$
 - Cost of such a swap is constant, independent of where the elements to be swapped are in the array
- Inserting or deleting a value is expensive
- Need to shift elements right or left, respectively, depending on the location of the modification

Lists

- Each location is a *cell*, consisting of a value and a link to the next cell
 - Think of a list as a train, made up of a linked sequence of cells
- The name of the list l gives us access to $l[0]$, the first cell
- To reach cell $l[i]$, we must traverse the links from $l[0]$ to $l[1]$ to $l[2]$... to $l[i-1]$ to $l[i]$
 - Takes time proportional to i
- Cost of swapping $l[i]$ and $l[j]$ varies, depending on values i and j
- On the other hand, if we are already at $l[i]$ modifying the list is easy
 - *Insert* - create a new cell and reroute the links
 - *Delete* - bypass the deleted cell by rerouting the links
- Each insert/delete requires a fixed amount of local "plumbing", independent of where in the list it is performed

Dictionaries

- Values are stored in a fixed block of size m
- Keys are mapped to $\{0, 1, \dots, m-1\}$
- Hash function $h: K \rightarrow S$ maps a large set of keys K to a small range S
- Simple hash function: interpret $k \in K$ as a bit sequence representing a number n_k in binary, and compute $n_k \bmod m$, where $|S| = m$
- Mismatch in sizes means that there will be *collisions* -- $k_1 \neq k_2$, but $h(k_1) = h(k_2)$
- A good hash function maps keys "randomly" to minimize collisions
- Hash can be used as a *signature* of authenticity
 - Modifying k slightly will drastically alter $h(k)$
 - No easy way to reverse engineer a k' to map to a given $h(k)$
 - Use to check that large files have not been tampered with in transit, either due to network errors or malicious intervention
- Dictionary uses a hash function to map key values to storage locations
- Lookup requires computing $h(k)$ which takes roughly the same time for any k
 - Compare with computing the offset $a[i]$ for any index i in an array
- Collisions are inevitable, different mechanisms to manage this, which we will not discuss now
- Effectively, a dictionary combines flexibility with random access

Lists in Python

- Flexible size, allow inserting/deleting elements in between
- However, implementation is an array, rather than a list
- Initially allocate a block of storage to the list
- When storage runs out, double the allocation
- $l.append(x)$ is efficient, moves the right end of the list one position forward within the array
- $l.insert(0, x)$ inserts a value at the start, expensive because it requires shifting all the elements by 1
- We will run experiments to validate these claims

Measuring execution time

- Call `time.perf_counter()`
- Actual return value is meaningless, but difference between two calls measures time in seconds

In [1]: `import time`

- 10^7 appends to an empty Python list

In [2]: `start = time.perf_counter()
l = []
for i in range(10000000):
 l.append(i)
elapsed = time.perf_counter() - start
print(elapsed)`

0.8226494059999823

- Doubling the work approximately doubles the time, linear

In [3]: `start = time.perf_counter()
l = []
for i in range(20000000):
 l.append(i)
elapsed = time.perf_counter() - start
print(elapsed)`

1.8082762559999992

- 10^5 inserts at the beginning of a Python list

In [4]: `start = time.perf_counter()
l = []
for i in range(100000):
 l.insert(0,i)
elapsed = time.perf_counter() - start
print(elapsed)`

3.061872225000002

- Doubling and tripling the work multiplies the time by 4 and 9, respectively, so quadratic

In [5]: `start = time.perf_counter()
l = []
for i in range(200000):
 l.insert(0,i)
elapsed = time.perf_counter() - start
print(elapsed)`

11.92333228299998

In [6]: `start = time.perf_counter()
l = []
for i in range(300000):
 l.insert(0,i)
elapsed = time.perf_counter() - start
print(elapsed)`

26.641093278000028

- Creating 10^7 entries in an empty dictionary

In [7]: `start = time.perf_counter()
d = {}
for i in range(10000000,0,-1):
 d[i] = i
elapsed = time.perf_counter() - start
print(elapsed)`

1.2130631100000073

- Doubling the operations, doubles the time, so linear
- Dictionaries are effectively random access

```
In [8]: start = time.perf_counter()
d = {}
for i in range(20000000,0,-1):
    d[i] = i
elapsed = time.perf_counter() - start
print(elapsed)
```

2.5414540689999967

Implementing a "real" list using dictionaries

```
In [9]: def createlist(): # Equivalent of l = [] is l = createlist()
return({})

def listappend(l,x):
    if l == {}:
        l["value"] = x
        l["next"] = {}
        return

    node = l
    while node["next"] != {}:
        node = node["next"]

    node["next"]["value"] = x
    node["next"]["next"] = {}
    return

def listinsert(l,x):
    if l == {}:
        l["value"] = x
        l["next"] = {}
        return

    newnode = {}
    newnode["value"] = l["value"]
    newnode["next"] = l["next"]
    l["value"] = x
    l["next"] = newnode
    return

def printlist(l):
    print("{",end="")

    if l == {}:
        print("}")
        return
    node = l

    print(node["value"],end="")
    while node["next"] != {}:
        node = node["next"]
        print(", ",node["value"],end="")
    print("}")
    return
```

- Display a small list as nested dictionaries

```
In [10]: start = time.perf_counter()
l = createlist()
for i in range(10):
    listappend(l,i)
elapsed = time.perf_counter() - start
print(elapsed)
print(l)
```

0.0030066140000144514

```
{'value': 0, 'next': {'value': 1, 'next': {'value': 2, 'next': {'value': 3, 'next': {'value': 4, 'next': {'value': 5, 'next': {'value': 6, 'next': {'value': 7, 'next': {'value': 8, 'next': {'value': 9, 'next': {}}}}}}}}}}}
```

- Insert 10^7 elements at the beginning in this implementation of a list

```
In [11]: start = time.perf_counter()
l = createlist()
for i in range(10000000):
    listinsert(l,i)
elapsed = time.perf_counter() - start
print(elapsed)
```

9.297652907999918

- Doubling the work doubles the time, so linear

```
In [12]: start = time.perf_counter()
l = createlist()
for i in range(20000000):
    listinsert(l,i)
elapsed = time.perf_counter() - start
print(elapsed)
```

19.07784092600002

- Append 10^4 elements in this implementation of a list

```
In [15]: start = time.perf_counter()
l = createlist()
for i in range(10000):
    listappend(l,i)
elapsed = time.perf_counter() - start
print(elapsed)
```

3.264379636000058

- Halving the work takes 1/4 of the time, so quadratic

```
In [16]: start = time.perf_counter()
l = createlist()
for i in range(5000):
    listappend(l,i)
elapsed = time.perf_counter() - start
print(elapsed)
```

0.7762869529999534

Defining our own data structures

- We have implemented a "linked" list using dictionaries
- The fundamental functions like `listappend`, `listinsert`, `listdelete` modify the underlying list
- Instead of `mylist = {}`, we wrote `mylist = createlist()`
- To check empty list, use a function `isempty()` rather than `mylist == {}`
- Can we clearly separate the **interface** from the **implementation**
- Define the data structure in a more "modular" way