
Hash Tables

CS206
Sep 22

Map

- A searchable collection of key-value pairs
- Multiple entries with the same key are not allowed
- Also known as dictionary (python), associative array (perl)

Map Implementation

```
public class Map206<K,V> {  
    private ArrayList<Pair<K,V>> underlying = new ArrayList<>();  
  
    private class Pair<L,W> {  
        public L key;  
        public W vl;  
        Pair(L key, W value) {  
            ky=key;  
            vl=value;  
        }  
    }  
  
    public V get(K key) {  
  
    }  
}
```

Map Interface

- <https://docs.oracle.com/javase/7/docs/api/java/util/Map.html>

```
public interface Map206Interface<K, V> {  
    public void put(K key, V val);  
    public V get(K key);  
    public boolean containsKey(K key);  
    public int size();  
    public Set<K> keySet();  
}
```

WordCount using maps

All changes invisible to external users

```
private Map206<String, Integer> counts = new Map206<>();
```

```
void countFile(String filename) {
    try (BufferedReader br = new BufferedReader(new FileReader(filename));) {
        String line;
        while (null != (line = br.readLine())) { // read line and test if there is a line to read
            line = line.toLowerCase().replace(".", "").replace(",", "").replace("?", "").replace("!", "").replace("-", "");
        }

        String[] ss = line.split("\\s+"); // split the line by spaces
        for (String token : ss) {
            if (token.length() > 0) {
                int tokencount = 0;
                if (counts.containsKey(token)) {
                    tokencount=counts.get(token);
                }
                counts.put(token, tokencount+1);
            }
        }
    } catch (FileNotFoundException e) {
        System.err.println("Error in opening the file:" + filename);
        System.exit(1);
    } catch (IOException ioe) {
        System.err.println("Error reading file " + ioe);
        System.exit(1);
    }
}
```

```
public String toString() {
    StringBuffer sb = new StringBuffer();
    for (String key : counts.keySet()) {
        sb.append(key + " " + counts.get(key));
        sb.append("\n");
    }
    sb.append("Distinct words: " + counts.size());
    return sb.toString();
}
```

HashTables

- A hash table is an array of size N
 - associated hash function h that maps a key to integers into $[0, N-1]$
 - item (k, v) is stored at index $h(k)$
- $h(x) = x \% N$ is such a function for integers

Simple Hashtable Implementation

```
public class SimpleHT {
    private String[] backingArray;
    public SimpleHT() {
        backingArray = new String[4];
    }
    private int h(int k) {
        return k%4;
    }
    public void put(Integer key, String value) {
        backingArray[h(key)] = value;
    }
    public String get(Integer key) {
        return backingArray[h(key)];
    }
}
```

HashTable Example

```
public static void main(String[] args) {
    SimpleHT sht = new SimpleHT();
    for (int i=0; i<10; i++) {
        System.out.println("adding item with key=" + i
+ " value=" + String.format("%c", 'a'+i));
        sht.put(i, String.format("%c", 'a'+i));
    }
    for (int i=0; i<10; i++)
        System.out.println("getting key="+i+"
value="+sht.get(i));
}
```

adding item with key=0 value=a

....

getting key=0 value=i

getting key=1 value=j

getting key=2 value=g

getting key=3 value=h

getting key=4 value=i

getting key=5 value=j

getting key=6 value=g

getting key=7 value=h

getting key=8 value=i

getting key=9 value=j

Hash Functions

- The goal is to “disperse” the keys in an appropriately random way
- A hash function is usually specified as the composition of two functions:
- hash code: $\text{key} \rightarrow \text{integers}$
- compression: $\text{integers} \rightarrow [0, N-1]$

see SepChainHT.java

Hash Codes

- Polynomial accumulation: partition bits of key into a sequence of components of fixed length $a_0a_1\dots a_{n-1}$

- Evaluate the polynomial

$$p(z) = a_0 + a_1z + a_2z^2 + \dots + a_{n-1}z^{n-1}$$

Polynomial accumulation on Strings

Recommended
by textbook

```
static int POLY_MULT=33;
public int stringHasher(String ss) {
    BigInteger ll = new BigInteger("0");
    for (int i=0; i<ss.length(); i++) {
        BigInteger bb =
        BigInteger.valueOf(POLY_MULT).pow(i).multiply(BigInteger.valueOf(
        (int)ss.charAt(i)));
        ll = ll.add(bb);
    }
    ll = ll.mod(BigInteger.valueOf(backingArray.length));
    return ll.intValue();
}
```

Handles really
large numbers

Array storing the
hashtable

$33^{15}=59938945498865420543457$

Collisions

drawing 500 unique words from Oliver Twist and assuming a hashtable size of 1009, get these collisions

16 probable child when
42 fagins xxix importance that xv administering
104 stage pledge near
132 surgeon can night
271 things fang birth
341 alone sequel life
415 maylie check circumstances
418 mentioning containing growth
625 meet she first
732 there affording encounters
749 possible out acquainted
761 never xviii after goaded where
833 marks jew gentleman
985 adventures inseparable experience

Realistic hash codes computation in Java

- Use the hashCode function defined on Java Object.
- So put into hashtable is just

```
private int h(Object k) {  
    return k.hashCode() % backingArray.length;  
}  
public void put(Object key, Object value) {  
    backingArray[h(key)] = value;  
}
```

Collisions

- Handling of collisions is one of the most important topics for hashtables
 - Rehashing
 - make the table bigger
 - $O(n)$ time so want to avoid
 - Alternative to rehashing
 - Separate Chaining
 - Probing

Separate Chaining

- Idea: each spot in hashtable holds a array list of key value pairs when the key maps to that hashvalue.
- Replace the item if the key is the same
- Otherwise, add to list
- Generally do not want more than about number of objects as size of table
- Chains can get long

Hash tables get crowded, chains get long

HT_SIZE=1009

Using unique words drawn from “Oliver Twist”.
Unique count at top of table

278

0	762
1	217
2	29
3	1
4	0
5	0
6	0
7	0
8	0
9	0

473

0	622
1	308
2	73
3	5
4	1
5	0
6	0
7	0
8	0
9	0

1550

0	210
1	342
2	252
3	136
4	55
5	9
6	4
7	1
8	0
9	0

2510

0	87
1	198
2	268
3	208
4	140
5	70
6	26
7	10
8	2
9	0

In class exercise

- Show the final contents of the hashtable using separate chaining assuming
 - table size is 7
 - $h(t) = t \% 7$
- Data: $\langle 0,a \rangle$ $\langle 32,b \rangle$ $\langle 39,c \rangle$ $\langle 12,d \rangle$
 $\langle 14,e \rangle$ $\langle 35,f \rangle$ $\langle 27,g \rangle$ $\langle 13,h \rangle$ $\langle 15,i \rangle$
 $\langle 5,j \rangle$ $\langle 12,k \rangle$ $\langle 13,l \rangle$ $\langle 4,m \rangle$ $\langle 0,n \rangle$
 $\langle 35,o \rangle$
- What is the longest chain?

Open Addressing Linear Probing

- Store only $\langle K, V \rangle$ at each location in array
 - No awkward lists
- If key is different and location is in use then go to next spot in array
 - repeat until free location found

Linear Probing Example

- Suppose
 - hashtable size is 7
 - $h(t) = t \% 7$
 - add:
 - $\langle 3, A \rangle$
 - $\langle 10, B \rangle$
 - $\langle 17, C \rangle$
 - $\langle 24, Z \rangle$
 - $\langle 3, D \rangle$
 - $\langle 4, E \rangle$

Linear Probing

- Store only $\langle K, V \rangle$ at each location in array
 - No awkward linked lists
- If key is different and location is in use then go to next spot in array
 - if key is same, replace value
 - repeat until free location found

Probing Distance

- Given a hash value $h(x)$, linear probing generates $h(x)$, $h(x) + 1$, $h(x) + 2$, ...
 - Primary clustering – the bigger the cluster gets, the faster it grows
- Quadratic probing – $h(x)$, $h(x) + 1$, $h(x) + 4$, $h(x) + 9$, ...
 - Quadratic probing leads to secondary clustering, more subtle, not as dramatic, but still systematic
- Double hashing
 - Use a second hash function to determine jumps

Performance Analysis for probing

- In the worst case, searches, insertions and removals take $O(n)$ time
 - when all the keys collide
- The load factor α affects the performance of a hash table
 - expected number of probes for an insertion with open addressing is $\frac{1}{1 - \alpha}$
- Expected time of all operations is $O(1)$ provided α is not close to 1
 - NOTE: cheating here $O()$ is about true worst case

Open Addressing vs Chaining

- Probing is significantly faster in practice
- locality of references – much faster to access a series of elements in an array than to follow the same number of pointers in a linked list
- Efficient probing requires soft/lazy deletions – tombstoning
 - de-tombstoning

In class exercise

- Show the final contents of the hashtable using linear probing assuming
 - table size is 13
 - $h(t) = t \% 13$
- Data: $\langle 0, a \rangle$ $\langle 32, b \rangle$ $\langle 39, c \rangle$ $\langle 12, d \rangle$
 $\langle 14, e \rangle$ $\langle 35, f \rangle$ $\langle 27, g \rangle$ $\langle 13, h \rangle$ $\langle 15, i \rangle$
 $\langle 5, j \rangle$ $\langle 12, k \rangle$ $\langle 13, l \rangle$ $\langle 4, m \rangle$ $\langle 0, n \rangle$ $\langle 35, o \rangle$
- What is the most number of steps you needed to take to find a free location?

Using Hashtables

- No worries about hashing functions, rehashing, ...
- Someone else responsibility
- Example: who is visiting my site, and how often?
 - for instance, hackers?
- web servers keep access logs
- `java.util.HashMap`