

Lecture 19: Hash Code & Implementations

generating ideas for hash tables:

→ how to define $h(x)$?

Terminology: a hashfunction h is a rule for assigning (numerical) values to object instances

— $h(x)$ = hash code of x or hash value

Goal: Fix a rule for generating hash codes where values appear random

— collisions are unlikely.

(x and y collide if $h(x) = h(y)$ but $x \neq y$).

In Java: Object class has a method

hashCode(): for any object, obj , calling $obj.hashCode()$ returns an int that is a hash code for obj .

Guarantees:

(1) if $obj.hashCode()$ involved multiple times, same execution, semantically same object gets same value.

(2) if $o.equals(p)$, then $o.hashCode() = p.hashCode()$ but reverse is not necessarily true

Non-Guarantee:

- "as much as reasonably possible", $hashCode()$ should avoid collisions

Conclusion: Defining hashCode() for a class is an art form.

art form.
→ should NOT assume that default implementation is great... "trust no one".

Two issues: toward implementing hashtable

- ① `hashCode()` → negative value (Index out of bounds)
- ② `hashCode` may have patterns (not random) and lead to poor performance

Goal: Use hashCode() method to generate indices for a hash table of size n.

- assume $n = 2^d$

Ideal: If $a.hashCode() \neq b.hashCode()$

then, $P(h(a) = h(b)) = \frac{1}{n} (= \frac{1}{2^d})$

Simple approach: $x = \text{hashCode}(\text{obj})$ (some int)

determine $h(x) = |x| \% n \leftarrow$ return between 0 and $n-1$
 \nwarrow make positive

Problem (Doesn't satisfy ideal: if $x=0$, $y=n$
 $\Rightarrow h(x)=0$, $h(y)=0$.)

Diff. way of hashing

- choose z at random (int)

define h as $h(x) = (z * x) \% n$

Better or worse than before?

- what if z is divisible by n ?
 \rightarrow everything hashes to 0!

Terrible... but...

A Tangent \Rightarrow int are 32 bit binary rep. of #s:

$\overbrace{001001110\dots01}^{32}$
 \swarrow \searrow
 sign bits 31 value bits
 1 \rightarrow neg # # from 0 to $(2^{31}-1)$
 0 \rightarrow pos #
 size is fixed

eg. $(2^{31}-1): 011111\dots1$
 $+ 1 : 000000\dots1$

"2s complement representation"

$2^{31} : 10000000\dots00$ or (-2^{31})

Recall: $n=2^d$

Fact: If z is a random odd number and x is any odd number, then $(z * x) \% n$ is equally likely to be any value in $0, 1, \dots, n-1$.

More generally: $x = \underbrace{01001\dots0100\dots0}_d \underbrace{\dots}_{32-d}$

then in binary rep of $z * x$ (as int),

$z * x = \underbrace{***\dots*}_{d} \underbrace{00\dots0}_{32-d}$

all 0/1 values of $*$ are equally likely. d $32-d$

Consequence: if $x \neq y$, and z is an odd number chosen at random, then the first d -bits of $z * x$ and $z * y$ are equal with probability of $\left(\frac{1}{2^d} = \frac{1}{n}\right)$ (randomness over choice of z)

Suggests: given $x = \text{obj.hashCode}()$
define $h(x) = \text{first } d \text{ bits of } z * x$.

In code: $\ggg = \text{bit shift operator}$.

$$\begin{cases} y = b_{31} b_{30} b_{29} \dots b_0 & (\text{bits of } y) \\ y \ggg j = \underbrace{00 \dots 0}_j b_{31} b_{30} \dots b_j \end{cases}$$

first d bits of y is $y \ggg (32-d)$