# CS6033 Lecture 3
# Slides/Notes

**Review of Math Background in Probability;
Hash Tables  (Notes, Ch 11 (skip Secs. 11.3.4,
11.3.5 and 11.5))**


By Prof. Yi-Jen Chiang

CSE Dept., Tandon School of Engineering

New York University

1

---

Abstract Data Type (ADT): Dictionary D
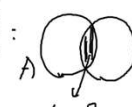
Support 3 operations:

1. Insert (X, D)    (x: item with a key)

2. Search (k, D)    (k: key)

3. Delete (K, D) : Delete the item with key K
                   from D.

Hash Table : Randomized Algorithms
             (Using probabilistic analysis)

Review of Math Background
in Probability

1. Boole's Inequality
   (Union Bound)

① Let A, B be 2 events.
   $Pr[A \cup B] \leq Pr[A] + Pr[B]$

Pf :         B    $Pr[A \cup B]$
   A              $= Pr[A] + Pr[B] - Pr[A \cap B]$
       $A \cap B$     $\leq Pr[A] + Pr[B]$.

② Let $A_1, A_2 \ldots A_n$ be $n$ events.
   $Pr[A_1 \cup \cdots \cup A_n] \leq \sum_{i=1}^{n} Pr[A_i]$

2

## 2. Linearity of Expectation

Def: Let $X$ be a random variable. Expected value of $X$, $E[X]$, is defined as

$$E[X] = \sum_x x \cdot Pr[X=x]$$

Intuition: weighted sum of value where the weights are probabilities.

eg. Exam score: 95, 92, 80, 75, 71

Average: $(95 + 92 + 80 + 75 + 71)/5$
$= \frac{1}{5} \cdot 95 + \frac{1}{5} \cdot 92 + \cdots$

eg. 95, 95, 80, 75, 75.

$Pr = \frac{2}{5} \quad \frac{1}{5} \quad \frac{2}{5}$

Average: $95 \cdot \frac{2}{5} + 80 \cdot \frac{1}{5} + 75 \cdot \frac{2}{5}$

3

---

## Linearity of Expectation

① Let $X, Y$ be 2 random variables — NOT necessarily independent (i.e. they can be dependent.)

$$E[X+Y] = E[X] + E[Y]$$

Pf: $E[X+Y] = \sum_{x,y} (x+y) \cdot Pr[X=x \text{ and } Y=y]$

$= \underbrace{\sum_{x,y} x \cdot Pr[X=x \text{ and } Y=y]}_{A} + \underbrace{\sum_{x,y} y \cdot Pr[X=x \text{ and } Y=y]}_{B}$   // $E[Y]$ by the same process.
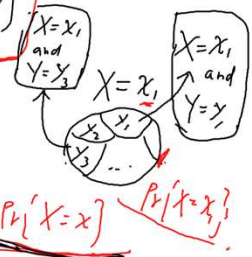
$\sum_x \sum_y x \cdot Pr[X=x \text{ and } Y=y] = \sum_x x \left( \sum_y Pr[X=x \text{ and } Y=y] \right)$

$= \sum_x x \cdot Pr[X=x]$

$= E[X]$.

$X=x_1$ and $Y=y_3$
$X=x_1$ and $Y=y_1$
$X=x_1$

$= Pr[X=x]$
$Pr[X=x_1]$

$A = E[X], \quad B = E[Y]$.

$E[X+Y] = A + B = E[X] + E[Y]$ ✗✗

4

② Let $X_1, X_2 \cdots X_n$ be $n$ random variables. (they can be dependent)

$E[X_1 + X_2 + \cdots + X_n] = E[X_1] + E[X_2] + \cdots + E[X_n]$.

Pf: $E[X_1 + X_2 + \cdots + X_n] = E((X_1 + \cdots + X_{n-1}) + X_n) = E[X_1 + \cdots + X_{n-1}] + E[X_n]$

$= E[X_1 + \cdots + X_{n-2}] + E(X_{n-1}) + E[X_n] = \cdots = E[X_1] + E[X_2] + \cdots + E[X_n]$.

3. <u>Bernoulli Trial</u> : Flip a coin (with $Pr['Head] = p$ and $Pr['Tail] = q = 1-p$) many times, each time is indept. of the others.

① Expected # of flips to get the <u>first head</u> $= \frac{1}{p}$. $\left( eg. \ p = \frac{1}{2} \Rightarrow \frac{1}{p} = 2 \right)$

---

Pf : <u>M1</u> : Let $X$ be random variable for the # flips to get the first head.

$E[X] = \sum_{i=1}^{\infty} i \cdot Pr['X=i]$   $\left[ \begin{array}{l} Pr\{X=i\} = Pr\{ \text{first } (i-1) \text{ flips are tails \&} \\ \qquad\qquad\quad \text{the last flip is head} \} \\ \qquad\quad = q^{i-1} p. \end{array} \right]$

$= \sum_{i=1}^{\infty} i \cdot (q^{i-1} \cdot p) = S$

$S = 1 \cdot p + 2 \cdot q p + 3 \cdot q^2 p + 4 q^3 p + \cdots$

$-) \ qS = \qquad\quad 1 \cdot qp + 2 \cdot q^2 p + 3 q^3 p + \cdots$

$(1-q)S = 1 \cdot p + 1 \cdot qp + 1 \cdot q^2 p + 1 \cdot q^3 p + \cdots$

$= p(1 + q + q^2 + q^3 + \cdots) = p \cdot \frac{1}{1-q} = p \cdot \frac{1}{p} = 1$

$\boxed{(1-q) S = 1}$

$\Rightarrow S = \frac{1}{1-q} = \frac{1}{p}.$

$S = E[X] = \frac{1}{p}.$ ※

<u>M2</u>: Let $E$ be the desired expected # of flips.

* <u>MA</u>:   $\underline{E = 1 + 0 \cdot p + E \cdot q} \leftrightarrow$ *If first flip is Head : $Pr = p$. 0 additional flips

first slip is Head : $Pr = p$. 0 additional flips

Tail : $Pr = q$. $E = $ :

$\boxed{(1-q) E = 1} \Rightarrow E = \frac{1}{1-q} = \frac{1}{p}$ ※

*MB: $E = P \cdot (1 + 0) + q(1 + E)$   # total flips is $\begin{cases}(1+0) & \text{if first flip is Head.} \\ \underbrace{\quad}_{P_r = P.} \\ (1+E) & \text{if first flip is Tail} \\ \underbrace{\quad}_{P_r = q.}\end{cases}$

$E = P + q + qE$

$(P + q = 1)$   $= 1 + qE \Rightarrow E = \frac{1}{1-q} = \frac{1}{P}$ ※

② Expected # of flips to get the

⚡ K-th head: $X$

$\begin{array}{c} A = 1 + q + q^2 + \cdots + \quad \\ -) qA = \quad q + q^2 + \cdots \quad \\ \hline (1-q)A = 1 \quad \Rightarrow A = \frac{1}{1-q}. \end{array}$

1st Head | 2nd Head | -- | k-th Head.

OOOO~◉ | O·--·◉ | --- | O·O·O~● |
T T T H | T·-- H | | T T~ H |

$X_1$ | $X_2$ | | $X_K$.

Def: $X_i$: # flips to get the $i$-th Head after getting the $(i-1)$st Head.

$X = X_1 + X_2 + \cdots + X_K$     $E[X] = E[X_1] + E[X_2] + \cdots + E[X_k] = \frac{1}{P} \cdot K = \frac{K}{P}$ ※

---

Hash Table     Array: keys are in $\mathbb{Z}^+$.        keys as indices
$k$

keys are from $U$, where $U$ is a large set of possible key values.

$n$ items. If we use an array: $\boxed{\begin{array}{c} O(1) \\ \text{worst-case} \\ \text{time} \end{array}} \rightarrow$ for insert, delete, search (Good)

But we need memory space $|U| \gg n$.

Bad!

Hash Table: Use memory space $m = O(n)$
and retain $O(1)$ time for insert, delete, search.

→ Typically expected time   [Worst-case time can be $\Theta(n)$]

[Look at details for each specific approach.]

Main Idea : Use a hash table of $m$ slots

slot indices are $0, 1, \cdots, (m-1)$.

Use a hash function $h : U \longrightarrow \{0, 1, 2, \cdots, (m-1)\}$

For each key $k \in U$. $h(k) \in \{0, 1, \cdots, (m-1)\}$

$h(k)$ ... key $(k)$

Issue : Different keys may be mapped/hashed to the same slot.

Called Collision

2 ways to handle collisions : $\underline{M1}$ : chaining

$\underline{M2}$ : Open Addressing

(next class)

$h(k_1)$

$h(k_2) \to$ ... $k_2$ ... $\times \to k_1$

Insert : $O(1)$ time worst-case

Always add to the front of the list/chain

9

Ideal Situation / Independent uniform hashing

Strong Assumption : (Simple uniform hashing in 3rd Ed.)

Assume : For each key $k$. $h(k)$ goes to one of the $m$ slots with equal prob $\frac{1}{m}$   i.e. $Pr\left[h(k) = i\right] = \frac{1}{m}$, $\forall k \in U$ and $i \in \{0, 1, \cdots, (m-1)\}$

$h(k)$ ... $0$ $1$ $2$ ... $m-1$

[meaning: the hash function $h()$ is perfect]

(1) Expected search time for an unsuccessful search

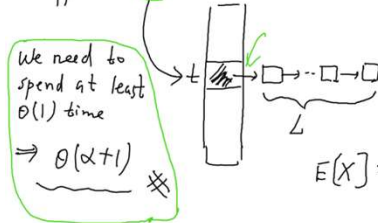(2) Expected search time for a successful search

10

5

Under independent uniform hashing assumption.

Thm A : Expected time for an unsuccessful search is
$\Theta(\alpha+1)$.

Pf : Let $k$ be the key of the unsuccessful search.
(k is NOT any of the n items in the hash table)

Suppose $h(k)=t$

Recall:
n : # items
m : # slots in Hash Table.
Def : Load factor
$\alpha = \frac{n}{m}$.
Assume : $\alpha$ is some constant.

Let $X$ be a random var. for the search time.

$X = X_1 + X_2 + \cdots + X_n$

where $X_i$ is a random var : $X_i = \begin{cases} 1 & \text{if } h(i)=h(k) \\ 0 & \text{if } h(i) \neq h(k) \end{cases}$

$Pr\{X_i = 1\} = \frac{1}{m}$.

We need to spend at least $\Theta(1)$ time
$\Rightarrow \Theta(\alpha+1)$ #

$E[X] = E[X_1] + E[X_2] + \cdots + E[X_n] = \left(1 \cdot \frac{1}{m} + 0 \cdot (1-\frac{1}{m})\right) \cdot n = \frac{n}{m} = \alpha$
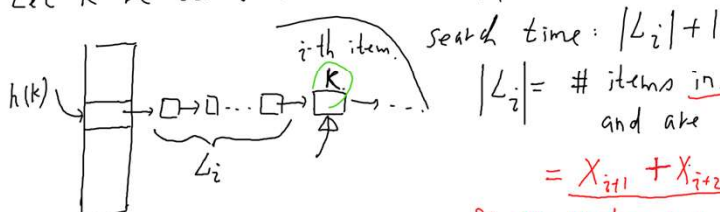
11

---

Thm B : The expected time for a successful search is $\Theta(\alpha+1)$.

Pf : The key $k$ is one of the n items inserted before.
(searched) and each of these n items has an equal prob of $\frac{1}{n}$
to be the searched key $k$.

Let $k$ be the $i$-th inserted item. (we will range $i$ from $1, 2, \cdots, n$ with equal prob $\frac{1}{n}$)

search time : $|L_i| + 1$

$|L_i| = $ # items inserted after the $i$-th item and are hashed into $h(k)$.

$= X_{i+1} + X_{i+2} + \cdots + X_n$

$X_j$ is random var for the $j$-th inserted item $\ell$

$X_j = \begin{cases} 1 & \text{if } h(\ell) = h(k) \\ 0 & \text{else} \end{cases}$    $Pr\{X_j = 1\} = \frac{1}{m}$    $(\ell \neq k)$

$k$ is the $i$-th inserted item.

Let $X$ be random var for search time.   $X = |L_i| + 1$

$E[X] = \frac{1}{n}\sum_{i=1}^{n} E[|L_i|+1] = \frac{1}{n}\sum_{i=1}^{n}\left(1 + E[X_{i+1}] + E[X_{i+2}] \cdots + E[X_n]\right) = \frac{1}{n}\sum_{i=1}^{n}\left(1 + \frac{1}{m}(n-i)\right]$

12

$$E[X] = \frac{1}{n}\sum_{i=1}^{n}\left(1+\frac{n-i}{m}\right) = \frac{1}{n}\left[n + \sum_{i=1}^{n}\frac{n-i}{m}\right]$$

$$= 1 + \frac{1}{n}\sum_{i=1}^{n}\frac{1}{m}(n-i) = 1 + \frac{1}{mn}\sum_{i=1}^{n}(n-i) = 1 + \frac{1}{mn}\left(0+1+2+\cdots+(n-1)\right)$$

$$= 1 + \frac{1}{mn}\cdot\frac{(n-1)\cdot n}{2} = 1 + \frac{1}{2}\frac{n}{m} - \frac{1}{2m} = \Theta\left(1+\frac{1}{2}\alpha\right) = \Theta(1+\alpha) \quad \#$$

13

## Hash Functions

* Static Hashing:

1. Division Method : $h(k) = k \bmod m$. ↗ Typicall we use a prime number for $m$.

2. Multiplication Method : Let $A$ be a real #. $A \in (0,1)$. $h(k) = \lfloor (A \cdot k \bmod 1)\cdot m \rfloor$

   ($A \cdot k \bmod 1$: $A \cdot k$. take the decimal part)

2'. Multiply-shift Method : (special case of multiplication method. See textbook)

* Random Hashing: We have a family $\mathcal{H}$ of hash functions. For each execution (a sequence of insert, search, delete ops), we randomly select one hash function from $\mathcal{H}$ to use.

★ Universal Hashing (most important one. See next) | Others: see textbook.

14

# Slide 15

✳ **Universal Hashing:** Use a family $H$ of hash functions.

For each execution, randomly choose one hash function $h()$ from $H$ (with equal prob. $\frac{1}{|H|}$) to use. (size $|H|$)

**Property of univ. hashing:** For any pair of keys $k \neq \ell$, there are at most $\boxed{t}$ hash functions $h()$ in $H$ s.t. $h(k) = h(\ell)$ i.e. $h()$ hashes $k, \ell$ into collision.

**Goal:** $\boxed{\forall \text{keys } k \neq \ell, \; Pr\{h(k) = h(\ell)\} = \boxed{\frac{1}{m}}}$ (i.e. $Pr\{\text{hashing into collision}\} = \frac{1}{m}$)

$Pr\{\text{chosen hash function is one of the } t \text{ functions to hash into collision}\}$

$= \frac{t}{|H|} \implies$ We want: $\frac{t}{|H|} = \boxed{\frac{1}{m}} \implies$ take $t = \frac{|H|}{m}$. ✳

✗ Using $\boxed{\forall \text{keys } k \neq \ell, \; Pr\{h(k) = h(\ell)\} = \frac{1}{m}}$ the pfs for Thm A & Thm B carry over. //

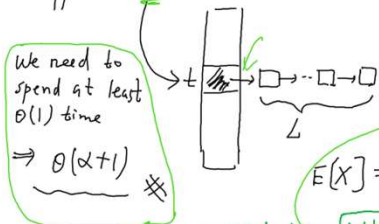No need for the assumption of independent uniform hashing !!

15

# Slide 16

Under independent uniform hashing assumption.

**Thm A:** Expected time for an unsuccessful search is $\Theta(\alpha + 1)$.

**Pf:** Let $k$ be the key of the unsuccessful search.
(k is NOT any of the $n$ items in the hash table)

Suppose $h(k) = t$

Let $X$ be a random var. for the search time.

$X = X_1 + X_2 + \cdots + X_n$

where $X_i$ is a random var: $X_i = \begin{cases} 1 & \text{if } h(i) = h(k) \\ 0 & \text{if } h(i) \neq h(k) \end{cases}$

$Pr\{X_i = 1\} = \frac{1}{m}$

$E[X] = E[X_1] + E[X_2] + \cdots + E[X_n] = \left(1 \cdot \frac{1}{m} + 0 \cdot (1 - \frac{1}{m})\right) \cdot n = \frac{n}{m} = \alpha$

**We need to spend at least $\Theta(1)$ time**
$\implies \Theta(\alpha + 1)$ ✳

equivalent to $\boxed{\forall \text{keys } k, \ell, \; k \neq \ell \quad Pr\{h(k) = h(\ell)\} = \frac{1}{m}}$ (Property obtained from universal hashing)

**Recall:**
$n$: # items
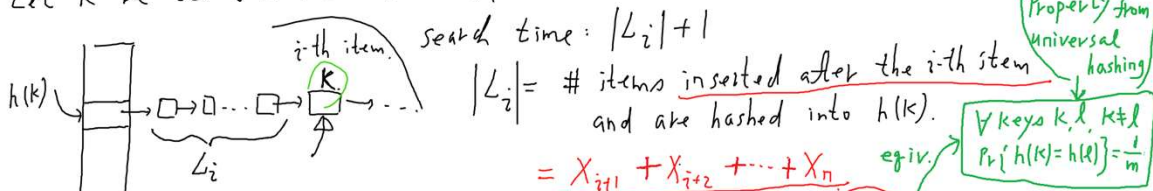$m$: # slots in Hash Table.
Ref: load factor
$\alpha = \frac{n}{m}$.
Assume: $\alpha$ is some constant.

16

8

Thm B: The expected time for a successful search is $\Theta(\alpha+1)$

Pf: The key $K$ is one of the $n$ items inserted before.
(searched) and each of these $n$ items has an equal prob. $\frac{1}{n}$
to be the searched key $k$.

Let $k$ be the $i$-th inserted item. (we will range $i$ from $1, 2, \cdots, n$ with equal prob $\frac{1}{n}$)

search time: $|L_i|+1$

$|L_i| = $ # items inserted after the $i$-th item and are hashed into $h(k)$.

$= X_{i+1} + X_{i+2} + \cdots + X_n$ equiv.

$X_j$ is random var for the $j$-th inserted item $\ell$

$X_j = \begin{cases} 1 & \text{if } h(\ell) = h(k) \\ 0 & \text{else.} \end{cases}$  $Pr\{X_j = 1\} = \left(\frac{1}{m}\right)$  $(\ell \neq K)$

$K$ is the $i$-th inserted item.

Let $X$ be random. var for search time. $X = |L_i|+1$

$E[X] = \frac{1}{n} \sum_{i=1}^{n} E[|L_i|+1] = \frac{1}{n} \sum_{i=1}^{n} \left(1+ E[X_{i+1}] + E[X_{i+2}]\cdots + E[X_n]\right) = \frac{1}{n} \sum_{i=1}^{n} \left(1+ \frac{1}{m}(n-i)\right) = \cdots$

(The rest is the same)

(Property from universal hashing)

$\forall$ keys $k, \ell, k \neq \ell$
$Pr\{h(k) = h(\ell)\} = \frac{1}{m}$

17

9