# CS7800: Advanced Algorithms
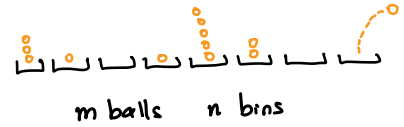
Class 22: Randomized Algorithms III
- Balls and Bins: Chernoff Bounds
- Universal Hashing

Jonathan Ullman

November 25, 2025
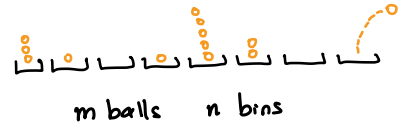
# Balls and Bins: Maximum Load


m balls    n bins

- Let $L_i$ be the number of balls in bin $i$

- Expected maximum load is $\mathbb{E}\left(\max_i L_i\right) = \sum_{k=1}^{\infty} \mathbb{P}\left(\max_i L_i \geq k\right)$

  want to bound this probability

So far:

① Trivial bound: $\mathbb{E}(\max L_i) \leq m$

② Markov's inequality: $\mathbb{E}(\max L_i) \leq \infty$

③ Chebyshev's inequality: $\mathbb{E}(\max L_i) \leq O\left(\frac{m}{n} + \sqrt{m}\right)$

Today: Tighter analysis with Chernoff Bounds

# Balls and Bins: Maximum Load

- Let $L_i$ be the number of balls in bin $i$

- Expected maximum load is $\mathbb{E}\left(\max\limits_i L_i\right) = \sum\limits_{k=1}^{\infty} \mathbb{P}\left(\max\limits_i L_i \geq k\right)$

- Let $L_{i,j} = \begin{cases} 1 & \text{if ball } j \text{ is in bin } i \\ 0 & \text{otherwise} \end{cases}$  $\Rightarrow$  $L_i = L_{i,1} + \ldots + L_{i,m}$

m balls    n bins

# Aside: Central Limit Theorem

# Chernoff Bounds

$$Z_i = \begin{cases} 1 \text{ with prob } p \\ 0 \text{ with prob } 1-p \end{cases}$$

$$Z = Z_1 + \dots + Z_m$$

$Z_1, \dots, Z_m$ independent

$$\mu = \mathbb{E}(Z) = pm$$

Thm:
$$\mathbb{P}\left(Z \geq (1+\varepsilon)\mu\right) \leq \left(\frac{e^\varepsilon}{(1+\varepsilon)^{1+\varepsilon}}\right)^\mu$$

$\varepsilon < 1$

$$e^{-\frac{\mu \varepsilon^2}{4}}$$

$\varepsilon > 1$

$$e^{-\mu\varepsilon}$$

# Chernoff Bounds

**Thm:** $\mathbb{P}\left( Z \geq (1+\varepsilon)\mu \right) \leq \left( \dfrac{e^{\varepsilon}}{(1+\varepsilon)^{1+\varepsilon}} \right)^{\mu}$

$Z_i = \begin{cases} 1 & \text{with prob } p \\ 0 & \text{with prob } 1-p \end{cases}$

$Z_1, \ldots, Z_m$ independent

$Z = Z_1 + \ldots + Z_m$

$\mu = \mathbb{E}(Z) = pm$

**Proof:**

$\mathbb{P}(Z \geq t\mu) = \mathbb{P}\left( e^{sZ} \geq e^{st\mu} \right)$ ← What is this dark magic?

$\leq e^{-st\mu} \cdot \mathbb{E}\left( e^{sZ} \right)$ ← Markov

$= e^{-st\mu} \cdot \mathbb{E}\left( \prod_i e^{sZ_i} \right)$

$= e^{-st\mu} \cdot \prod_i \mathbb{E}\left( e^{sZ_i} \right)$ ← Independence

$Z_i \in \{0,1\}$ so this is something "simple"

# Chernoff Bounds

Thm: $\mathbb{P}\left( Z \geqslant (1+\varepsilon)\mu \right) \leqslant \left( \dfrac{e^{\varepsilon}}{(1+\varepsilon)^{1+\varepsilon}} \right)^{\mu}$

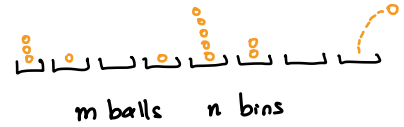$Z_i = \begin{cases} 1 & \text{with prob } p \\ 0 & \text{with prob } 1-p \end{cases}$

$Z_1, \ldots, Z_m$ independent

$Z = Z_1 + \ldots + Z_m$

$\mu = \mathbb{E}(Z) = pm$

## Proof Cont'd:

# Balls and Bins: Maximum Load

- Let $L_i$ be the number of balls in bin $i$

- Expected maximum load is $\mathbb{E}\left(\max_i L_i\right) = \sum_{k=1}^{\infty} \mathbb{P}\left(\max_i L_i \geq k\right)$

- Let $L_{i,j} = \begin{cases} 1 & \text{if ball } j \text{ is in bin } i \\ 0 & \text{otherwise} \end{cases}$   $\Rightarrow$   $L_i = L_{i,1} + \ldots + L_{i,m}$

Apply Chernoff Bound

m balls    n bins

# Application: Hash Tables

Goal:   Store a set of $\underline{m\ elements}$  $S \subseteq U$,
such that we can efficiently check if $x \in S$

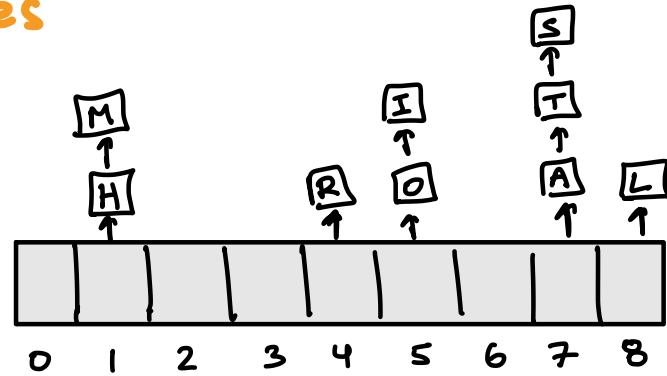$\hookrightarrow$ A "dictionary" also lets us associate a value
with each key $x$

- A $\underline{hash\ table}$  $T[1:n]$ stores the elements in $n$ $\underline{bins}$

- A $\underline{hash\ function}$  $h: U \to \{0, 1, \ldots, n-1\}$ maps elements
to bins  $x \mapsto T[h(x)]$

# Application: Hash Tables

Linear chaining:
a common way to
deal with collisions
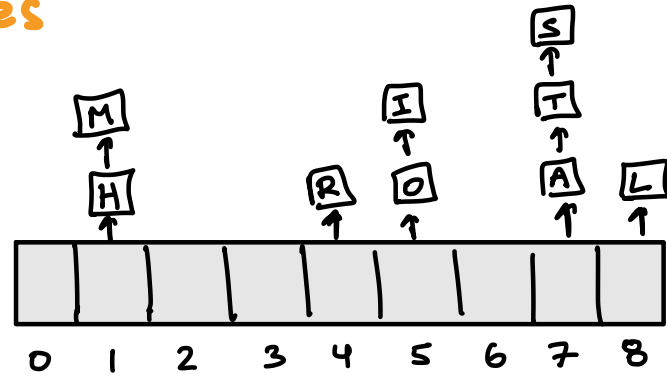


Looks a lot like balls in bins!

- Load factor $= \dfrac{m}{n}$

- Let $\ell(x) =$ #of elements
  in the same bin as $x$
  $\#\{y \in S : h(y) = h(x)\}$

"collisions"

Worst-case lookup time $= \max\limits_{x \in \mathcal{U}} \ell(x)$

- Time to lookup $x \in \mathcal{U}$ is $O(\ell(x))$

# Application: Hash Tables

How should we choose the hash function $h: U \to \{0, 1, \ldots, n-1\}$ to have small maximum load?



Looks a lot like balls in bins!

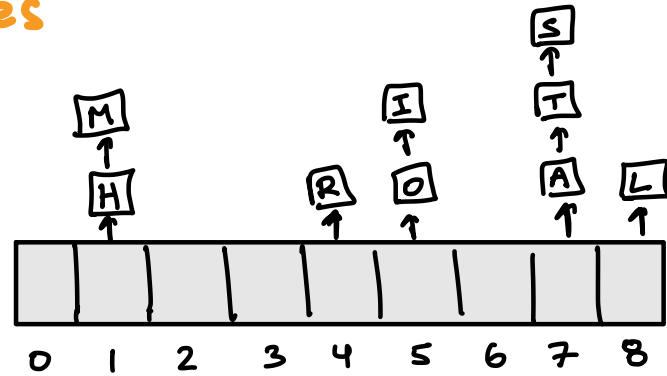Randomized hash function:

- Model $h$ as a uniformly random function $U \to \{0, 1, \ldots, n-1\}$

- Fix the set $S$ and study $\underset{\sim}{\mathbb{E}} \left( \max_x \ell(x) \right)$

  ← expectation over random choice of $h$

# Application: Hash Tables

How should we choose
the hash function
$h: \mathcal{U} \to \{0, 1, \ldots, n-1\}$ to have
small maximum load?



Looks a lot like balls in bins!

Uniformly random hash functions: $\left(\text{load factor } \frac{m}{n} = 1\right)$

 - Expected max load is $\Theta\left(\frac{\log n}{\log\log n}\right)$ [Will be true 99.999% of the time]

 - For any $x \in \mathcal{U}$, expected lookup time is
 $$\mathbb{E}\left(\ell(x)\right)$$

# Universal Hash Families

A **hash family** is a set of hash functions

$$\mathcal{H} \subseteq \{ h: U \longrightarrow \{0, 1, \ldots, n-1\} \}$$

<u>Definition</u>: $\mathcal{H}$ is 2-universal if for every distinct $x, y \in U$

$$\mathbb{P}_h \left( h(x) = h(y) \right) \leq \frac{1}{n}$$

Behaves like a uniformly
random hash fn if we only
look at pairs of points

# Constructing Universal Hashing

Construction:

- Fix a prime $p \geq |\mathcal{U}|$, bins $n$

- Let $h_{a,b}(x) \equiv (ax + b \bmod p) \bmod n$

$$\mathcal{H}_{p,n} = \left\{ h_{a,b} : a \in \mathbb{Z}_p^{\neq 0}, b \in \mathbb{Z}_p \right\}$$

Theorem: $\mathcal{H}_{p,n}$ is a 2-universal hash family

# Constructing Universal Hashing

$\mathcal{H}$ is 2-universal if for every distinct $x, y \in \mathcal{U}$

$$\underset{h}{\mathbb{P}}\left( h(x) = h(y) \right) \leq \frac{1}{n}$$

$h_{a,b}(x) = \left( ax + b \mod p \right) \mod n$

$a, b$ random
$a \neq 0$

$p$ fixed

**Lemma 1:** For every prime $p$

and $a \neq 0$ there is a unique

$a^{-1} \in \{1, 2, \ldots, p-1\}$ such that $a^{-1} \cdot a = 1 \mod p$

**Proof:** ① $az = c \mod p$ has at most one solution

if $az = az' \mod p$ $\implies$ $a(z - z') = 0 \mod p$

for $z, z' \in \{1, 2, \ldots, p-1\}$ $\implies$ $z - z'$ divisible by $p$

$\implies$ $z - z' = 0$

② $az = 0 \mod p$ has no solutions

# Constructing Universal Hashing

$$h_{a,b}(x) = \left(ax + b \mod p\right) \mod n$$

$a, b$ random
$a \neq 0$

$p$ fixed

**Lemma 2:** If $x \neq y$ and $r \neq s$ then the system

$$ax + b = r \mod p$$
$$ay + b = s \mod p$$

has a unique solution

**Proof:**

# Constructing Universal Hashing

$H$ is 2-universal if for every distinct $x, y \in U$
$$\underset{n}{P}(h(x) = h(y)) \leq \frac{1}{n}$$

$$h_{a,b}(x) = (ax+b \mod p) \mod n$$

$a, b$ random
$a \neq 0$

$p$ fixed

**Thm:** $H_{p,n}$ is 2-universal

**Proof:** By Lemma 2,

$$\underset{a,b}{P}\left(ax+b = r \mod p \wedge ay+b = s \mod p\right) = \frac{1}{p(p-1)}$$

$$\underset{a,b}{P}\left(h_{a,b}(x) = h_{a,b}(y)\right) = \frac{N}{p(p-1)}$$

where $N$ is number of $r \neq s$ in $\mathbb{Z}_p$ so that $r = s \mod n$

$$\leq \frac{p(p-1)}{m \, p(p-1)}$$

$$N \leq p \cdot \frac{p-1}{n}$$

choices of $r$

choices of $s$ for a given $r$

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

# Application: Hash Tables

How should we choose
the hash function
$h: U \to \{0, 1, \ldots, n-1\}$ to have
small maximum load?



Looks a lot like balls in bins!

~~Uniformly random~~ Universal hash functions: (load factor $\frac{m}{n} = 1$)

- Expected max load is $\Theta\left(\frac{\log n}{\log \log n}\right)$ $\Theta(\sqrt{n})$

- For any $x \in U$, expected lookup time is
$$\mathbb{E}(\ell(x)) = \sum_{y \in S} \mathbb{P}(h(x) = h(y)) \leq n \cdot \frac{1}{n} = 1$$

Still true for universal hash families!