

CS7800: Advanced Algorithms

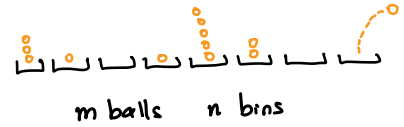
Class 22 : Randomized Algorithms III

- Balls and Bins: Chernoff Bounds
- Universal Hashing

Jonathan Ullman

November 25, 2025

Balls and Bins: Maximum Load



- Let L_i be the number of balls in bin i
- Expected maximum load is $\mathbb{E}(\max_i L_i) = \sum_{k=1}^{\infty} \underbrace{\mathbb{P}(\max_i L_i \geq k)}_{\text{want to bound this probability}}$

So far:

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

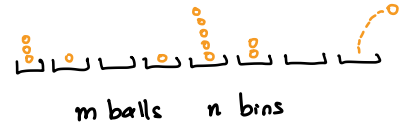
② Markov's inequality: $\mathbb{E}(\max L_i) \leq \infty$

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

Today: Tighter analysis with Chernoff Bounds $\mathbb{E}(\max L_i) = O\left(\frac{m}{n} + \frac{\log n}{\log \log n}\right)$

Union
bound

Balls and Bins: Maximum Load



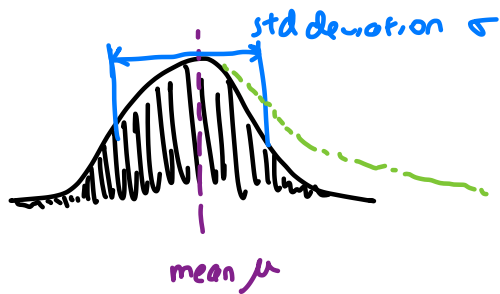
- Let L_i be the number of balls in bin i
- Expected maximum load is $\mathbb{E}(\max_i L_i) = \sum_{k=1}^{\infty} \mathbb{P}(\max_i L_i \geq k)$
- Let $L_{i,j} = \begin{cases} 1 & \text{if ball } j \text{ is in bin } i \\ 0 & \text{otherwise} \end{cases} \Rightarrow L_i = L_{i,1} + \dots + L_{i,m}$

Want to bound $\mathbb{P}(\max_i L_i \geq k) \leq n \cdot \mathbb{P}(L_i \geq k)$

L_i is a sum of m independent
simple random variables

Aside: Central Limit Theorem

Gaussian distribution Z



$$P(Z \geq \mu + 3\sigma) \leq .003$$

$$P(Z \geq \mu + 6\sigma) \leq .000001 ?$$

$$P(Z \geq \mu + t\sigma) \leq e^{-t^2}$$

$$\sigma^2 = \text{Var}(Z)$$

$$\text{Chebyshev } P(Z \geq \mu + t\sigma) \leq 1/t^2$$

If Z_1, \dots, Z_m are independent random variables with $\mu = E(Z_i)$, $\sigma^2 = \text{Var}(Z_i)$ then $Z = \frac{Z_1 + \dots + Z_m}{\sqrt{m}}$ then

$Z \xrightarrow{m \rightarrow \infty}$ Gaussian with mean μ and variance σ^2

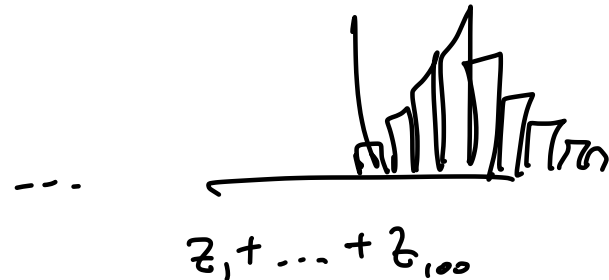
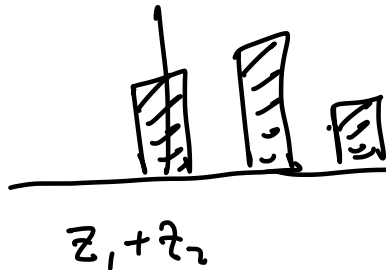
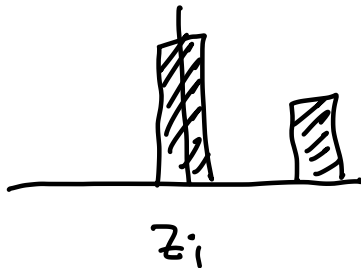
Aside: Central Limit Theorem

CLT:

If z_1, \dots, z_m are independent random variables with $\mu = \mathbb{E}(z_i)$, $\sigma^2 = \text{Var}(z_i)$ then $Z = \frac{(z_1 - \mu) + (z_2 - \mu) + \dots + (z_m - \mu)}{\sqrt{z_m}}$

$Z \xrightarrow{m \rightarrow \infty}$ Gaussian with mean 0 and variance 1

e.g. $z_i = \begin{cases} 1 & \text{up } 1/3 \\ 0 & \text{up } 2/3 \end{cases}$



Chernoff Bounds

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

Z_1, \dots, Z_m independent

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

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

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

$\mathbb{P}(Z - \mu \geq \epsilon\mu)$

$\epsilon < 1$

$-\frac{\mu\epsilon^2}{4}$

e

$\epsilon > 1$

$-\mu\epsilon$

$\frac{e^\epsilon}{(1+\epsilon)^{1+\epsilon}} \leq \left(\frac{e}{\epsilon} \right)^{\epsilon\mu}$

Chernoff Bounds

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

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

Z_1, \dots, Z_n independent

$$Z = Z_1 + \dots + Z_n$$

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

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

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

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

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

$$\begin{aligned} \mathbb{E}(e^{sZ_i}) &= pe^s + (1-p) \\ &= 1 + p(e^s - 1) \end{aligned}$$

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

Chernoff Bounds

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

Proof Cont'd:

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

Z_1, \dots, Z_m independent

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

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

Balls and Bins: Maximum Load



- Let L_i be the number of balls in bin i
- Expected maximum load is $\mathbb{E}(\max_i L_i) = \sum_{k=1}^{\infty} \mathbb{P}(\max_i L_i \geq k)$
- Let $L_{i,j} = \begin{cases} 1 & \text{if ball } j \text{ is in bin } i \\ 0 & \text{otherwise} \end{cases} \Rightarrow L_i = L_{i,1} + \dots + L_{i,m}$

Apply Chernoff Bound

$$L_{i,j} = \begin{cases} 1 & \text{w.p. } 1/n \\ 0 & \text{w.p. } 1 - 1/n \end{cases} \quad L_i = L_{i,1} + \dots + L_{i,m} \quad \mathbb{E}(L_i) = \frac{m}{n}$$

Chernoff: $\mathbb{P}(L_i \geq (1+\epsilon)\frac{m}{n}) \leq \left(\frac{e^\epsilon}{(1+\epsilon)^{1+\epsilon}}\right)^{m/n}$

$\frac{m}{n}$ "big" $\frac{m}{n} \geq \log n$

$$\mathbb{P}(L_i \geq (1+\epsilon)\frac{m}{n}) \leq \frac{e^\epsilon}{(1+\epsilon)^{1+\epsilon}} \leq e^{-\frac{\epsilon^2 m/n}{4}} \leadsto \mathbb{E}(\max_i L_i) = O\left(\frac{m}{n}\right)$$

Balls and Bins: Maximum Load



- Let L_i be the number of balls in bin i
- Expected maximum load is $\mathbb{E}(\max_i L_i) = \sum_{k=1}^{\infty} \mathbb{P}(\max_i L_i \geq k)$
- Let $L_{i,j} = \begin{cases} 1 & \text{if ball } j \text{ is in bin } i \\ 0 & \text{otherwise} \end{cases} \Rightarrow L_i = L_{i,1} + \dots + L_{i,m}$

Apply Chernoff Bound

$$L_{i,j} = \begin{cases} 1 & \text{w.p. } 1/n \\ 0 & \text{w.p. } 1 - 1/n \end{cases} \quad L_i = L_{i,1} + \dots + L_{i,m} \quad \mathbb{E}(L_i) = \frac{m}{n}$$

$$\text{Chernoff: } \mathbb{P}(L_i \geq (1+\epsilon) \frac{m}{n}) \leq \left(\frac{e^\epsilon}{(1+\epsilon)^{1+\epsilon}} \right)^{m/n}$$

$$m/n = 1 \quad \left(\frac{m}{n} \text{ "small"} \right)$$

$$\mathbb{P}(L_i \geq 1+\epsilon) \leq \frac{e^\epsilon}{(1+\epsilon)^{1+\epsilon}} \leq \left(\frac{e}{\epsilon} \right)^\epsilon \rightsquigarrow \mathbb{E}(\max_i L_i) = O\left(\frac{\log n}{\log \log n} \right)$$

Application: Hash Tables

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

"universe"

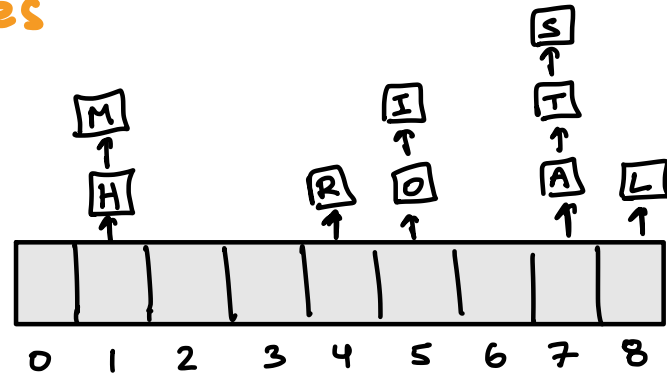


↳ A "dictionary" also lets us associate a value
with each key x

- A hash table $T[1:n]$ stores the elements in n bins
- A hash function $h: \mathcal{U} \rightarrow \{0, 1, \dots, n-1\}$ maps elements
to bins $x \mapsto T[h(x)]$

Application: Hash Tables

Linear chaining:
a common way to
deal with collisions



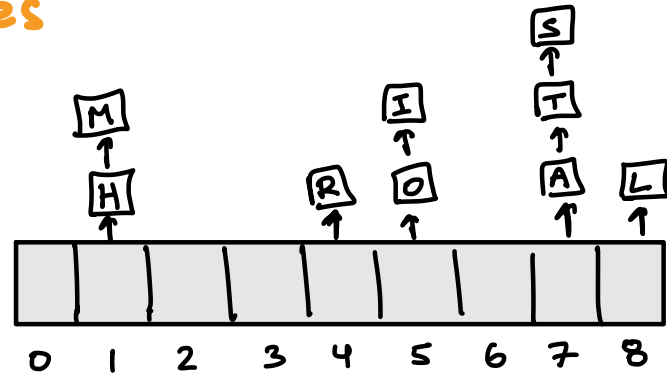
Looks a lot like balls in bins!

- Load factor = $\frac{m}{n}$
- Let $l(x) = \# \text{ of elements in the same bin as } x$
 $\# \{y \in S : h(y) = h(x)\}$
"collisions"
Worst-case lookup time = $\max_{x \in U} l(x)$
- Time to lookup $x \in U$ is $O(l(x))$

Application: Hash Tables

How should we choose
the hash function

$h: \mathcal{U} \rightarrow \{0, 1, \dots, n-1\}$ to have
small maximum load?



Looks a lot like balls in bins!

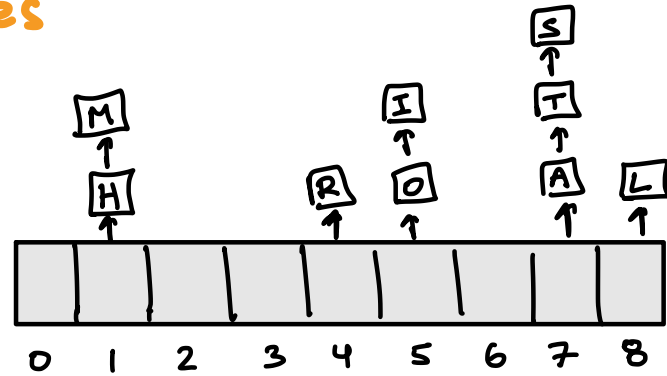
Randomized hash function:

- Model h as a uniformly random function $\mathcal{U} \rightarrow \{0, 1, \dots, n-1\}$
- Fix the set S and study $\mathbb{E}(\max_x l(x))$
↖ expectation over random choice of h

Application: Hash Tables

How should we choose the hash function

$h: \mathcal{U} \rightarrow \{0, 1, \dots, n-1\}$ to have small maximum load?



Looks a lot like balls in bins!

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

- 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

$$\mathbb{E}(\ell(x)) = \mathbb{E}\left(\sum_{y \in \mathcal{U}} \mathbb{1}_{h(y)=h(x)}\right) = \sum_{y \in \mathcal{U}} \mathbb{E}\left(\mathbb{1}_{h(y)=h(x)}\right) = n \cdot \left(\frac{1}{n}\right) = 1$$

$\mathbb{P}(h(x)=h(y)) = 1/n$

Universal Hash Families

A hash family is a set of hash functions

$$\mathcal{H} \subseteq \{ h: \mathcal{U} \rightarrow \{0, 1, \dots, n-1\} \}$$

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

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

h

Choose h
from \mathcal{H} with
equal probability

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

Constructing Universal Hashing

Construction:

- Fix a prime $p \geq |U|$, bms n
- Let $h_{a,b}(x) \equiv (ax + b \bmod p) \bmod n$

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

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

$$\Pr_x(h(x) = h(y)) \leq \frac{1}{n}$$

$$|\mathcal{H}_{all}| = n^{|U|} \approx n^p$$

$$|\mathcal{H}_{p,n}| = (p-1)p$$

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

Constructing Universal Hashing

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

Lemma 1: For every prime p

and $a \neq 0$ there is a unique

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

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

$$\text{if } az = az' \bmod p \Rightarrow a(z - z') = 0 \bmod p$$

$$\text{for } z, z' \in \{1, 2, \dots, p-1\} \Rightarrow z - z' \text{ divisible by } p$$

$$\Rightarrow z - z' = 0$$

② $az = 0 \bmod p$ has no solutions

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

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

\nwarrow
 p fixed

Constructing Universal Hashing

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

$$ax + b = r \pmod{p}$$

$$ay + b = s \pmod{p}$$

has a unique solution

Proof:

$$a = x^{-1}(r - b)$$
$$x^{-1}(r - b)y + b = s$$
$$b = s - x^{-1}(r - b)y$$

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

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

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

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

\nwarrow
 p fixed

Constructing Universal Hashing

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

Thm. $\mathcal{H}_{p,n}$ is 2-universal

Proof: By Lemma 2,

$$\mathbb{P}_{a,b}(ax+b=r \bmod p \wedge ay+b=s \bmod p) = \frac{1}{p(p-1)}$$

$$\mathbb{P}_{a,b}(h_{a,b}(x) = h_{a,b}(y)) = \frac{N}{p(p-1)}$$

$$\leq \frac{p(p-1)}{n \cdot p(p-1)}$$

$$= \frac{1}{n}$$

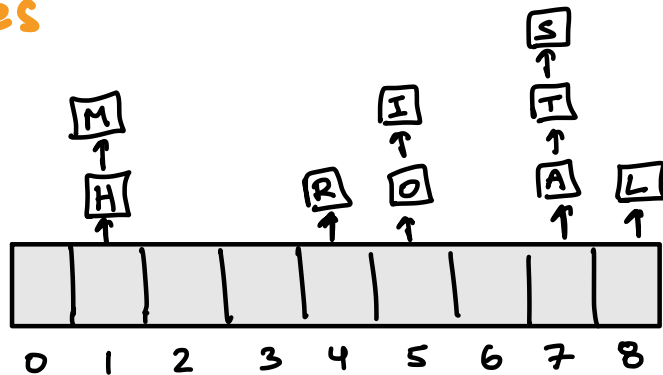
where N is number of $r \neq s$ in \mathbb{Z}_p
so that $r = s \bmod n$

$$N \leq \underbrace{p}_{\text{choices of } r} \cdot \underbrace{\frac{p-1}{n}}_{\text{choices of } s \text{ for a given } r}$$

Application: Hash Tables

How should we choose the hash function

$h: \mathcal{U} \rightarrow \{0, 1, \dots, n-1\}$ to have small maximum load?



Looks a lot like balls in bins!

Universal

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

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

- For any $x \in \mathcal{U}$, expected lookup time is

$$\mathbb{E}(l(x)) = \sum_{y \in \mathcal{U}} \mathbb{P}(h(x) = h(y)) \leq n \cdot \frac{1}{n} = 1$$

Still true for universal hash families!