CS 473: Algorithms

Ruta Mehta

University of Illinois, Urbana-Champaign

Spring 2021

High Probability Analysis & Universal Hashing

Lecture 09 Feb 23, 2021

Outline

Randomized QuickSort w.h.p. (any questions?)

What is the probability that the algorithm will terminate in $O(n \log n)$ time?

Balls & Bins

- Expected bin size.
- Expected max bin size \rightarrow max size w.h.p.
- Analogy to hashing

Hashing

Part I

Randomized **QuickSort** (Contd.)

Randomized QuickSort: Recall

Input: Array **A** of **n** distinct numbers. **Output:** Numbers in sorted order.

Randomized QuickSort

- Pick a pivot element uniformly at random from A.
- Split array into 2 subarrays: those smaller than pivot (L), and those larger than pivot (R).
- Recursively sort the subarrays, and concatenate them.

Randomized QuickSort: Recall

Input: Array **A** of **n** distinct numbers. **Output:** Numbers in sorted order.

Randomized QuickSort

- Pick a pivot element uniformly at random from A.
- Split array into 2 subarrays: those smaller than pivot (L), and those larger than pivot (R).
- Recursively sort the subarrays, and concatenate them.

Note: On *every* input randomized **QuickSort** takes $O(n \log n)$ time in expectation. On *every* input it may take $\Omega(n^2)$ time with some small probability.

Randomized QuickSort: Recall

Input: Array **A** of **n** distinct numbers. **Output:** Numbers in sorted order.

Randomized QuickSort

- Pick a pivot element uniformly at random from A.
- Split array into 2 subarrays: those smaller than pivot (L), and those larger than pivot (R).
- Recursively sort the subarrays, and concatenate them.

Note: On *every* input randomized **QuickSort** takes $O(n \log n)$ time in expectation. On *every* input it may take $\Omega(n^2)$ time with some small probability.

Question: With what probability it takes $O(n \log n)$ time?

Informal Statement

Random variable Q(A) = # comparisons done by the algorithm.

We will show that $Pr[Q(A) \leq 32n \ln n] \geq 1 - 1/n^3$.

Informal Statement

Random variable Q(A) = # comparisons done by the algorithm.

We will show that $Pr[Q(A) \leq 32n \ln n] \geq 1 - 1/n^3$.

If n = 100 then this gives $Pr[Q(A) \le 32n \ln n] \ge 0.99999$.

Informal Statement

We will show that $Pr[Q(A) \leq 32n \ln n] \geq 1 - 1/n^3$.

Outline of the proof

- If depth of recursion is k then $Q(A) \leq kn$.
- Prove that depth of recursion $\leq 32 \ln n$ with high probability (w.h.p.) . This will imply the result.

Informal Statement

We will show that $Pr[Q(A) \leq 32n \ln n] \geq 1 - 1/n^3$.

Outline of the proof

- If depth of recursion is k then $Q(A) \leq kn$.
- Prove that depth of recursion $\leq 32 \ln n$ with high probability (w.h.p.) . This will imply the result.
 - Focus on a single element. Prove that it "participates" in $> 32 \ln n$ levels with probability (w.p.) at most $1/n^4$.
 - 2 By union bound, any of the *n* elements participates in > 32 ln *n* levels w.p. at most

Informal Statement

We will show that $Pr[Q(A) \leq 32n \ln n] \geq 1 - 1/n^3$.

Outline of the proof

- If depth of recursion is k then $Q(A) \leq kn$.
- Prove that depth of recursion $\leq 32 \ln n$ with high probability (w.h.p.) . This will imply the result.
 - Focus on a single element. Prove that it "participates" in $> 32 \ln n$ levels with probability (w.p.) at most $1/n^4$.
 - 2 By union bound, any of the n elements participates in $> 32 \ln n$ levels w.p. at most $1/n^3$.

Informal Statement

We will show that $Pr[Q(A) \leq 32n \ln n] \geq 1 - 1/n^3$.

Outline of the proof

- If depth of recursion is k then $Q(A) \leq kn$.
- Prove that depth of recursion $\leq 32 \ln n$ with high probability (w.h.p.) . This will imply the result.
 - Focus on a single element. Prove that it "participates" in $> 32 \ln n$ levels with probability (w.p.) at most $1/n^4$.
 - 2 By union bound, any of the n elements participates in $> 32 \ln n$ levels w.p. at most $1/n^3$.
 - **3** Therefore, all elements participate in $\leq 32 \ln n$ w.p. $(1 1/n^3)$.

Informal Statement

An element participates in $> 32 \ln n$ w.p. $\le 1/n^4$.

Intuition

• When we pick a pivot from an array of size n uniformly at random, what is the probability that its rank is between n/4 and 3n/4?

Informal Statement

An element participates in $> 32 \ln n$ w.p. $\le 1/n^4$.

Intuition

• When we pick a pivot from an array of size n uniformly at random, what is the probability that its rank is between n/4 and 3n/4? 1/2.

Informal Statement

An element participates in $> 32 \ln n$ w.p. $\le 1/n^4$.

Intuition

- When we pick a pivot from an array of size n uniformly at random, what is the probability that its rank is between n/4 and 3n/4? 1/2.
- ② If we pick such a pivot then the size of L and R is at most?

Informal Statement

An element participates in $> 32 \ln n$ w.p. $\le 1/n^4$.

Intuition

- When we pick a pivot from an array of size n uniformly at random, what is the probability that its rank is between n/4 and 3n/4? 1/2.
- If we pick such a pivot then the size of L and R is at most? 3n/4. (Balanced split)

Informal Statement

An element participates in $> 32 \ln n$ w.p. $\le 1/n^4$.

Intuition

- When we pick a pivot from an array of size n uniformly at random, what is the probability that its rank is between n/4 and 3n/4? 1/2.
- If we pick such a pivot then the size of L and R is at most? 3n/4. (Balanced split)
- If an array is reduced to at least its 3/4th size every time, then after how many rounds only one element remains?

Informal Statement

An element participates in $> 32 \ln n$ w.p. $\le 1/n^4$.

Intuition

- When we pick a pivot from an array of size n uniformly at random, what is the probability that its rank is between n/4 and 3n/4? 1/2.
- If we pick such a pivot then the size of L and R is at most? 3n/4. (Balanced split)
- If an array is reduced to at least its 3/4th size every time, then after how many rounds only one element remains? $\leq 4 \ln n$.

Informal Statement

An element participates in $> 32 \ln n$ w.p. $\le 1/n^4$.

Intuition

- When we pick a pivot from an array of size n uniformly at random, what is the probability that its rank is between n/4 and 3n/4? 1/2.
- If we pick such a pivot then the size of L and R is at most? 3n/4. (Balanced split)
- If an array is reduced to at least its 3/4th size every time, then after how many rounds only one element remains? $< 4 \ln n$.
- If 32 In n splits, then **E**[Balanced-split] = 16 In n. Out of these there are < 4 In n balanced split w.p. $\le 1/n^4$.

• If k levels of recursion then kn comparisons.

- If k levels of recursion then kn comparisons.
- Fix an element $s \in A$. We will track it at each level.
- Let S_i be the partition containing s at i^{th} level.
- $S_1 = A$ and $S_k = \{s\}$.

- If k levels of recursion then kn comparisons.
- Fix an element $s \in A$. We will track it at each level.
- Let S_i be the partition containing s at i^{th} level.
- $S_1 = A$ and $S_k = \{s\}$.
- We call s lucky in i^{th} iteration, if balanced split: $|S_{i+1}| \leq (3/4)|S_i|$ and $|S_i \setminus S_{i+1}| \leq (3/4)|S_i|$.

- If k levels of recursion then kn comparisons.
- Fix an element $s \in A$. We will track it at each level.
- Let S_i be the partition containing s at i^{th} level.
- $S_1 = A$ and $S_k = \{s\}$.
- We call s lucky in i^{th} iteration, if balanced split: $|S_{i+1}| \leq (3/4)|S_i|$ and $|S_i \setminus S_{i+1}| \leq (3/4)|S_i|$.
- If $\rho = \#$ lucky rounds in first k rounds, then $|S_k| \leq (3/4)^{\rho} n$.

- If **k** levels of recursion then **kn** comparisons.
- Fix an element $s \in A$. We will track it at each level.
- Let S_i be the partition containing s at i^{th} level.
- $S_1 = A$ and $S_k = \{s\}$.
- We call s lucky in i^{th} iteration, if balanced split: $|S_{i+1}| \leq (3/4)|S_i|$ and $|S_i \setminus S_{i+1}| \leq (3/4)|S_i|$.
- If $\rho = \#$ lucky rounds in first k rounds, then $|S_k| \leq (3/4)^{\rho} n$.
- For $|S_k| = 1$, $\rho = 4 \ln n \ge \log_{4/3} n$ suffices.

• $X_i = 1$ if s is lucky in i^{th} iteration.

- $X_i = 1$ if s is lucky in i^{th} iteration.
- Observation: X_1, \ldots, X_k are independent variables.
- $\Pr[X_i = 1] = \frac{1}{2}$ Why?

- $X_i = 1$ if s is lucky in i^{th} iteration.
- Observation: X_1, \ldots, X_k are independent variables.
- $\Pr[X_i = 1] = \frac{1}{2}$ Why?
- Clearly, $\rho = \sum_{i=1}^k X_i$. Let $\mu = \mathbf{E}[\rho] = \frac{k}{2}$.

- $X_i = 1$ if s is lucky in i^{th} iteration.
- Observation: X_1, \ldots, X_k are independent variables.
- $\Pr[X_i = 1] = \frac{1}{2}$ Why?
- Clearly, $\rho = \sum_{i=1}^k X_i$. Let $\mu = \mathbf{E}[\rho] = \frac{k}{2}$.
- Set $k = 32 \ln n$ and $\delta = \frac{3}{4}$. $(1 \delta) = \frac{1}{4}$.

- $X_i = 1$ if s is lucky in i^{th} iteration.
- Observation: X_1, \ldots, X_k are independent variables.
- $\Pr[X_i = 1] = \frac{1}{2}$ Why?
- Clearly, $\rho = \sum_{i=1}^k X_i$. Let $\mu = \mathbf{E}[\rho] = \frac{k}{2}$.
- Set $k = 32 \ln n$ and $\delta = \frac{3}{4}$. $(1 \delta) = \frac{1}{4}$.

Probability of $\leq 4 \ln n$ lucky rounds out of $32 \ln n$ rounds is,

- $X_i = 1$ if s is lucky in i^{th} iteration.
- **Observation:** X_1, \ldots, X_k are independent variables.
- $\Pr[X_i = 1] = \frac{1}{2}$ Why?
- Clearly, $\rho = \sum_{i=1}^k X_i$. Let $\mu = \mathbf{E}[\rho] = \frac{k}{2}$.
- Set $k = 32 \ln n$ and $\delta = \frac{3}{4}$. $(1 \delta) = \frac{1}{4}$.

Probability of $\leq 4 \ln n$ lucky rounds out of $32 \ln n$ rounds is,

$$Pr[\rho \le 4 \ln n] = Pr[\rho \le \frac{k}{8}]$$

= $Pr[\rho \le (1 - \delta)\mu]$

- $X_i = 1$ if s is lucky in i^{th} iteration.
- **Observation:** X_1, \ldots, X_k are independent variables.
- $\Pr[X_i = 1] = \frac{1}{2}$ Why?
- Clearly, $\rho = \sum_{i=1}^k X_i$. Let $\mu = \mathbf{E}[\rho] = \frac{k}{2}$.
- Set $k = 32 \ln n$ and $\delta = \frac{3}{4}$. $(1 \delta) = \frac{1}{4}$.

Probability of $\leq 4 \ln n$ lucky rounds out of $32 \ln n$ rounds is,

$$\begin{array}{lll} \Pr[\rho \leq 4 \ln n] & = & \Pr[\rho \leq {}^{k}/8] \\ & = & \Pr[\rho \leq (1-\delta)\mu] \\ (\textit{Chernoff}) & \leq & 2e^{\frac{-\delta^{2}\mu}{2}} \\ & = & 2e^{-\frac{9k}{64}} \\ & = & 2e^{-4.5 \ln n} \leq \frac{1}{n^{4}} \end{array}$$

Randomized **QuickSort** w.h.p. Analysis

• n input elements. Probability that depth of recursion in **QuickSort** $> 32 \ln n$ is at most $\frac{1}{n^4} * n = \frac{1}{n^3}$.

Randomized **QuickSort** w.h.p. Analysis

• n input elements. Probability that depth of recursion in **QuickSort** $> 32 \ln n$ is at most $\frac{1}{n^4} * n = \frac{1}{n^3}$.

Theorem

With high probability (i.e., $1 - \frac{1}{n^3}$) the depth of the recursion of **QuickSort** is $\leq 32 \ln n$. Due to n comparisons in each level, with high probability, the running time of **QuickSort** is $O(n \ln n)$.

Randomized **QuickSort** w.h.p. Analysis

• n input elements. Probability that depth of recursion in **QuickSort** $> 32 \ln n$ is at most $\frac{1}{n^4} * n = \frac{1}{n^3}$.

Theorem

With high probability (i.e., $1 - \frac{1}{n^3}$) the depth of the recursion of **QuickSort** is $\leq 32 \ln n$. Due to n comparisons in each level, with high probability, the running time of **QuickSort** is $O(n \ln n)$.

Q: How to increase the probability?

Part II

Balls and Bins

Problem

If n balls are thrown independently and uniformly into n bins, how many balls lend in a bin in expectation (expected size of a bin)?

Problem

If n balls are thrown independently and uniformly into n bins, how many balls lend in a bin in expectation (expected size of a bin)?

Solution

• Fix a bin, say **j**.

Problem

If n balls are thrown independently and uniformly into n bins, how many balls lend in a bin in expectation (expected size of a bin)?

- Fix a bin, say j.
- Random variable X_{ij} is 1 if ith balls falls in jth bin, otherwise 0.

Problem

If n balls are thrown independently and uniformly into n bins, how many balls lend in a bin in expectation (expected size of a bin)?

- Fix a bin, say j.
- Random variable X_{ij} is 1 if *i*th balls falls in *j*th bin, otherwise 0.
- $E[X_{ij}] = Pr[X_{ij} = 1] =$

Problem

If n balls are thrown independently and uniformly into n bins, how many balls lend in a bin in expectation (expected size of a bin)?

- Fix a bin, say j.
- Random variable X_{ij} is 1 if ith balls falls in jth bin, otherwise 0.
- $E[X_{ij}] = Pr[X_{ij} = 1] = 1/n$.

Problem

If n balls are thrown independently and uniformly into n bins, how many balls lend in a bin in expectation (expected size of a bin)?

- Fix a bin, say j.
- Random variable X_{ij} is 1 if *i*th balls falls in *j*th bin, otherwise 0.
- $E[X_{ij}] = Pr[X_{ij} = 1] = 1/n$.
- R.V. $Y_j = \#$ balls in jth bin $= \sum_{i=1}^n X_{ij}$.

Problem

If n balls are thrown independently and uniformly into n bins, how many balls lend in a bin in expectation (expected size of a bin)?

- Fix a bin, say j.
- Random variable X_{ij} is 1 if ith balls falls in jth bin, otherwise 0.
- $E[X_{ij}] = Pr[X_{ij} = 1] = 1/n$.
- R.V. $Y_j = \#$ balls in jth bin $= \sum_{i=1}^n X_{ij}$.
- $E[Y_j] = \sum_{i=1}^n E[X_{ij}] = n \cdot 1/n = 1$.

Problem

If n balls are thrown independently and uniformly into n bins, what is the expected "maximum" bin size?

Problem

If n balls are thrown independently and uniformly into n bins, what is the expected "maximum" bin size?

$$\mathbf{E}\left[\max_{j=1}^{n} Y_{j}\right]$$
?

Problem

If n balls are thrown independently and uniformly into n bins, what is the expected "maximum" bin size?

$$\mathsf{E}\Big[\mathsf{max}_{j=1}^n Y_j\Big]?$$

Possible Solution

• R.V.
$$Z = \max_{j=1}^{n} Y_{j}$$
. $E[Z] = \sum_{k=1}^{n} Pr[Z = k] k$.

Problem

If n balls are thrown independently and uniformly into n bins, what is the expected "maximum" bin size?

$$\mathsf{E}\Big[\mathsf{max}_{j=1}^n Y_j\Big]?$$

Possible Solution

- R.V. $Z = \max_{j=1}^{n} Y_{j}$. $E[Z] = \sum_{k=1}^{n} Pr[Z = k] k$.
- How to compute Pr[Z = k], i.e., count configurations where no bin has more than k balls and at least one has k balls.

Problem

If n balls are thrown independently and uniformly into n bins, what is the expected "maximum" bin size?

$$\mathsf{E}\Big[\mathsf{max}_{j=1}^n Y_j\Big]?$$

Possible Solution

- R.V. $Z = \max_{j=1}^{n} Y_{j}$. $E[Z] = \sum_{k=1}^{n} Pr[Z = k] k$.
- How to compute Pr[Z = k], i.e., count configurations where no bin has more than k balls and at least one has k balls.
- Too many to count!!

Problem

What is the expected maximum bin size?

R.V.
$$Z = \max_{i=1}^n Y_i$$
. Show $\mathbf{E}[Z] \leq O(\frac{\ln n}{\ln \ln n})$?

Possible Solution

• If $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le 1/n^2$, then: define $A = \frac{8 \ln n}{\ln \ln n}$.

Problem

What is the expected maximum bin size?

R.V.
$$Z = \max_{j=1}^n Y_j$$
. Show $\mathbf{E}[Z] \leq O\left(\frac{\ln n}{\ln \ln n}\right)$?

Possible Solution

• If $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le 1/n^2$, then: define $A = \frac{8 \ln n}{\ln \ln n}$.

$$E[Z] = \sum_{k=1}^{n} \Pr[Z = k] k$$

$$\leq \sum_{k=1}^{A} \Pr[Z = k] A + \sum_{k=A+1}^{n} \Pr[Z = k] n$$

Problem

What is the expected maximum bin size?

R.V.
$$Z = \max_{j=1}^n Y_j$$
. Show $\mathbf{E}[Z] \leq O\left(\frac{\ln n}{\ln \ln n}\right)$?

Possible Solution

• If $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le 1/n^2$, then: define $A = \frac{8 \ln n}{\ln \ln n}$.

$$E[Z] = \sum_{k=1}^{n} \Pr[Z = k] k$$

$$\leq \sum_{k=1}^{A} \Pr[Z = k] A + \sum_{k=A+1}^{n} \Pr[Z = k] n$$

$$\leq A \cdot \Pr[Z \leq A] + n \cdot \Pr[Z > A]$$

$$\leq A \cdot (1) + n \cdot (1/n^{2}) = O(A) = O\left(\frac{\ln n}{\ln \ln n}\right)$$

Problem

What is the expected maximum bin size?

R.V.
$$Z = \max_{j=1}^n Y_j$$
. Show $\mathbf{E}[Z] \leq O\left(\frac{\ln n}{\ln \ln n}\right)$?

Possible Solution

• If $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le 1/n^2$, then: define $A = \frac{8 \ln n}{\ln \ln n}$.

$$\begin{split} \mathsf{E}[Z] &= \sum_{k=1}^n \mathsf{Pr}[Z=k] \, k \\ &\leq \sum_{k=1}^A \mathsf{Pr}[Z=k] \, A + \sum_{k=A+1}^n \mathsf{Pr}[Z=k] \, n \\ &\leq A \cdot \mathsf{Pr}[Z \leq A] + n \cdot \mathsf{Pr}[Z > A] \\ &\leq A \cdot (1) + n \cdot (1/n^2) = O(A) = O\left(\frac{\ln n}{\ln \ln n}\right) \end{split}$$

Bound $\Pr[Z > \frac{8 \ln n}{\ln \ln n}]$.

Bound $\Pr[Z > \frac{8 \ln n}{\ln \ln n}]$ using Chernoff inequality.

Chernoff Ineq. We Saw

 X_1, \ldots, X_k independent binary R.V., and $X = \sum_{i=1}^k X_i$, $\mu = \mathbf{E}[X]$, then for $0 < \delta < 1$

$$\Pr[X \geq (1+\delta)\mu] \leq e^{-\delta^2\mu/3}$$
 & $\Pr[X \leq (1-\delta)\mu] \leq e^{-\delta^2\mu/2}$

Bound $\Pr[Z > \frac{8 \ln n}{\ln \ln n}]$ using Chernoff inequality.

Chernoff Ineq. We Saw

$$X_1, \ldots, X_k$$
 independent binary R.V., and $X = \sum_{i=1}^k X_i$, $\mu = \mathbf{E}[X]$, then for $0 < \delta < 1$

$$\mathsf{Pr}[X \geq (1+\delta)\mu] \leq \mathrm{e}^{-\delta^2\mu/3}$$
 & $\mathsf{Pr}[X \leq (1-\delta)\mu] \leq \mathrm{e}^{-\delta^2\mu/2}$

Stronger Versions

- ullet For $\delta>0$, $\Pr[X>(1+\delta)\mu]<\left(rac{e^\delta}{(1+\delta)^{(1+\delta)}}
 ight)^\mu.$
- ullet For $0<\delta<1$ Pr $[X<(1-\delta)\mu]<\left(rac{e^{-\delta}}{(1-\delta)^{(1-\delta)}}
 ight)^{\mu}$

Problem

What is the expected maximum bin size? Let $Z = \max_{j=1}^n Y_j$. Show $\mathbf{E}[Z] \leq O(\frac{\ln n}{\ln \ln n})$. \to Show $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \leq 1/n^2$.

Problem

What is the expected maximum bin size? Let $Z = \max_{j=1}^n Y_j$. Show $\mathbf{E}[Z] \leq O(\frac{\ln n}{\ln \ln n})$. \to Show $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \leq 1/n^2$.

Solution

• Recall: $Y_j = \#$ balls in bin j, $\mathsf{E}[Y_j] = 1$, and $A = \frac{8 \ln n}{\ln \ln n}$

$$\Pr[Y_j > A] = \Pr[Y_j \ge A \operatorname{E}[Y]] < \left(\frac{e^{A-1}}{A^A}\right) < \left(\frac{n^{6/\ln\ln n}}{A^A}\right)$$

Problem

What is the expected maximum bin size? Let $Z = \max_{j=1}^n Y_j$. Show $\mathbf{E}[Z] \leq O(\frac{\ln n}{\ln \ln n})$. \to Show $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \leq 1/n^2$.

Solution

• Recall: $Y_j = \#$ balls in bin j, $\mathbf{E}[Y_j] = 1$, and $A = \frac{8 \ln n}{\ln \ln n}$

$$\Pr[Y_j > A] = \Pr[Y_j \ge A \operatorname{E}[Y]] < \left(\frac{e^{A-1}}{A^A}\right) < \left(\frac{n^{6/\ln \ln n}}{A^A}\right)$$

$$A^{A} = \left(\frac{8 \ln n}{\ln \ln n}\right)^{\frac{6 \ln n}{\ln \ln n}} \ge (\sqrt{\ln n})^{\frac{8 \ln n}{\ln \ln n}} = (\ln n)^{\frac{4 \ln n}{\ln \ln n}} = e^{4 \lg n} = n^4$$

Problem

What is the expected maximum bin size? Let $Z = \max_{j=1}^n Y_j$. Show $\mathbf{E}[Z] \leq O(\frac{\ln n}{\ln \ln n})$. \to Show $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \leq 1/n^2$.

Solution

• Recall: $Y_j = \#$ balls in bin j, $\mathbf{E}[Y_j] = 1$, and $A = \frac{8 \ln n}{\ln \ln n}$

$$\Pr[Y_j > A] = \Pr[Y_j \ge A \operatorname{E}[Y]] < \left(\frac{e^{A-1}}{A^A}\right) < \left(\frac{n^{6/\ln \ln n}}{A^A}\right)$$

$$A^{A} = \left(\frac{8 \ln n}{\ln \ln n}\right)^{\frac{8 \ln n}{\ln \ln n}} \ge \left(\sqrt{\ln n}\right)^{\frac{8 \ln n}{\ln \ln n}} = (\ln n)^{\frac{4 \ln n}{\ln \ln n}} = e^{4 \lg n} = n^{4}$$

$$\Pr\left[Y_{j} > \frac{8 \ln n}{\ln \ln n}\right] < 1/n^{3}$$

Problem

What is the expected maximum bin size? Let $Z = \max_{j=1}^{n} Y_j$.

Show
$$E[Z] \le O(\frac{\ln n}{\ln \ln n}) \to \text{Show } Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le 1/n^2$$
.

Solution

• Recall: $Y_j = \#$ balls in bin j. $E[Y_j] = 1$.

$$\Pr[Y_j > 8 \ln n / \ln \ln n] \le 1/n^3$$
 (Using Chernoff)

Problem

What is the expected maximum bin size? Let $Z = \max_{j=1}^{n} Y_j$.

Show
$$E[Z] \le O(\frac{\ln n}{\ln \ln n}) \to \text{Show } Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le 1/n^2$$
.

Solution

- Recall: $Y_j = \#$ balls in bin j. $E[Y_j] = 1$.
 - $\Pr[Y_j > 8 \ln n / \ln \ln n] \le 1/n^3$ (Using Chernoff)
- (Union bound) $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le \sum_{j=1}^{n} \Pr[Y_j > \frac{8 \ln n}{\ln \ln n}] \le n \cdot 1/n^3 = 1/n^2.$

Problem

What is the expected maximum bin size? Let $Z = \max_{j=1}^{n} Y_j$.

Show
$$E[Z] \le O(\frac{\ln n}{\ln \ln n}) \to \text{Show } Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le 1/n^2$$
.

Solution

- Recall: $Y_j = \#$ balls in bin j. $\mathsf{E}[Y_j] = 1$. $\mathsf{Pr}[Y_i > 8 \ln n / \ln \ln n] \le 1/n^3$ (Using Chernoff)
- (Union bound) $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le \sum_{j=1}^{n} \Pr[Y_j > \frac{8 \ln n}{\ln \ln n}] \le n \cdot 1/n^3 = 1/n^2.$
- Max bin size is at most $O(\frac{\ln n}{\ln \ln n})$ with probability $1 1/n^2$.

Problem

What is the expected maximum bin size? Let $Z = \max_{j=1}^{n} Y_j$.

Show
$$E[Z] \le O(\frac{\ln n}{\ln \ln n}) \to \text{Show } Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le 1/n^2$$
.

Solution

- Recall: $Y_j = \#$ balls in bin j. $E[Y_j] = 1$. $Pr[Y_i > 8 \ln n / \ln \ln n] < 1/n^3$ (Using Chernoff)
- (Union bound) $\Pr[Z > \frac{8 \ln n}{\ln \ln n}] \le \sum_{i=1}^{n} \Pr[Y_i > \frac{8 \ln n}{\ln \ln n}] \le n \cdot 1/n^3 = 1/n^2.$
- Max bin size is at most $O(\frac{\ln n}{\ln \ln n})$ with probability $1 1/n^2$.

 $\Omega(\frac{\ln n}{\ln \ln n})$ is a lower bound as well!

Hashing

Storing elements in a table such that look up is O(1)-time.

Hashing

Storing elements in a table such that look up is O(1)-time.

Throwing numbered balls

Imagine that n balls have numbers coming from a universe \mathcal{U} . $|\mathcal{U}| \gg n$.

Hashing

Storing elements in a table such that look up is O(1)-time.

Throwing numbered balls

Imagine that n balls have numbers coming from a universe \mathcal{U} . $|\mathcal{U}|\gg n$.

Hashing: throw balls (elements) randomly into *n* bins such that **bin** sizes are small

Hashing

Storing elements in a table such that look up is O(1)-time.

Throwing numbered balls

Imagine that n balls have numbers coming from a universe \mathcal{U} . $|\mathcal{U}|\gg n$.

Hashing: throw balls (elements) randomly into n bins such that bin sizes are small and also lookup is easy!.

Part III

Hash Tables

Dictionary Data Structure

- $oldsymbol{0}$ $oldsymbol{\mathcal{U}}$: universe of keys with total order: numbers, strings, etc.
- ② Data structure to store a subset $S \subseteq \mathcal{U}$
- Operations:
 - **o** Search/lookup: given $x \in \mathcal{U}$ is $x \in S$?
 - **2** Insert: given $x \notin S$ add x to S.
 - **3 Delete**: given $x \in S$ delete x from S

Dictionary Data Structure

- $oldsymbol{0}$ $oldsymbol{\mathcal{U}}$: universe of keys with total order: numbers, strings, etc.
- ② Data structure to store a subset $S \subseteq \mathcal{U}$
- Operations:
 - **o** Search/lookup: given $x \in \mathcal{U}$ is $x \in S$?
 - **2** Insert: given $x \notin S$ add x to S.
 - **3 Delete**: given $x \in S$ delete x from S
- Static structure: S given in advance or changes very infrequently, main operations are lookups.

Dictionary Data Structure

- $oldsymbol{0}$ $oldsymbol{\mathcal{U}}$: universe of keys with total order: numbers, strings, etc.
- ② Data structure to store a subset $S \subseteq \mathcal{U}$
- Operations:
 - **1** Search/lookup: given $x \in \mathcal{U}$ is $x \in S$?
 - **2** Insert: given $x \notin S$ add x to S.
 - **3 Delete**: given $x \in S$ delete x from S
- Static structure: S given in advance or changes very infrequently, main operations are lookups.
- Oynamic structure: S changes rapidly so inserts and deletes as important as lookups.

Dictionary Data Structures

Common solutions:

- Static:
 - Store S as a sorted array
 - **2** Lookup: Binary search in $O(\log |S|)$ time (comparisons)
- ② Dynamic:
 - Store S in a balanced binary search tree
 - 2 Lookup, Insert, Delete in $O(\log |S|)$ time (comparisons)

Dictionary Data Structures

Question: "Should Tables be Sorted?" (also title of famous paper by Turing award winner Andy Yao)

Dictionary Data Structures

Question: "Should Tables be Sorted?" (also title of famous paper by Turing award winner Andy Yao)

Hashing is a widely used & powerful technique for dictionaries.

Motivation:

- Universe \mathcal{U} may not be (naturally) totally ordered.
- Keys correspond to large objects (images, graphs etc) for which comparisons are very expensive.
- **3** Want to improve "average" performance of lookups to O(1) even at cost of extra space or errors with small probability: many applications for fast lookups in networking, security, etc.

Hash Table data structure:

- **1** A (hash) table/array T of size m (the table size).
- ② A hash function $h: \mathcal{U} \to \{0, \dots, m-1\}$.
- 1 Item $x \in \mathcal{U}$ hashes to slot h(x) in T.

Hash Table data structure:

- **1** A (hash) table/array T of size m (the table size).
- ② A hash function $h: \mathcal{U} \to \{0, \dots, m-1\}$.
- 1 Item $x \in \mathcal{U}$ hashes to slot h(x) in T.

Given $S \subseteq \mathcal{U}$. How do we store S and how do we do lookups?

Hash Table data structure:

- A (hash) table/array T of size m (the table size).
- lacktriangledown A hash function $h: \mathcal{U}
 ightarrow \{0, \ldots, m-1\}$.
- 1 Item $x \in \mathcal{U}$ hashes to slot h(x) in T.

Given $S \subseteq \mathcal{U}$. How do we store S and how do we do lookups?

Ideal situation:

- **9** Each element $x \in S$ hashes to a distinct slot in T. Store x in slot h(x)
- **2** Lookup: Given $y \in \mathcal{U}$ check if T[h(y)] = y. O(1) time!

Hash Table data structure:

- A (hash) table/array T of size m (the table size).
- ② A hash function $h: \mathcal{U} \to \{0, \dots, m-1\}$.
- 1 Item $x \in \mathcal{U}$ hashes to slot h(x) in T.

Given $S \subseteq \mathcal{U}$. How do we store S and how do we do lookups?

Ideal situation:

- **9** Each element $x \in S$ hashes to a distinct slot in T. Store x in slot h(x)
- **2** Lookup: Given $y \in \mathcal{U}$ check if T[h(y)] = y. O(1) time!

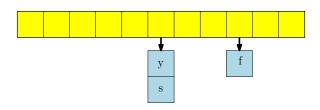
Collisions unavoidable if $|T| < |\mathcal{U}|$. Several techniques to handle them.

Handling Collisions: Chaining

Collision: h(x) = h(y) for some $x \neq y$.

Chaining to handle collisions:

- For each slot i store all items hashed to slot i in a linked list.
 T[i] points to the linked list
- **2** Lookup: to find if $y \in \mathcal{U}$ is in T, check the linked list at T[h(y)]. Time proportion to size of linked list.



This is also known as **Open hashing**.

Handling Collisions

Several other techniques:

• Cuckoo hashing. Every value has two possible locations. When inserting, insert in one of the locations, otherwise, kick stored value to its other location. Repeat till stable. if no stability then rebuild table.

- **2** ...
- Others.

Does hashing give O(1) time per operation for dictionaries?

Does hashing give O(1) time per operation for dictionaries?

Complexity of evaluating h on a given element?

Does hashing give O(1) time per operation for dictionaries?

- Complexity of evaluating **h** on a given element? Should be small.
- 2 Relative sizes of the universe \mathcal{U} and the set to be stored S.

Does hashing give O(1) time per operation for dictionaries?

- Complexity of evaluating **h** on a given element? Should be small.
- **2** Relative sizes of the universe \mathcal{U} and the set to be stored S. Typically $|\mathcal{U}| \gg |S|$.

Does hashing give O(1) time per operation for dictionaries?

- Complexity of evaluating **h** on a given element? Should be small.
- **2** Relative sizes of the universe \mathcal{U} and the set to be stored S. Typically $|\mathcal{U}| \gg |S|$.
- Size of table T relative to size of S.

Does hashing give O(1) time per operation for dictionaries?

- Complexity of evaluating **h** on a given element? Should be small.
- **②** Relative sizes of the universe \mathcal{U} and the set to be stored S. Typically $|\mathcal{U}| \gg |S|$.
- 3 Size of table T relative to size of S. |S|/|T| is a constant smaller than 1.

Does hashing give O(1) time per operation for dictionaries?

- Complexity of evaluating **h** on a given element? Should be small.
- **2** Relative sizes of the universe \mathcal{U} and the set to be stored S. Typically $|\mathcal{U}| \gg |S|$.
- 3 Size of table T relative to size of S. |S|/|T| is a constant smaller than 1.
 - **Load factor/Fill factor:** The ratio n/m where n = |S| and m = |T|.

Does hashing give O(1) time per operation for dictionaries?

- Complexity of evaluating **h** on a given element? Should be small.
- **2** Relative sizes of the universe \mathcal{U} and the set to be stored S. Typically $|\mathcal{U}| \gg |S|$.
- 3 Size of table T relative to size of S. |S|/|T| is a constant smaller than 1.

```
Load factor/Fill factor: The ratio n/m where n = |S| and m = |T|.
```

Main and interrelated questions:

- Worst-case vs average-case vs randomized (expected) time?
- 2 How do we choose h?

- **1** \mathcal{U} : universe (very large).
- ② Assume $N = |\mathcal{U}| \gg m$ where m is size of table T. In particular assume $N \geq m^2$ (very conservative).

- **1** U: universe (very large).
- Assume $N = |\mathcal{U}| \gg m$ where m is size of table T. In particular assume $N > m^2$ (very conservative).
- **3** Fix hash function $h: \mathcal{U} \to \{0, \dots, m-1\}$.

CS473 Spring 2021 28 / 52

- U: universe (very large).
- ② Assume $N = |\mathcal{U}| \gg m$ where m is size of table T. In particular assume $N \geq m^2$ (very conservative).
- $lacksquare{0}$ Fix hash function $h:\mathcal{U} o \{0,\ldots,m-1\}$.
- N items hashed to m slots. Minimize the max load. Howmuch is it?

- **1** U: universe (very large).
- ② Assume $N = |\mathcal{U}| \gg m$ where m is size of table T. In particular assume $N > m^2$ (very conservative).
- \bullet Fix hash function $h: \mathcal{U} \to \{0, \ldots, m-1\}$.
- N items hashed to m slots. Minimize the max load. Howmuch is it? By pigeon hole principle, N/m > m!.

CS473 28 Spring 2021 28 / 52

- \mathcal{U} : universe (very large).
- ② Assume $N = |\mathcal{U}| \gg m$ where m is size of table T. In particular assume $N \geq m^2$ (very conservative).
- $lackbox{0}$ Fix hash function $h:\mathcal{U} o \{0,\ldots,m-1\}$.
- N items hashed to m slots. Minimize the max load. Howmuch is it? By pigeon hole principle, N/m ≥ m!.
- Implies that there is a set $S \subseteq \mathcal{U}$ where |S| = m such that all of S hashes to same slot. Ooops.

- U: universe (very large).
- ② Assume $N = |\mathcal{U}| \gg m$ where m is size of table T. In particular assume $N \geq m^2$ (very conservative).
- $lackbox{0}$ Fix hash function $h:\mathcal{U} o \{0,\ldots,m-1\}$.
- N items hashed to m slots. Minimize the max load. Howmuch is it? By pigeon hole principle, N/m ≥ m!.
- Implies that there is a set $S \subseteq \mathcal{U}$ where |S| = m such that all of S hashes to same slot. Ooops.

Lesson: For every hash function there is a very bad set. Bad set. Bad.

How many hash functions are there, anyway?

Let $\mathcal H$ be the set of all functions from $\mathcal U=\{1,\ldots,U\}$ to $\{1,\ldots,m\}$. The number of functions in $\mathcal H$ is

- (A) U + m.
- (B) Um.
- (C) U^m.
- (D) m^U .
- (E) $\binom{U+m}{m}$.
- **(F)** The answer is blowing in the wind.

How many bits one need?

Let $\mathcal H$ be a set of functions from $\mathcal U=\{1,\ldots,U\}$ to $\{1,\ldots,m\}$. Specifying a function in $\mathcal H$ requires:

- (A) O(U+m) bits.
- (B) O(Um) bits.
- (C) $O(U^m)$ bits.
- (D) $O(m^U)$ bits.
- (E) $O(\log |\mathcal{H}|)$ bits.
- (F) Many many bits. At least two.

- Hash function are often chosen in an ad hoc fashion. Implicit assumption is that input behaves well.
- May work well for aircraft control. Susceptible to denial of service attack in routing.

- Hash function are often chosen in an ad hoc fashion. Implicit assumption is that input behaves well.
- May work well for aircraft control. Susceptible to denial of service attack in routing.

Parameters: $N = |\mathcal{U}|$, m = |T|, n = |S|

9 \mathcal{H} is a **family** of hash functions: each function $h \in \mathcal{H}$ should be efficient to evaluate (that is, to compute h(x)).

- Hash function are often chosen in an ad hoc fashion. Implicit assumption is that input behaves well.
- May work well for aircraft control. Susceptible to denial of service attack in routing.

Parameters: $N = |\mathcal{U}|, m = |\mathcal{T}|, n = |\mathcal{S}|$

- **1** \mathcal{H} is a **family** of hash functions: each function $h \in \mathcal{H}$ should be efficient to evaluate (that is, to compute h(x)).
- \bullet h is chosen randomly from \mathcal{H} (typically uniformly at random). Implicitly assumes that \mathcal{H} allows an efficient sampling.

Ruta (UIUC) CS473 31 Spring 2021 31 / 52

- Hash function are often chosen in an ad hoc fashion. Implicit assumption is that input behaves well.
- May work well for aircraft control. Susceptible to denial of service attack in routing.

Parameters: $N = |\mathcal{U}|, m = |T|, n = |S|$

- **9** \mathcal{H} is a **family** of hash functions: each function $h \in \mathcal{H}$ should be efficient to evaluate (that is, to compute h(x)).
- **2** h is chosen randomly from \mathcal{H} (typically uniformly at random). Implicitly assumes that \mathcal{H} allows an efficient sampling.
- ② Randomized guarantee: should have the property that for any fixed set $S \subseteq \mathcal{U}$ of size m the expected number of collisions for a function chosen from \mathcal{H} should be "small". Here the expectation is over the randomness in choice of h.

Question: Why not let \mathcal{H} be the set of *all* functions from \mathcal{U} to $\{0,1,\ldots,m-1\}$?

Question: Why not let $\mathcal H$ be the set of *all* functions from $\mathcal U$ to $\{0,1,\ldots,m-1\}$?

1 Too many functions! A random function has high complexity! # of functions: $M = m^{|\mathcal{U}|}$. Bits to encode such a function $\approx \log M = |\mathcal{U}| \log m$.

Question: Why not let \mathcal{H} be the set of *all* functions from \mathcal{U} to $\{0,1,\ldots,m-1\}$?

① Too many functions! A random function has high complexity! # of functions: $M = m^{|\mathcal{U}|}$. Bits to encode such a function $\approx \log M = |\mathcal{U}| \log m$.

Question: Are there good and compact families \mathcal{H} ?

Question: Why not let \mathcal{H} be the set of *all* functions from \mathcal{U} to $\{0,1,\ldots,m-1\}$?

● Too many functions! A random function has high complexity! # of functions: $M = m^{|\mathcal{U}|}$. Bits to encode such a function $\approx \log M = |\mathcal{U}| \log m$.

Question: Are there good and compact families \mathcal{H} ?

lacktriangle Yes... But what it means for ${\cal H}$ to be good and compact.

Question: What are good properties of \mathcal{H} in distributing data?

Question: What are good properties of \mathcal{H} in distributing data?

• Consider any element $x \in \mathcal{U}$. If $h \in \mathcal{H}$ is picked randomly then x should go into a random slot in T. In other words $\Pr[h(x) = i] = 1/m$ for every $0 \le i < m$. (Uniform)

Question: What are good properties of \mathcal{H} in distributing data?

- Consider any element $x \in \mathcal{U}$. If $h \in \mathcal{H}$ is picked randomly then x should go into a random slot in T. In other words $\Pr[h(x) = i] = 1/m$ for every $0 \le i < m$. (Uniform)
- ② Consider any two distinct elements $x, y \in \mathcal{U}$. Then if $h \in \mathcal{H}$ is picked randomly then the probability of a collision between x and y should be at most 1/m. In other words $\Pr[h(x) = h(y)] = 1/m$ (cannot be smaller).

Question: What are good properties of \mathcal{H} in distributing data?

- Consider any element $x \in \mathcal{U}$. If $h \in \mathcal{H}$ is picked randomly then x should go into a random slot in T. In other words $\Pr[h(x) = i] = 1/m$ for every $0 \le i < m$. (Uniform)
- ② Consider any two distinct elements $x, y \in \mathcal{U}$. Then if $h \in \mathcal{H}$ is picked randomly then the probability of a collision between x and y should be at most 1/m. In other words $\Pr[h(x) = h(y)] = 1/m$ (cannot be smaller).
- Second property is stronger than the first and is crucial.

Definition

A family of hash function \mathcal{H} is (2-)universal if for all distinct $x, y \in \mathcal{U}$, $\Pr_h[h(x) = h(y)] = 1/m$ where m is the table size.

Analyzing Universal Hashing

- T is hash table of size m.
- **2** $S \subseteq \mathcal{U}$ is a **fixed** set of size $\leq m$.
- **1** Is chosen randomly from a universal hash family \mathcal{H} .
- x is a fixed element of U.

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

• The time to look up x is the size of the list at T[h(x)]: same as the number of elements in S that collide with x under h.

Ruta (UIUC) CS473 35 Spring 2021 35 / 52

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

- ① The time to look up x is the size of the list at T[h(x)]: same as the number of elements in S that collide with x under h.
- ② Let $\ell(x)$ be this number. We want $E[\ell(x)]$

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

- ① The time to look up x is the size of the list at T[h(x)]: same as the number of elements in S that collide with x under h.
- ② Let $\ell(x)$ be this number. We want $E[\ell(x)]$
- **3** For $y \in S$ let A_y be the event that x, y collide and D_y be the corresponding indicator variable.

Ruta (UIUC) CS473 35 Spring 2021 35 / 52

Continued...

Number of elements colliding with x: $\ell(x) = \sum_{y \in S} D_y$.

Ruta (UIUC) CS473 36 Spring 2021 36 / 52

Continued...

Number of elements colliding with x: $\ell(x) = \sum_{y \in S} D_y$.

$$\Rightarrow \mathbb{E}[\ell(x)] = \sum_{y \in S} \mathbb{E}[D_y] \quad \text{linearity of expectation}$$

$$= \sum_{y \in S} \Pr[h(x) = h(y)]$$

$$= \sum_{y \in S} \frac{1}{m} \quad \text{(since } \mathcal{H} \text{ is a universal hash family)}$$

$$= |S|/m$$

$$= \frac{n}{m}$$

$$< 1 \quad \text{(if } |S| < m)$$

Ruta (UIUC) CS473 36 Spring 2021 36 / 52

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

Answer: O(n/m).

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

Answer: O(n/m).

Comments:

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

Answer: O(n/m).

Comments:

- 0 O(1) expected time also holds for insertion.
- ② Analysis assumes static set S but holds as long as S is a set formed with at most O(m) insertions and deletions.
- Worst-case: look up time can be large! How large?

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

Answer: O(n/m).

Comments:

- 0 O(1) expected time also holds for insertion.
- ② Analysis assumes static set S but holds as long as S is a set formed with at most O(m) insertions and deletions.
- Worst-case: look up time can be large! How large? $\Omega(\log n/\log \log n)$ [Lower bound holds even under stronger assumptions.]

Universal: \mathcal{H} such that $\Pr[h(x) = h(y)] = 1/m$.

Ruta (UIUC) CS473 38 Spring 2021 38 / 52

Universal: \mathcal{H} such that $\Pr[h(x) = h(y)] = 1/m$.

All functions

 $\mathcal{H}:$ Set of all possible functions $h:\mathcal{U} o \{0,\ldots,m-1\}.$

Universal.

Universal: \mathcal{H} such that $\Pr[h(x) = h(y)] = 1/m$.

All functions

 $\mathcal{H}:$ Set of all possible functions $h:\mathcal{U} \to \{0,\ldots,m-1\}$.

- Universal.
- $\bullet |\mathcal{H}| = m^{|\mathcal{U}|}$
- representing h requires $|\mathcal{U}| \log m$ Not O(1)!

Universal: \mathcal{H} such that $\Pr[h(x) = h(y)] = 1/m$.

All functions

 $\mathcal{H}:$ Set of all possible functions $h:\mathcal{U} o \{0,\ldots,m-1\}.$

- Universal.
- $\bullet |\mathcal{H}| = m^{|\mathcal{U}|}$
- representing h requires $|\mathcal{U}| \log m$ Not O(1)!

We need compactly representable universal family.

Parameters:
$$N = |\mathcal{U}|, m = |T|, n = |S|$$

① Choose a **prime** number p > N. Define function $h_{a,b}(x) = ((ax + b) \mod p) \mod m$.

Ruta (UIUC) CS473 39 Spring 2021 39 / 52

Parameters: $N = |\mathcal{U}|, m = |T|, n = |S|$

Choose a prime number p > N. Define function h_{a,b}(x) = ((ax + b) mod p) mod m.
 Let H = {h_{a,b} | a, b ∈ Z_p, a ≠ 0} (Z_p = {0,1,...,p-1}).

Ruta (UIUC) CS473 39 Spring 2021 39 / 52

Parameters: $N = |\mathcal{U}|, m = |T|, n = |S|$

- Choose a **prime** number p > N. Define function $h_{a,b}(x) = ((ax + b) \mod p) \mod m$.
- ② Let $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ $(\mathbb{Z}_p = \{0, 1, \dots, p-1\})$. Note that $|\mathcal{H}| = p(p-1)$.

Ruta (UIUC) CS473 39 Spring 2021 39 / 52

- Parameters: $N = |\mathcal{U}|, m = |T|, n = |S|$
 - Choose a **prime** number p > N. Define function $h_{a,b}(x) = ((ax + b) \mod p) \mod m$.
 - ② Let $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ $(\mathbb{Z}_p = \{0, 1, \dots, p-1\})$. Note that $|\mathcal{H}| = p(p-1)$.

Theorem

H is a universal hash family.

Parameters:
$$N = |\mathcal{U}|, m = |T|, n = |S|$$

- Choose a **prime** number p > N. Define function $h_{a,b}(x) = ((ax + b) \mod p) \mod m$.
- ② Let $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\} \ (\mathbb{Z}_p = \{0, 1, \dots, p-1\}).$ Note that $|\mathcal{H}| = p(p-1)$.

$\mathsf{Theorem}$

 \mathcal{H} is a universal hash family.

Comments:

- **1** $h_{a,b}$ can be evaluated in O(1) time.
- 2 Easy to store, i.e., just store a, b. Easy to sample.

CS473 39 Spring 2021 39 / 52

Some math required...

Lemma (LemmaUnique)

Let p be a prime number, and $\mathbb{Z}_p = \{0, 1, \dots, p-1\}$. x: an integer number in \mathbb{Z}_p , $x \neq 0$ \implies There exists a unique $y \in \mathbb{Z}_p$ s.t. $xy = 1 \mod p$.

In other words: For every element there is a unique inverse. \implies set $\mathbb{Z}_p = \{0, 1, \dots, p-1\}$ when working modulo p is a field.

Ruta (UIUC) CS473 40 Spring 2021 40 / 5

Claim

Let p be a prime number. For any $x, y, z \in \{1, ..., p-1\}$ s.t. $y \neq z$, we have that $xy \mod p \neq xz \mod p$.

Claim

Let **p** be a prime number. For any $x, y, z \in \{1, \dots, p-1\}$ s.t. $y \neq z$, we have that $xy \mod p \neq xz \mod p$.

Proof.

Assume for the sake of contradiction $xy \mod p = xz \mod p$. Then

$$x(y-z) = 0 \mod p$$

 $\implies p \text{ divides } x(y-z)$
 $\implies p \text{ divides } y-z$
 $\implies y-z=0 \implies y=z$

And that is a contradiction.

41

Lemma (LemmaUnique)

```
Let p be a prime number,

x: an integer number in \{1, \ldots, p-1\}.

\implies There exists a unique y s.t. xy = 1 \mod p.
```

Proof.

By the above claim if $xy = 1 \mod p$ and $xz = 1 \mod p$ then y = z. Hence uniqueness follows.

Lemma (LemmaUnique)

```
Let p be a prime number,

x: an integer number in \{1, \ldots, p-1\}.

\implies There exists a unique y s.t. xy = 1 \mod p.
```

Proof.

By the above claim if $xy = 1 \mod p$ and $xz = 1 \mod p$ then y = z. Hence uniqueness follows.

Existence. For any $x \in \{1, \dots, p-1\}$ we have that $\{x*1 \mod p, x*2 \mod p, \dots, x*(p-1) \mod p\} =$

Lemma (LemmaUnique)

```
Let p be a prime number,

x: an integer number in \{1, \ldots, p-1\}.

\implies There exists a unique y s.t. xy = 1 \mod p.
```

Proof.

By the above claim if $xy = 1 \mod p$ and $xz = 1 \mod p$ then y = z. Hence uniqueness follows.

```
Existence. For any x \in \{1, \ldots, p-1\} we have that \{x*1 \mod p, x*2 \mod p, \ldots, x*(p-1) \mod p\} = \{1, 2, \ldots, p-1\}. \Longrightarrow There exists a number y \in \{1, \ldots, p-1\} such that xy = 1 \mod p.
```

Proof of the Theorem: Outline

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

Theorem

 $\mathcal{H} = \{h_{a,b} \mid a,b \in \mathbb{Z}_p, a \neq 0\}$ is universal.

Proof.

Fix $x, y \in \mathcal{U}$. We need to show that

$$\Pr_{h_{a,b}\sim\mathcal{H}}[h_{a,b}(x)=h_{a,b}(y)]\leq 1/m$$
. Note that $|\mathcal{H}|=p(p-1)$.

Ruta (UIUC) CS473 43 Spring 2021 43 / 52

Proof of the Theorem: Outline

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

Theorem

 $\mathcal{H} = \{h_{a,b} \mid a,b \in \mathbb{Z}_p, a \neq 0\}$ is universal.

Proof.

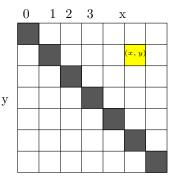
Fix $x, y \in \mathcal{U}$. We need to show that

$$\Pr_{h_{a,b}\sim\mathcal{H}}[h_{a,b}(x)=h_{a,b}(y)]\leq 1/m$$
. Note that $|\mathcal{H}|=p(p-1)$.

- Let (a, b) (equivalently $h_{a,b}$) be bad for x, y if $h_{a,b}(x) = h_{a,b}(y)$.
- 2 Claim: Number of bad (a, b) is at most p(p-1)/m.
- **3** Total number of hash functions is p(p-1) and hence probability of a collision is $\leq 1/m$.

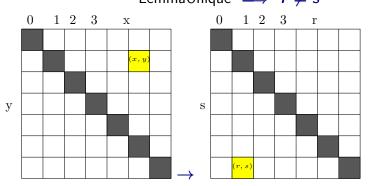
Ruta (UIUC) CS473 43 Spring 2021 43 / 52

$$g_{a,b}(x) = (ax + b) \mod p$$
, $h_{a,b}(x) = (g_{a,b}(x)) \mod m$
First map $x \neq y$ to $r = g_{a,b}(x)$ and $s = g_{a,b}(y)$.
LemmaUnique $\implies r \neq s$



Ruta (UIUC) CS473 44 Spring 2021 44 / 52

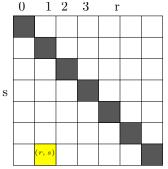
$$g_{a,b}(x) = (ax + b) \mod p$$
, $h_{a,b}(x) = (g_{a,b}(x)) \mod m$
First map $x \neq y$ to $r = g_{a,b}(x)$ and $s = g_{a,b}(y)$.
LemmaUnique $\implies r \neq s$



As (a, b) varies, (r, s) takes all possible p(p - 1) values. Since (a, b) is picked u.a.r., every value of (r, s) has equal probability.

Ruta (UIUC) CS473 44 Spring 2021 44 / 52

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

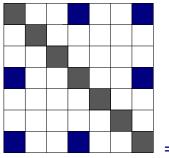






Ruta (UIUC) CS473 44 Spring 2021 44 / 5:

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



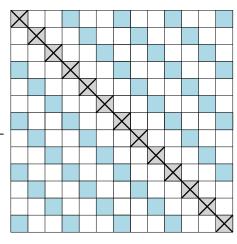




Ruta (UIUC) CS473 44 Spring 2021 44 / 52

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

- First part of mapping maps (x, y) to a random location $(g_{a,b}(x), g_{a,b}(y))$ in the "matrix".
- $(g_{a,b}(x), g_{a,b}(y))$ is not on main diagonal.
- All blue locations are "bad" map by mod m to a location of collision.
- But... at most 1/m fraction of allowable locations in the matrix are bad.



We need

to show at most 1/m fraction of bad $h_{a,b}$

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

2 lemmas ...

Fix
$$x \neq y \in \mathbb{Z}_p$$
, and let $r = (ax + b) \mod p$ and $s = (ay + b) \mod p$.

Ruta (UIUC) CS473 45 Spring 2021 45 / 52

to show at most 1/m fraction of bad $h_{a,b}$

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

2 lemmas ...

Fix $x \neq y \in \mathbb{Z}_p$, and let $r = (ax + b) \mod p$ and $s = (ay + b) \mod p$.

1-to-1 correspondence between p(p-1) pairs of (a,b) (equivalently $h_{a,b}$) and p(p-1) pairs of (r,s).

Ruta (UIUC) CS473 45 Spring 2021 45 / 52

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

2 lemmas ...

Fix $x \neq y \in \mathbb{Z}_p$, and let $r = (ax + b) \mod p$ and $s = (ay + b) \mod p$.

- **1** 1-to-1 correspondence between p(p-1) pairs of (a,b)(equivalently $h_{a,b}$) and p(p-1) pairs of (r,s).
- ② Out of all possible p(p-1) pairs of (r,s), at most p(p-1)/m fraction satisfies $r \mod m = s \mod m$.

CS473 45 Spring 2021 45 / 52

Some Lemmas

Lemma

If $x \neq y$ then for any $a, b \in \mathbb{Z}_p$ such that $a \neq 0$, we have $ax + b \mod p \neq ay + b \mod p$.

Some Lemmas

Lemma

If $x \neq y$ then for any $a, b \in \mathbb{Z}_p$ such that $a \neq 0$, we have $ax + b \mod p \neq ay + b \mod p$.

Proof.

Suppose not

$$ax + b \mod p = ay + b \mod p \Rightarrow a(x - y) \mod p = 0$$

Ruta (UIUC) CS473 46 Spring 2021 46 / 52

Lemma

If $x \neq y$ then for any $a, b \in \mathbb{Z}_p$ such that $a \neq 0$, we have $ax + b \mod p \neq ay + b \mod p$.

Proof.

Suppose not

$$ax + b \mod p = ay + b \mod p \Rightarrow a(x - y) \mod p = 0$$

But, $a \neq 0$ and $(x - y) \neq 0$.

Lemma

If $x \neq y$ then for any $a, b \in \mathbb{Z}_p$ such that $a \neq 0$, we have $ax + b \mod p \neq ay + b \mod p$.

Proof.

Suppose not

$$ax + b \mod p = ay + b \mod p \Rightarrow a(x - y) \mod p = 0$$

But, $a \neq 0$ and $(x - y) \neq 0$. And a and (x - y) cannot divide p since p is prime and a < p and (x - y) < p. Contradiction!

Ruta (UIUC) CS473 46 Spring 2021 46 / 52

Lemma

If $x \neq y$ then for each (r, s) such that $r \neq s$ and 0 < r, s < p-1 there is exactly one a, b such that $ax + b \mod p = r$ and $ay + b \mod p = s$

Proof.

Solve the two equations:

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

Lemma

If $x \neq y$ then for each (r, s) such that $r \neq s$ and $0 \leq r, s \leq p-1$ there is exactly one a, b such that $ax + b \mod p = r$ and $ay + b \mod p = s$

Proof.

Solve the two equations:

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

We get
$$a = \frac{r-s}{x-y} \mod p$$
 and $b = r - ax \mod p$.

One-to-one correspondence between (a, b) and (r, s)

Understanding the hashing

Once we fix a and b, and we are given a value x, we compute the hash value of x in two stages:

- **1** Compute: $r \leftarrow (ax + b) \mod p$.
- **2** Fold: $r' \leftarrow r \mod m$

Collision...

Given two distinct values x and y they might collide only because of folding.

Lemma

not equal pairs (r, s) of $\mathbb{Z}_p \times \mathbb{Z}_p$ that are folded to the same number is p(p-1)/m.

Folding numbers

Lemma

pairs $(r, s) \in \mathbb{Z}_p \times \mathbb{Z}_p$ such that $r \neq s$ and $r \mod m = s$ mod m (folded to the same number) is p(p-1)/m.

Proof.

Consider a pair $(r,s) \in \{0,1,\ldots,p-1\}^2$ s.t. $r \neq s$. Fix r:

 $\mathbf{0}$ $a = r \mod m$.

Folding numbers

Lemma

pairs $(r, s) \in \mathbb{Z}_p \times \mathbb{Z}_p$ such that $r \neq s$ and $r \mod m = s$ mod m (folded to the same number) is p(p-1)/m.

Proof.

Consider a pair $(r, s) \in \{0, 1, \dots, p-1\}^2$ s.t. $r \neq s$. Fix r:

- $\mathbf{0}$ $a = r \mod m$.
- ② There are $\lceil p/m \rceil$ values of s that fold into a. That is

 $r \mod m = s \mod m$.

- 3 One of them is when r = s.
- $\bullet \implies \#$ of colliding pairs

Folding numbers

Lemma

pairs $(r, s) \in \mathbb{Z}_p \times \mathbb{Z}_p$ such that $r \neq s$ and $r \mod m = s$ mod m (folded to the same number) is p(p-1)/m.

Proof.

Consider a pair $(r, s) \in \{0, 1, \dots, p-1\}^2$ s.t. $r \neq s$. Fix r:

- $\mathbf{0}$ $a = r \mod m$.
- ② There are $\lceil p/m \rceil$ values of s that fold into a. That is

 $r \mod m = s \mod m$.

- **3** One of them is when r = s.
- \implies # of colliding pairs $(\lceil p/m \rceil 1)p \le (p-1)p/m$

Proof of Claim

of bad pairs is p(p-1)/m

Proof.

Let $a, b \in \mathbb{Z}_p$ such that $a \neq 0$ and $h_{a,b}(x) = h_{a,b}(y)$.

- ② Collision if and only if $r \mod m = s \mod m$.
- (Folding error): Number of pairs (r, s) such that $r \neq s$ and $0 \leq r, s \leq p-1$ and $r \mod m = s \mod m$ is p(p-1)/m.
- From previous lemma there is one-to-one correspondence between (a, b) and (r, s). Hence total number of bad (a, b) pairs is p(p-1)/m.

Ruta (UIUC) CS473 50 Spring 2021 50 / 52

Proof of Claim

of bad pairs is p(p-1)/m

Proof.

Let $a, b \in \mathbb{Z}_p$ such that $a \neq 0$ and $h_{a,b}(x) = h_{a,b}(y)$.

- ② Collision if and only if $r \mod m = s \mod m$.
- **③** (Folding error): Number of pairs (r, s) such that $r \neq s$ and $0 \leq r, s \leq p-1$ and $r \mod m = s \mod m$ is p(p-1)/m.
- From previous lemma there is one-to-one correspondence between (a, b) and (r, s). Hence total number of bad (a, b) pairs is p(p-1)/m.

Prob of x and y to collide: $\frac{\# \text{ bad } (a,b) \text{ pairs}}{\#(a,b) \text{ pairs}} = \frac{p(p-1)/m}{p(p-1)} = \frac{1}{m}$.

Ruta (UIUC) CS473 50 Spring 2021 50 / 52

Rehashing, amortization and...

.. making the hash table dynamic

So far we assumed fixed S of size $\simeq m$.

Question: What happens as items are inserted and deleted?

- ① If |S| grows to more than cm for some constant c then hash table performance clearly degrades.
- ② If |S| stays around $\simeq m$ but incurs many insertions and deletions then the initial random hash function is no longer random enough!

Ruta (UIUC) CS473 51 Spring 2021 51 / 52

Rehashing, amortization and...

... making the hash table dynamic

So far we assumed fixed **S** of size $\simeq m$.

Question: What happens as items are inserted and deleted?

- If |S| grows to more than cm for some constant c then hash table performance clearly degrades.
- ② If |S| stays around $\simeq m$ but incurs many insertions and deletions then the initial random hash function is no longer random enough!

Solution: Rebuild hash table periodically!

- Choose a new table size based on current number of elements in table.
- Choose a new random hash function and rehash the elements.
- Oiscard old table and hash function.

Question: When to rebuild? How expensive?

Ruta (UIUC) CS473 51 Spring 2021 51 / 52

Rebuilding the hash table

- **9** Start with table size m where m is some estimate of |S| (can be some large constant).
- ② If |S| grows to more than twice current table size, build new hash table (choose a new random hash function) with double the current number of elements. Can also use similar trick if table size falls below quarter the size.
- If |S| stays roughly the same but more than c|S| operations on table for some chosen constant c (say 10), rebuild.

The **amortize** cost of rebuilding to previously performed operations. Rebuilding ensures O(1) expected analysis holds even when S changes. Hence O(1) expected look up/insert/delete time *dynamic* data dictionary data structure!