
Algorithms for Big Data (II)

Chihao Zhang

Shanghai Jiao Tong University

Sept. 27, 2019

REVIEW OF LAST LECTURE

REVIEW OF LAST LECTURE

Last time, we met the **streaming model**.

REVIEW OF LAST LECTURE

Last time, we met the **streaming model**.

We studied Morris' algorithm for counting the number of elements in a data stream.

REVIEW OF LAST LECTURE

Last time, we met the **streaming model**.

We studied Morris' algorithm for counting the number of elements in a data stream.

We used **Averaging trick** and **Median trick** to boost the quality of Morris' algorithm.

REVIEW OF LAST LECTURE

Last time, we met the **streaming model**.

We studied Morris' algorithm for counting the number of elements in a data stream.

We used **Averaging trick** and **Median trick** to boost the quality of Morris' algorithm.

Today we will take a closer look at the mathematical tools needed in the course.

MARKOV'S INEQUALITY

MARKOV'S INEQUALITY

Markov's inequality

For every **nonnegative** random variable X and every $a \geq 0$, it holds that

$$\Pr [X \geq a] \leq \frac{\mathbf{E}[X]}{a}.$$

MARKOV'S INEQUALITY

Markov's inequality

For every **nonnegative** random variable X and every $a \geq 0$, it holds that

$$\Pr [X \geq a] \leq \frac{\mathbf{E}[X]}{a}.$$

Proof.

Let $\mathbf{1}_{X \geq a}$ be the indicator random variable such that $\mathbf{1}_{X \geq a}(x) = \begin{cases} 1, & \text{if } x \geq a, \\ 0, & \text{otherwise.} \end{cases}$

Then it holds that $X \geq a \cdot \mathbf{1}_{X \geq a}$. Take the expectation on both sides, we obtain

$$\mathbf{E}[X] \geq a \cdot \mathbf{E}[\mathbf{1}_{X \geq a}] = a \cdot \Pr[X \geq a].$$



CHEBYSHEV'S INEQUALITY

CHEBYSHEV'S INEQUALITY

Chebyshev's inequality

For every random variable X and every $a \geq 0$, it holds that

$$\Pr [|X - \mathbf{E} [X]| \geq a] \leq \frac{\mathbf{Var} [X]}{a^2}.$$

CHEBYSHEV'S INEQUALITY

Chebyshev's inequality

For every random variable X and every $a \geq 0$, it holds that

$$\Pr [|X - \mathbf{E}[X]| \geq a] \leq \frac{\mathbf{Var}[X]}{a^2}.$$

Proof.

$$\begin{aligned} \Pr [|X - \mathbf{E}[X]| \geq a] &= \Pr \left[(X - \mathbf{E}[X])^2 \geq a^2 \right] \\ &\leq \frac{\mathbf{E} \left[(X - \mathbf{E}[X])^2 \right]}{a^2} \quad (\text{Markov's inequality}) \\ &= \frac{\mathbf{Var}[X]}{a^2}. \end{aligned}$$

□

CHERNOFF'S BOUND

CHERNOFF'S BOUND

Chernoff bound

Let X_1, \dots, X_n be **independent** Bernoulli trials with $\mathbf{E}[X_i] = p_i$ for every $i = 1, \dots, n$. Let $X = \sum_{i=1}^n X_i$. Then for every $0 < \varepsilon < 1$, it holds that

$$\Pr[|X - \mathbf{E}[X]| > \varepsilon \cdot \mathbf{E}[X]] \leq 2 \exp\left(-\frac{\varepsilon^2 \mathbf{E}[X]}{3}\right).$$

CHERNOFF'S BOUND

Chernoff bound

Let X_1, \dots, X_n be **independent** Bernoulli trials with $\mathbf{E}[X_i] = p_i$ for every $i = 1, \dots, n$. Let $X = \sum_{i=1}^n X_i$. Then for every $0 < \varepsilon < 1$, it holds that

$$\Pr[|X - \mathbf{E}[X]| > \varepsilon \cdot \mathbf{E}[X]] \leq 2 \exp\left(-\frac{\varepsilon^2 \mathbf{E}[X]}{3}\right).$$

The main tool to prove Chernoff bound is the **moment generating function** e^{tX} for a random variable X .

CHERNOFF'S BOUND

Chernoff bound

Let X_1, \dots, X_n be **independent** Bernoulli trials with $\mathbf{E}[X_i] = p_i$ for every $i = 1, \dots, n$. Let $X = \sum_{i=1}^n X_i$. Then for every $0 < \varepsilon < 1$, it holds that

$$\Pr[|X - \mathbf{E}[X]| > \varepsilon \cdot \mathbf{E}[X]] \leq 2 \exp\left(-\frac{\varepsilon^2 \mathbf{E}[X]}{3}\right).$$

The main tool to prove Chernoff bound is the **moment generating function** e^{tX} for a random variable X .

It holds that

$$\begin{aligned} \mathbf{E}[e^{tX}] &= \mathbf{E}\left[e^{t \sum_{i=1}^n X_i}\right] = \prod_{i=1}^n \mathbf{E}[e^{tX_i}] = \prod_{i=1}^n ((1 - p_i) + p_i e^t) \\ &= \prod_{i=1}^n (1 - (1 - e^t)p_i) \leq \prod_{i=1}^n e^{-(1-e^t)p_i} = e^{-(1-e^t)\mathbf{E}[X]}. \end{aligned}$$

PROOF OF CHERNOFF BOUND

PROOF OF CHERNOFF BOUND

For every $t > 0$, we have

$$\Pr [X \geq (1 + \varepsilon)\mathbf{E}[X]] = \Pr \left[e^{tX} \geq e^{t(1+\varepsilon)\mathbf{E}[X]} \right] \leq \frac{\mathbf{E} [e^{tX}]}{e^{t(1+\varepsilon)\mathbf{E}[X]}} \leq \frac{e^{-(1-e^t)\mathbf{E}[X]}}{e^{t(1+\varepsilon)\mathbf{E}[X]}}.$$

PROOF OF CHERNOFF BOUND

For every $t > 0$, we have

$$\Pr [X \geq (1 + \varepsilon)\mathbf{E}[X]] = \Pr \left[e^{tX} \geq e^{t(1+\varepsilon)\mathbf{E}[X]} \right] \leq \frac{\mathbf{E} [e^{tX}]}{e^{t(1+\varepsilon)\mathbf{E}[X]}} \leq \frac{e^{-(1-e^t)\mathbf{E}[X]}}{e^{t(1+\varepsilon)\mathbf{E}[X]}}.$$

To find an optimal t , we calculate the derivative of above and obtain for $t = \log(1 + \varepsilon)$,

$$\Pr [X \geq (1 + \varepsilon)\mathbf{E}[X]] \leq \left(\frac{e^\varepsilon}{(1 + \varepsilon)^{1+\varepsilon}} \right)^{\mathbf{E}[X]} \leq e^{-\varepsilon^2\mathbf{E}[X]/3}.$$

PROOF OF CHERNOFF BOUND

For every $t > 0$, we have

$$\Pr [X \geq (1 + \varepsilon)\mathbf{E}[X]] = \Pr \left[e^{tX} \geq e^{t(1+\varepsilon)\mathbf{E}[X]} \right] \leq \frac{\mathbf{E} [e^{tX}]}{e^{t(1+\varepsilon)\mathbf{E}[X]}} \leq \frac{e^{-(1-e^t)\mathbf{E}[X]}}{e^{t(1+\varepsilon)\mathbf{E}[X]}}.$$

To find an optimal t , we calculate the derivative of above and obtain for $t = \log(1 + \varepsilon)$,

$$\Pr [X \geq (1 + \varepsilon)\mathbf{E}[X]] \leq \left(\frac{e^\varepsilon}{(1 + \varepsilon)^{1+\varepsilon}} \right)^{\mathbf{E}[X]} \leq e^{-\varepsilon^2\mathbf{E}[X]/3}.$$

We can similarly prove that

$$\Pr [X \leq (1 - \varepsilon)\mathbf{E}[X]] \leq e^{-\varepsilon^2\mathbf{E}[X]/2}.$$

PROOF OF CHERNOFF BOUND

For every $t > 0$, we have

$$\Pr [X \geq (1 + \varepsilon)\mathbf{E}[X]] = \Pr \left[e^{tX} \geq e^{t(1+\varepsilon)\mathbf{E}[X]} \right] \leq \frac{\mathbf{E} [e^{tX}]}{e^{t(1+\varepsilon)\mathbf{E}[X]}} \leq \frac{e^{-(1-e^t)\mathbf{E}[X]}}{e^{t(1+\varepsilon)\mathbf{E}[X]}}.$$

To find an optimal t , we calculate the derivative of above and obtain for $t = \log(1 + \varepsilon)$,

$$\Pr [X \geq (1 + \varepsilon)\mathbf{E}[X]] \leq \left(\frac{e^\varepsilon}{(1 + \varepsilon)^{1+\varepsilon}} \right)^{\mathbf{E}[X]} \leq e^{-\varepsilon^2\mathbf{E}[X]/3}.$$

We can similarly prove that

$$\Pr [X \leq (1 - \varepsilon)\mathbf{E}[X]] \leq e^{-\varepsilon^2\mathbf{E}[X]/2}.$$

Combining the bounds for both lower and upper tails, we finish the proof.

BALLS-INTO-BINS

Balls-into-Bins is a simple yet important probabilistic model.

BALLS-INTO-BINS

Balls-into-Bins is a simple yet important probabilistic model.

Suppose we throw m ball into n bins **uniformly** and **independently**, what is the (expected) **maxload** of the bins?

BALLS-INTO-BINS

Balls-into-Bins is a simple yet important probabilistic model.

Suppose we throw m ball into n bins **uniformly** and **independently**, what is the (expected) **maxload** of the bins?

When $m = n$, the answer is $\Theta\left(\frac{\log n}{\log \log n}\right)$.

BALLS-INTO-BINS

Balls-into-Bins is a simple yet important probabilistic model.

Suppose we throw m ball into n bins **uniformly** and **independently**, what is the (expected) **maxload** of the bins?

When $m = n$, the answer is $\Theta\left(\frac{\log n}{\log \log n}\right)$.

It models an important object, the Hash functions.

INDEPENDENCE

INDEPENDENCE

A set of random variables X_1, \dots, X_n are **mutually independent** if for every index set $I \subseteq [n]$ and values $\{x_i\}_{i \in I}$,

$$\Pr \left[\bigwedge_{i \in I} X_i = x_i \right] = \prod_{i=1}^n \Pr [X_i = x_i].$$

INDEPENDENCE

A set of random variables X_1, \dots, X_n are **mutually independent** if for every index set $I \subseteq [n]$ and values $\{x_i\}_{i \in I}$,

$$\Pr \left[\bigwedge_{i \in I} X_i = x_i \right] = \prod_{i=1}^n \Pr [X_i = x_i].$$

k-WISE INDEPENDENCE

k-WISE INDEPENDENCE

A weaker notion of independence is the *k-wise independence*.

k -WISE INDEPENDENCE

A weaker notion of independence is the k -wise independence.

A set of random variables X_1, \dots, X_n are k -wise independent if for every index set $I \subseteq [n]$ with $|I| \leq k$ and values $\{x_i\}_{i \in I}$,

$$\Pr \left[\bigwedge_{i \in I} X_i = x_i \right] = \prod_{i=1}^n \Pr [X_i = x_i].$$

k -WISE INDEPENDENCE

A weaker notion of independence is the k -wise independence.

A set of random variables X_1, \dots, X_n are k -wise independent if for every index set $I \subseteq [n]$ with $|I| \leq k$ and values $\{x_i\}_{i \in I}$,

$$\Pr \left[\bigwedge_{i \in I} X_i = x_i \right] = \prod_{i=1}^n \Pr [X_i = x_i].$$

We call X_1, \dots, X_n $\text{pairwise independent}$ if they are 2-wise independent.

EXAMPLES

Suppose we have n independent bits $X_1, \dots, X_n \in \{0, 1\}$.

EXAMPLES

Suppose we have n independent bits $X_1, \dots, X_n \in \{0, 1\}$.

For every $I \in [n]$, define $Y_I = \left(\sum_{j \in I} X_j \right) \bmod 2$.

EXAMPLES

Suppose we have n independent bits $X_1, \dots, X_n \in \{0, 1\}$.

For every $I \in [n]$, define $Y_I = \left(\sum_{j \in I} X_j \right) \bmod 2$.

The random bits $\{Y_I\}_{I \subseteq [n]}$ are pairwise independent.

EXAMPLES

Suppose we have n independent bits $X_1, \dots, X_n \in \{0, 1\}$.

For every $I \in [n]$, define $Y_I = \left(\sum_{j \in I} X_j \right) \bmod 2$.

The random bits $\{Y_I\}_{I \subseteq [n]}$ are pairwise independent.

But they are not mutually independent!

PROPERTY OF PAIRWISE INDEPENDENCE

PROPERTY OF PAIRWISE INDEPENDENCE

Theorem

For pairwise independent X_1, \dots, X_n , we have

$$\mathbf{Var} [X_1 + \dots + X_n] = \mathbf{Var} [X_1] + \dots + \mathbf{Var} [X_n].$$

Proof.

$$\begin{aligned}\mathbf{Var} [X_1 + \dots + X_n] &= \mathbf{E} [(X_1 + \dots + X_n)^2] - (\mathbf{E} [X_1 + \dots + X_n])^2 \\ &= \sum_{i=1}^n \mathbf{E} [X_i^2] + 2 \sum_{1 \leq i < j \leq n} \mathbf{E} [X_i X_j] - \left(\sum_{i=1}^n \mathbf{E} [X_i]^2 + 2 \sum_{1 \leq i < j \leq n} \mathbf{E} [X_i] \mathbf{E} [X_j] \right) \\ &= \sum_{i=1}^n (\mathbf{E} [X_i^2] - \mathbf{E} [X_i]^2) = \sum_{i=1}^n \mathbf{Var} [X_i].\end{aligned}$$

□

HASH FUNCTIONS

HASH FUNCTIONS

In Balls-into-Bins, we distribute balls **uniformly** and **independently**.

HASH FUNCTIONS

In Balls-into-Bins, we distribute balls **uniformly** and **independently**.

This can be implemented using **Hash functions**

HASH FUNCTIONS

In Balls-into-Bins, we distribute balls **uniformly** and **independently**.

This can be implemented using **Hash functions**

Hash functions are important data structures that have been widely used in computer science.

HASH FUNCTIONS

In Balls-into-Bins, we distribute balls **uniformly** and **independently**.

This can be implemented using **Hash functions**

Hash functions are important data structures that have been widely used in computer science.

We will construct Hash functions with **theoretical guarantees**.

UNIVERSAL HASH FUNCTION FAMILIES

UNIVERSAL HASH FUNCTION FAMILIES

Let \mathcal{H} be a family of functions from $[m]$ to $[n]$ where $m \geq n$.

UNIVERSAL HASH FUNCTION FAMILIES

Let \mathcal{H} be a family of functions from $[m]$ to $[n]$ where $m \geq n$.

We call \mathcal{H} **k -universal** if for every distinct $x_1, \dots, x_k \in [m]$, we have

$$\Pr_{h \in \mathcal{H}} [h(x_1) = h(x_2) = \dots = h(x_k)] \leq \frac{1}{n^{k-1}}.$$

UNIVERSAL HASH FUNCTION FAMILIES

Let \mathcal{H} be a family of functions from $[m]$ to $[n]$ where $m \geq n$.

We call \mathcal{H} **k -universal** if for every distinct $x_1, \dots, x_k \in [m]$, we have

$$\Pr_{h \in \mathcal{H}} [h(x_1) = h(x_2) = \dots = h(x_k)] \leq \frac{1}{n^{k-1}}.$$

We call \mathcal{H} **strongly k -universal** if for every distinct $x_1, \dots, x_k \in [m]$, $y_1, \dots, y_k \in [n]$, we have

$$\Pr_{h \in \mathcal{H}} \left[\bigwedge_{i=1}^k h(x_i) = y_i \right] = \frac{1}{n^k}.$$

BALLS-INTO-BINS WITH 2-UNIVERSAL HASH FAMILY

BALLS-INTO-BINS WITH 2-UNIVERSAL HASH FAMILY

Let X_{ij} be the indicator of the event: i -th ball and j -th ball fall into the same bin.

BALLS-INTO-BINS WITH 2-UNIVERSAL HASH FAMILY

Let X_{ij} be the indicator of the event: i -th ball and j -th ball fall into the same bin.

Let $X = \sum_{1 \leq i < j \leq m} X_{ij}$ be the total number of collisions. Then

$$\mathbf{E}[X] = \sum_{1 \leq i < j \leq m} \mathbf{E}[X_{ij}] \leq \binom{m}{2} \frac{1}{n} < \frac{m^2}{2n}.$$

BALLS-INTO-BINS WITH 2-UNIVERSAL HASH FAMILY

Let X_{ij} be the indicator of the event: i -th ball and j -th ball fall into the same bin.

Let $X = \sum_{1 \leq i < j \leq m} X_{ij}$ be the total number of collisions. Then

$$\mathbf{E}[X] = \sum_{1 \leq i < j \leq m} \mathbf{E}[X_{ij}] \leq \binom{m}{2} \frac{1}{n} < \frac{m^2}{2n}.$$

Assume the maxload is Y , which causes $\binom{Y}{2} \leq X$ collisions. Then

$$\Pr \left[\binom{Y}{2} \geq \frac{m^2}{n} \right] \leq \Pr \left[X \geq \frac{m^2}{n} \right] \leq \frac{1}{n}.$$

BALLS-INTO-BINS WITH 2-UNIVERSAL HASH FAMILY

Let X_{ij} be the indicator of the event: i -th ball and j -th ball fall into the same bin.

Let $X = \sum_{1 \leq i < j \leq m} X_{ij}$ be the total number of collisions. Then

$$\mathbf{E}[X] = \sum_{1 \leq i < j \leq m} \mathbf{E}[X_{ij}] \leq \binom{m}{2} \frac{1}{n} < \frac{m^2}{2n}.$$

Assume the maxload is Y , which causes $\binom{Y}{2} \leq X$ collisions. Then

$$\Pr \left[\binom{Y}{2} \geq \frac{m^2}{n} \right] \leq \Pr \left[X \geq \frac{m^2}{n} \right] \leq \frac{1}{n}.$$

Therefore, $\Pr \left[Y - 1 \geq m\sqrt{2/n} \right] \leq \frac{1}{2}$. The maxload is $1 + \sqrt{2n}$ when $m = n$ with probability at least $1/2$.

CONSTRUCTION OF 2-UNIVERSAL FAMILY

CONSTRUCTION OF 2-UNIVERSAL FAMILY

Now we explicitly construct a universal family of Hash functions from $[m]$ to $[n]$.

CONSTRUCTION OF 2-UNIVERSAL FAMILY

Now we explicitly construct a universal family of Hash functions from $[m]$ to $[n]$.

Let $p \geq m$ be a prime and let

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

CONSTRUCTION OF 2-UNIVERSAL FAMILY

Now we explicitly construct a universal family of Hash functions from $[m]$ to $[n]$.

Let $p \geq m$ be a prime and let

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

The family is

$$\mathcal{H} = \{h_{a,b} : 1 \leq a \leq p-1, 0 \leq b \leq p-1\}.$$

PROOF

PROOF

We show that \mathcal{H} constructed above is indeed 2-universal.

PROOF

We show that \mathcal{H} constructed above is indeed 2-universal.

We compute the colliding probability

$$\Pr_{h_{a,b} \in \mathcal{H}} [h_{a,b}(x) = h_{a,b}(y)]$$

for $x \neq y$.

PROOF

We show that \mathcal{H} constructed above is indeed 2-universal.

We compute the colliding probability

$$\Pr_{h_{a,b} \in \mathcal{H}} [h_{a,b}(x) = h_{a,b}(y)]$$

for $x \neq y$.

First, we have if $x \neq y$, then $ax + b \neq ay + b \pmod{p}$.

Moreover $(a, b) \rightarrow (ax + b, ay + b)$ is a **bijection** from $\{1, \dots, p-1\} \times \{0, \dots, p-1\}$ to $\{(u, v) : 0 \leq u, v \leq p-1, u \neq v\}$.

PROOF

We show that \mathcal{H} constructed above is indeed 2-universal.

We compute the colliding probability

$$\Pr_{h_{a,b} \in \mathcal{H}} [h_{a,b}(x) = h_{a,b}(y)]$$

for $x \neq y$.

First, we have if $x \neq y$, then $ax + b \neq ay + b \pmod p$.

Moreover $(a, b) \rightarrow (ax + b, ay + b)$ is a **bijection** from $\{1, \dots, p-1\} \times \{0, \dots, p-1\}$ to $\{(u, v) : 0 \leq u, v \leq p-1, u \neq v\}$.

This is because
$$\begin{cases} ax + b = u \pmod p \\ ay + b = v \pmod p \end{cases} \text{ has a unique solution } \begin{cases} a = \frac{v-u}{y-x} \pmod p \\ b = u - ax \pmod p. \end{cases}$$

PROOF (CONT'D)

PROOF (CONT'D)

Therefore,

$$\Pr_{h_{a,b} \in \mathcal{H}} [h_{a,b}(x) = h_{a,b}(y)] = \Pr_{(u,v) \in \mathbb{F}_p^2: u \neq v} [u = v \pmod n].$$

PROOF (CONT'D)

Therefore,

$$\Pr_{h_{a,b} \in \mathcal{H}} [h_{a,b}(x) = h_{a,b}(y)] = \Pr_{(u,v) \in \mathbb{F}_p^2: u \neq v} [u = v \pmod{n}].$$

The number of (u, v) with $u \neq v$ is $p(p-1)$.

PROOF (CONT'D)

Therefore,

$$\Pr_{h_{a,b} \in \mathcal{H}} [h_{a,b}(x) = h_{a,b}(y)] = \Pr_{(u,v) \in \mathbb{F}_p^2: u \neq v} [u = v \pmod n].$$

The number of (u, v) with $u \neq v$ is $p(p-1)$.

For each u , the number of values of v with $u = v \pmod n$ is at most $\lceil p/n \rceil - 1$.

PROOF (CONT'D)

Therefore,

$$\Pr_{h_{a,b} \in \mathcal{H}} [h_{a,b}(x) = h_{a,b}(y)] = \Pr_{(u,v) \in \mathbb{F}_p^2: u \neq v} [u = v \pmod n].$$

The number of (u, v) with $u \neq v$ is $p(p-1)$.

For each u , the number of values of v with $u = v \pmod n$ is at most $\lceil p/n \rceil - 1$.

The probability is therefore at most

$$\frac{p(\lceil p/n \rceil - 1)}{p(p-1)} \leq \frac{1}{n}.$$