
Algorithms for Big Data (III)

Chihao Zhang

Shanghai Jiao Tong University

Sept. 29, 2019

REVIEW OF THE LAST LECTURE

REVIEW OF THE LAST LECTURE

Last time, we proved a few useful **concentration inequalities**.

REVIEW OF THE LAST LECTURE

Last time, we proved a few useful **concentration inequalities**.

We introduced the notion of **universal families of Hash functions**.

REVIEW OF THE LAST LECTURE

Last time, we proved a few useful **concentration inequalities**.

We introduced the notion of **universal families of Hash functions**.

We constructed a 2-universal universal family of Hash functions.

REVIEW: THE CONSTRUCTION

Let m, n be two integer and $p \geq m$ be a prime.

REVIEW: THE CONSTRUCTION

Let m, n be two integer and $p \geq m$ be a prime.

The family

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

where each $h_{a,b} : [m] \rightarrow [n]$ is defined as

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

REVIEW: THE CONSTRUCTION

Let m, n be two integer and $p \geq m$ be a prime.

The family

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

where each $h_{a,b} : [m] \rightarrow [n]$ is defined as

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

We proved that for every $x \neq y$,

$$\Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{n}.$$

REVIEW: THE CONSTRUCTION

Let m, n be two integer and $p \geq m$ be a prime.

The family

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

where each $h_{a,b} : [m] \rightarrow [n]$ is defined as

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

We proved that for every $x \neq y$,

$$\Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{n}.$$

So \mathcal{H} is a **2-universal Hash function family**.

STRONGLY 2-UNIVERSAL HASH FAMILY

STRONGLY 2-UNIVERSAL HASH FAMILY

Recall that if we further require that for any u, v ,

$$\Pr_{h \in \mathcal{H}} [h(x) = u \wedge h(y) = v] = \frac{1}{n^2},$$

then \mathcal{H} is called **strongly 2-universal family of Hash functions**.

STRONGLY 2-UNIVERSAL HASH FAMILY

Recall that if we further require that for any u, v ,

$$\Pr_{h \in \mathcal{H}} [h(x) = u \wedge h(y) = v] = \frac{1}{n^2},$$

then \mathcal{H} is called **strongly 2-universal family of Hash functions**.

When $m = n = p$ are primes, then we can modify the previously constructed \mathcal{H} to get a strong 2-universal family.

STRONGLY 2-UNIVERSAL HASH FAMILY

Recall that if we further require that for any u, v ,

$$\Pr_{h \in \mathcal{H}} [h(x) = u \wedge h(y) = v] = \frac{1}{n^2},$$

then \mathcal{H} is called **strongly 2-universal family of Hash functions**.

When $m = n = p$ are primes, then we can modify the previously constructed \mathcal{H} to get a strong 2-universal family.

In this case, we have

$$\mathcal{H} = \{h_{a,b}(x) = ax + b \pmod{p} : 0 \leq a, b \leq p - 1\}.$$

PROOF

PROOF

Lemma

The equation $ax + b = 0 \pmod p$ has unique solution (in \mathbb{F}_p) if $a \neq 0$ and p is a prime.

PROOF

Lemma

The equation $ax + b = 0 \pmod p$ has unique solution (in \mathbb{F}_p) if $a \neq 0$ and p is a prime.

The equations $h_{a,b}(x_1) = y_1$ and $h_{a,b}(x_2) = y_2$ are equivalent to

$$ax_1 + b = y_1 \pmod p, \quad ax_2 + b = y_2 \pmod p.$$

PROOF

Lemma

The equation $ax + b = 0 \pmod p$ has unique solution (in \mathbb{F}_p) if $a \neq 0$ and p is a prime.

The equations $h_{a,b}(x_1) = y_1$ and $h_{a,b}(x_2) = y_2$ are equivalent to

$$ax_1 + b = y_1 \pmod p, \quad ax_2 + b = y_2 \pmod p.$$

They have a unique solution

$$a = \frac{y_2 - y_1}{x_2 - x_1} \pmod p, \quad b = y_1 - ax_1 \pmod p.$$

PROOF

Lemma

The equation $ax + b = 0 \pmod p$ has unique solution (in \mathbb{F}_p) if $a \neq 0$ and p is a prime.

The equations $h_{a,b}(x_1) = y_1$ and $h_{a,b}(x_2) = y_2$ are equivalent to

$$ax_1 + b = y_1 \pmod p, \quad ax_2 + b = y_2 \pmod p.$$

They have a unique solution

$$a = \frac{y_2 - y_1}{x_2 - x_1} \pmod p, \quad b = y_1 - ax_1 \pmod p.$$

Therefore,

$$\Pr_{h_{a,b} \in H} [h_{a,b}(x_1) = y_1 \wedge h_{a,b}(x_2) = y_2] = \frac{1}{p^2}.$$

THE GENERAL CASE

THE GENERAL CASE

The Hash family we just constructed has the restriction that $m = n$

THE GENERAL CASE

The Hash family we just constructed has the restriction that $m = n$

We can naturally generalize $m = p$ to $m = p^k$.

THE GENERAL CASE

The Hash family we just constructed has the restriction that $m = n$

We can naturally generalize $m = p$ to $m = p^k$.

Write every number x in base p :

$$x = x_0 + x_1 \cdot p + x_2 \cdot p^2 + \dots x_{k-1} \cdot p^{k-1}.$$

THE GENERAL CASE

The Hash family we just constructed has the restriction that $m = n$

We can naturally generalize $m = p$ to $m = p^k$.

Write every number x in base p :

$$x = x_0 + x_1 \cdot p + x_2 \cdot p^2 + \dots x_{k-1} \cdot p^{k-1}.$$

For every $\bar{a} = (a_0, a_1, \dots, a_{k-1})$, with $0 \leq a_i \leq p - 1$ and $0 \leq b \leq p - 1$, define

$$h_{\bar{a}, b}(x) = \left(\sum_{i=0}^{k-1} a_i x_i + b \right) \bmod p.$$

THE GENERAL CASE

The Hash family we just constructed has the restriction that $m = n$

We can naturally generalize $m = p$ to $m = p^k$.

Write every number x in base p :

$$x = x_0 + x_1 \cdot p + x_2 \cdot p^2 + \dots x_{k-1} \cdot p^{k-1}.$$

For every $\bar{a} = (a_0, a_1, \dots, a_{k-1})$, with $0 \leq a_i \leq p - 1$ and $0 \leq b \leq p - 1$, define

$$h_{\bar{a}, b}(x) = \left(\sum_{i=0}^{k-1} a_i x_i + b \right) \text{ mod } p.$$

Then

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

PROOF

PROOF

Assuming $x \neq y$ and they differ on the position i ($x_i \neq y_i$).

PROOF

Assuming $x \neq y$ and they differ on the position i ($x_i \neq y_i$).

For every $u, v \in \{0, 1, \dots, p-1\}$, we have equations

$$\begin{cases} a_i x_i + b = \left(u - \sum_{j \neq i} a_j x_j \right) \pmod{p} \\ a_i y_i + b = \left(v - \sum_{j \neq i} a_j y_j \right) \pmod{p} \end{cases}$$

PROOF

Assuming $x \neq y$ and they differ on the position i ($x_i \neq y_i$).

For every $u, v \in \{0, 1, \dots, p-1\}$, we have equations

$$\begin{cases} a_i x_i + b = \left(u - \sum_{j \neq i} a_j x_j \right) \pmod{p} \\ a_i y_i + b = \left(v - \sum_{j \neq i} a_j y_j \right) \pmod{p} \end{cases}$$

For fixed x, y, u, v and $\{a_j\}_{j \neq i}$, a unique pair (a_i, b) (out of p^2 pairs) is determined.

Therefore,

$$\Pr_{h_{\bar{a}, b} \in \mathcal{H}} [h_{\bar{a}, b}(x) = u \wedge h_{\bar{a}, b}(y) = v] = \frac{1}{p^2}.$$

COUNTING DISTINCT ELEMENTS

COUNTING DISTINCT ELEMENTS

Back to the streaming model, we are given a sequence of numbers $\langle a_1, \dots, a_m \rangle$ where each $a_i \in [n]$.

COUNTING DISTINCT ELEMENTS

Back to the streaming model, we are given a sequence of numbers $\langle a_1, \dots, a_m \rangle$ where each $a_i \in [n]$.

It defines a **frequency vector** $\mathbf{f} = (f_1, \dots, f_n)$ where $f_i = |\{k \in [m] : a_k = i\}|$.

COUNTING DISTINCT ELEMENTS

Back to the streaming model, we are given a sequence of numbers $\langle a_1, \dots, a_m \rangle$ where each $a_i \in [n]$.

It defines a **frequency vector** $\mathbf{f} = (f_1, \dots, f_n)$ where $f_i = |\{k \in [m] : a_k = i\}|$.

We want to compute the number $d = |\{i \in [n] : f_i > 0\}|$.

COUNTING DISTINCT ELEMENTS

Back to the streaming model, we are given a sequence of numbers $\langle a_1, \dots, a_m \rangle$ where each $a_i \in [n]$.

It defines a **frequency vector** $\mathbf{f} = (f_1, \dots, f_n)$ where $f_i = |\{k \in [m] : a_k = i\}|$.

We want to compute the number $d = |\{i \in [n] : f_i > 0\}|$.

The value d is the number of distinct elements in the stream.

THE AMS ALGORITHM

The algorithm is named after Alon, Matias and Szegedy.

THE AMS ALGORITHM

The algorithm is named after Alon, Matias and Szegedy.

For every integer $p > 0$, we use $\text{zero}(p)$ to denote **number of trailing zeros** of p in binary.

THE AMS ALGORITHM

The algorithm is named after Alon, Matias and Szegedy.

For every integer $p > 0$, we use $\text{zero}(p)$ to denote **number of trailing zeros** of p in binary.

$$\text{zero}(p) \triangleq \max \{i : 2^i \text{ divides } p\}.$$

Algorithm AMS Algorithm for Counting Distinct Elements

Init:

A random Hash function $h : [n] \rightarrow [n]$ from a 2-universal family

$Z \leftarrow 0$

On Input y :

if $\text{zero}(h(y)) > Z$ **then**

$Z \leftarrow \text{zero}(h(y))$

end if

Output:

$\hat{d} = 2^{Z + \frac{1}{2}}$.

INTUITION

INTUITION

After applying the Hash function, $h(y)$ is uniform is $[n]$.

INTUITION

After applying the Hash function, $h(y)$ is uniform is $[n]$.

The probability that it has more than t trailing zeros is **at most 2^{-t}**

INTUITION

After applying the Hash function, $h(y)$ is uniform is $[n]$.

The probability that it has more than t trailing zeros is **at most 2^{-t}** .

Therefore, at least in expectation, if we have d distinct numbers, one of them may have $\log_2 d$ trailing zeros.

INTUITION

After applying the Hash function, $h(y)$ is uniform is $[n]$.

The probability that it has more than t trailing zeros is **at most 2^{-t}** .

Therefore, at least in expectation, if we have d distinct numbers, one of them may have $\log_2 d$ trailing zeros.

We now turn this intuition into a rigorous proof.

For every $0 \leq r \leq n$, we use a random variable Y_r to denote the number of $h(a_j)$ with trailing zero at least r .

For every $0 \leq r \leq n$, we use a random variable Y_r to denote **the number of $h(a_j)$ with trailing zero at least r .**

The sequence of variables $\{Y_r\}_{0 \leq r \leq n}$ determines the variable Z since $Z = \max_r \{Y_r > 0\}$.

For every $0 \leq r \leq n$, we use a random variable Y_r to denote **the number of $h(a_j)$ with trailing zero at least r .**

The sequence of variables $\{Y_r\}_{0 \leq r \leq n}$ determines the variable Z since $Z = \max_r \{Y_r > 0\}$.

This motivates us to understand the behavior of Y_r .

For every $0 \leq r \leq n$, we use a random variable Y_r to denote **the number of $h(a_j)$ with trailing zero at least r .**

The sequence of variables $\{Y_r\}_{0 \leq r \leq n}$ determines the variable Z since $Z = \max_r \{Y_r > 0\}$.

This motivates us to understand the behavior of Y_r .

For every $k \in [n]$, we denote $X_{k,r}$ as the indicator that $h(k)$ has at least r trailing zeros, then $Y_r = \sum_{k \in [n]: f_k > 0} X_{k,r}$.

For every $0 \leq r \leq n$, we use a random variable Y_r to denote **the number of $h(a_j)$ with trailing zero at least r .**

The sequence of variables $\{Y_r\}_{0 \leq r \leq n}$ determines the variable Z since $Z = \max_r \{Y_r > 0\}$.

This motivates us to understand the behavior of Y_r .

For every $k \in [n]$, we denote $X_{k,r}$ as the indicator that $h(k)$ has at least r trailing zeros, then $Y_r = \sum_{k \in [n]: f_k > 0} X_{k,r}$.

Using this decomposition, it is not hard to see that $\mathbf{E}[Y_r] = \frac{d}{2^r}$ and $\mathbf{Var}[Y_r] \leq \frac{d}{2^r}$.

Applying Markov's inequality, we obtain

$$\Pr [Y_r > 0] = \Pr [Y_r \geq 1] \leq \mathbf{E} [Y_r] = \frac{d}{2^r}.$$

Applying Markov's inequality, we obtain

$$\Pr [Y_r > 0] = \Pr [Y_r \geq 1] \leq \mathbf{E} [Y_r] = \frac{d}{2^r}.$$

Applying Chebyshev's inequality, we obtain

$$\Pr [Y_r = 0] \leq \Pr \left[|Y_r - \mathbf{E} [Y_r]| \geq \frac{d}{2^r} \right] \leq \frac{2^r}{d}.$$

Applying Markov's inequality, we obtain

$$\Pr [Y_r > 0] = \Pr [Y_r \geq 1] \leq \mathbf{E} [Y_r] = \frac{d}{2^r}.$$

Applying Chebyshev's inequality, we obtain

$$\Pr [Y_r = 0] \leq \Pr \left[|Y_r - \mathbf{E} [Y_r]| \geq \frac{d}{2^r} \right] \leq \frac{2^r}{d}.$$

We know that $Y_r > 0$ for all $r \leq Z$ and $Y_r = 0$ for all $r > Z$.

Applying Markov's inequality, we obtain

$$\Pr [Y_r > 0] = \Pr [Y_r \geq 1] \leq \mathbf{E} [Y_r] = \frac{d}{2^r}.$$

Applying Chebyshev's inequality, we obtain

$$\Pr [Y_r = 0] \leq \Pr \left[|Y_r - \mathbf{E} [Y_r]| \geq \frac{d}{2^r} \right] \leq \frac{2^r}{d}.$$

We know that $Y_r > 0$ for all $r \leq Z$ and $Y_r = 0$ for all $r > Z$.

Therefore, Z cannot be too far from $\log_2 d$:

Applying Markov's inequality, we obtain

$$\Pr [Y_r > 0] = \Pr [Y_r \geq 1] \leq \mathbf{E} [Y_r] = \frac{d}{2^r}.$$

Applying Chebyshev's inequality, we obtain

$$\Pr [Y_r = 0] \leq \Pr \left[|Y_r - \mathbf{E} [Y_r]| \geq \frac{d}{2^r} \right] \leq \frac{2^r}{d}.$$

We know that $Y_r > 0$ for all $r \leq Z$ and $Y_r = 0$ for all $r > Z$.

Therefore, Z cannot be too far from $\log_2 d$:

- ▶ if $Z \ll \log_2 d$, we can find a small r with $Y_r = 0$, which happens with small probability;

Applying Markov's inequality, we obtain

$$\Pr [Y_r > 0] = \Pr [Y_r \geq 1] \leq \mathbf{E} [Y_r] = \frac{d}{2^r}.$$

Applying Chebyshev's inequality, we obtain

$$\Pr [Y_r = 0] \leq \Pr \left[|Y_r - \mathbf{E} [Y_r]| \geq \frac{d}{2^r} \right] \leq \frac{2^r}{d}.$$

We know that $Y_r > 0$ for all $r \leq Z$ and $Y_r = 0$ for all $r > Z$.

Therefore, Z cannot be too far from $\log_2 d$:

- ▶ if $Z \ll \log_2 d$, we can find a small r with $Y_r = 0$, which happens with small probability;
- ▶ if $Z \gg \log_2 d$, we can find a big r with $Y_r > 0$, which happens with small probability.

If $\widehat{d} \leq \frac{d}{3}$, let r be the **largest integer** with $2^{r+\frac{1}{2}} \leq \frac{d}{3}$.

If $\widehat{d} \leq \frac{d}{3}$, let r be the **largest integer** with $2^{r+\frac{1}{2}} \leq \frac{d}{3}$.

$$\Pr \left[\widehat{d} \leq \frac{d}{3} \right] = \Pr [Z \leq r] = \Pr [Y_{r+1} = 0] \leq \frac{2^{r+1}}{d} \leq \frac{\sqrt{2}}{3}.$$

If $\widehat{d} \leq \frac{d}{3}$, let r be the **largest integer** with $2^{r+\frac{1}{2}} \leq \frac{d}{3}$.

$$\Pr \left[\widehat{d} \leq \frac{d}{3} \right] = \Pr [Z \leq r] = \Pr [Y_{r+1} = 0] \leq \frac{2^{r+1}}{d} \leq \frac{\sqrt{2}}{3}.$$

If $\widehat{d} \geq 3d$, let r be the **smallest integer** with $2^{r+\frac{1}{2}} \geq 3d$.

If $\widehat{d} \leq \frac{d}{3}$, let r be the **largest integer** with $2^{r+\frac{1}{2}} \leq \frac{d}{3}$.

$$\Pr \left[\widehat{d} \leq \frac{d}{3} \right] = \Pr [Z \leq r] = \Pr [Y_{r+1} = 0] \leq \frac{2^{r+1}}{d} \leq \frac{\sqrt{2}}{3}.$$

If $\widehat{d} \geq 3d$, let r be the **smallest integer** with $2^{r+\frac{1}{2}} \geq 3d$.

$$\Pr \left[\widehat{d} \geq 3d \right] = \Pr [Z \geq r] = \Pr [Y_r > 0] \leq \frac{d}{2^r} \leq \frac{\sqrt{2}}{3}.$$

If $\widehat{d} \leq \frac{d}{3}$, let r be the **largest integer** with $2^{r+\frac{1}{2}} \leq \frac{d}{3}$.

$$\Pr \left[\widehat{d} \leq \frac{d}{3} \right] = \Pr [Z \leq r] = \Pr [Y_{r+1} = 0] \leq \frac{2^{r+1}}{d} \leq \frac{\sqrt{2}}{3}.$$

If $\widehat{d} \geq 3d$, let r be the **smallest integer** with $2^{r+\frac{1}{2}} \geq 3d$.

$$\Pr [\widehat{d} \geq 3d] = \Pr [Z \geq r] = \Pr [Y_r > 0] \leq \frac{d}{2^r} \leq \frac{\sqrt{2}}{3}.$$

The algorithm costs $O(\log n)$ bits of memory.

MEDIAN

MEDIAN

We can apply the standard **Median trick** to the AMS algorithm. (**Exercise**)

MEDIAN

We can apply the standard **Median trick** to the AMS algorithm. (**Exercise**)

Using $O(\log \frac{1}{\delta} \log n)$ bits of memory, we can obtain

$$\Pr \left[\frac{d}{3} \leq \widehat{d} \leq 3d \right] \geq 1 - \delta.$$

THE BJKST ALGORITHM

The following improvement of AMS is due to Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan.

Algorithm BJKST Algorithm for Counting Distinct Elements

Init: Random Hash functions $h : [n] \rightarrow [n]$, $g : [n] \rightarrow [b\epsilon^{-4} \log^2 n]$, both from 2-universal families; $Z \leftarrow 0$, $B \leftarrow \emptyset$

On Input y :

if $\text{zero}(h(y)) \geq Z$ **then**

$B \leftarrow B \cup \{(g(y), \text{zeros}(h(y)))\}$

while $|B| \geq c/\epsilon^2$ **do**

$Z \leftarrow Z + 1$

 Remove all (α, β) with $\beta < Z$ from B

end while

end if

Output: $\hat{d} = |B| 2^Z$