

Distinct Elements

- We are back to our basic streaming model:
 $i_1, \dots, i_n \in [d] = \{1, \dots, d\}$.
- The frequency vector $x \in \mathbb{Z}^d$: $x_j = |\{t : i_t = j\}|$.

Distinct Elements

- We are back to our basic streaming model:
 $i_1, \dots, i_n \in [d] = \{1, \dots, d\}$.
- The frequency vector $x \in \mathbb{Z}^d$: $x_j = |\{t : i_t = j\}|$.
- *Counting distinct elements*: estimate $\|x\|_0 := |j : x_j > 0|$ up to $(1 + \epsilon)$ -factor approximation.

Distinct Elements

- We are back to our basic streaming model:
 $i_1, \dots, i_n \in [d] = \{1, \dots, d\}$.
- The frequency vector $x \in \mathbb{Z}^d$: $x_j = |\{t : i_t = j\}|$.
- *Counting distinct elements*: estimate $\|x\|_0 := |j : x_j > 0|$ up to $(1 + \epsilon)$ -factor approximation.
- Again, we must use space $O(\log d, \frac{1}{\epsilon})$.

An Ideal Algorithm

- If we can make the distribution of $\{i_t\}$ uniform, then it is easier to estimate its size.

An Ideal Algorithm

- If we can make the distribution of $\{i_t\}$ uniform, then it is easier to estimate its size.
- E.g., if we can have an ideal hash function h from $[d]$ to $[0, 1]$, let X be the minimum $h(i_t)$, then $\frac{1}{X}$ seems a reasonable estimate.

An Ideal Algorithm

- If we can make the distribution of $\{i_t\}$ uniform, then it is easier to estimate its size.
- E.g., if we can have an ideal hash function h from $[d]$ to $[0, 1]$, let X be the minimum $h(i_t)$, then $\frac{1}{X}$ seems a reasonable estimate.
- Indeed, suppose we have i.i.d. X_1, \dots, X_ℓ uniformly distributed on $[0, 1]$, let the smallest be $X_{(1)}$.

An Ideal Algorithm

- If we can make the distribution of $\{i_t\}$ uniform, then it is easier to estimate its size.
- E.g., if we can have an ideal hash function h from $[d]$ to $[0, 1]$, let X be the minimum $h(i_t)$, then $\frac{1}{X}$ seems a reasonable estimate.
- Indeed, suppose we have i.i.d. X_1, \dots, X_ℓ uniformly distributed on $[0, 1]$, let the smallest be $X_{(1)}$.
 - $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(\ell)}$ are called *order statistics*.

An Ideal Algorithm

- If we can make the distribution of $\{i_t\}$ uniform, then it is easier to estimate its size.
- E.g., if we can have an ideal hash function h from $[d]$ to $[0, 1]$, let X be the minimum $h(i_t)$, then $\frac{1}{X}$ seems a reasonable estimate.
- Indeed, suppose we have i.i.d. X_1, \dots, X_ℓ uniformly distributed on $[0, 1]$, let the smallest be $X_{(1)}$.
 - $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(\ell)}$ are called *order statistics*.
 - The distribution of $X_{(1)}$ is a so-called *Beta distribution* $B(1, \ell)$. We have $\mathbf{E}[X_{(1)}] = \frac{1}{\ell+1}$.

An Ideal Algorithm

- If we can make the distribution of $\{i_t\}$ uniform, then it is easier to estimate its size.
- E.g., if we can have an ideal hash function h from $[d]$ to $[0, 1]$, let X be the minimum $h(i_t)$, then $\frac{1}{X}$ seems a reasonable estimate.
- Indeed, suppose we have i.i.d. X_1, \dots, X_ℓ uniformly distributed on $[0, 1]$, let the smallest be $X_{(1)}$.
 - $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(\ell)}$ are called *order statistics*.
 - The distribution of $X_{(1)}$ is a so-called *Beta distribution* $B(1, \ell)$. We have
$$\mathbf{E}[X_{(1)}] = \frac{1}{\ell+1}.$$
- Therefore, $\frac{1}{X} - 1$ is an unbiased estimator of $\|x\|_0$.

An Ideal Algorithm

- If we can make the distribution of $\{i_t\}$ uniform, then it is easier to estimate its size.
- E.g., if we can have an ideal hash function h from $[d]$ to $[0, 1]$, let X be the minimum $h(i_t)$, then $\frac{1}{X}$ seems a reasonable estimate.
- Indeed, suppose we have i.i.d. X_1, \dots, X_ℓ uniformly distributed on $[0, 1]$, let the smallest be $X_{(1)}$.
 - $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(\ell)}$ are called *order statistics*.
 - The distribution of $X_{(1)}$ is a so-called *Beta distribution* $B(1, \ell)$. We have

$$\mathbf{E}[X_{(1)}] = \frac{1}{\ell+1}.$$
- Therefore, $\frac{1}{X} - 1$ is an unbiased estimator of $\|x\|_0$.
- $\text{Var}[X_{(1)}] = \frac{\ell}{(\ell+1)^2(\ell+2)} \leq \frac{1}{(\ell+1)^2}.$

An Ideal Algorithm

- If we can make the distribution of $\{i_t\}$ uniform, then it is easier to estimate its size.
- E.g., if we can have an ideal hash function h from $[d]$ to $[0, 1]$, let X be the minimum $h(i_t)$, then $\frac{1}{X}$ seems a reasonable estimate.
- Indeed, suppose we have i.i.d. X_1, \dots, X_ℓ uniformly distributed on $[0, 1]$, let the smallest be $X_{(1)}$.
 - $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(\ell)}$ are called *order statistics*.
 - The distribution of $X_{(1)}$ is a so-called *Beta distribution* $B(1, \ell)$. We have

$$\mathbf{E}[X_{(1)}] = \frac{1}{\ell+1}.$$
- Therefore, $\frac{1}{X} - 1$ is an unbiased estimator of $\|x\|_0$.
- $\text{Var}[X_{(1)}] = \frac{\ell}{(\ell+1)^2(\ell+2)} \leq \frac{1}{(\ell+1)^2}$.
- We can apply the Chebyshev bound, although the variance is a bit too large for our purpose.

An Algorithm Assuming Ideal Hash

- Maintain k independent, ideal hash functions.

An Algorithm Assuming Ideal Hash

- Maintain k independent, ideal hash functions.
- For each hash function h_i , store $Z_i = \min_t \{h_i(i_t)\}$, the smallest address used throughout the stream.

An Algorithm Assuming Ideal Hash

- Maintain k independent, ideal hash functions.
- For each hash function h_i , store $Z_i = \min_t \{h_i(i_t)\}$, the smallest address used throughout the stream.
- Take the median of Z_1, \dots, Z_k . Use it to estimate $\|x\|_0$. (See reading material for details.)
- This algorithm assumes we have access to ideal hash functions.

An Algorithm Assuming Ideal Hash

- Maintain k independent, ideal hash functions.
- For each hash function h_i , store $Z_i = \min_t \{h_i(i_t)\}$, the smallest address used throughout the stream.
- Take the median of Z_1, \dots, Z_k . Use it to estimate $\|x\|_0$. (See reading material for details.)
- This algorithm assumes we have access to ideal hash functions.
- Ideas for improvement:
 - Use real hash functions. Discretize the range. Possibly use k -wise independent hash family for appropriate k .

An Algorithm Assuming Ideal Hash

- Maintain k independent, ideal hash functions.
- For each hash function h_i , store $Z_i = \min_t \{h_i(i_t)\}$, the smallest address used throughout the stream.
- Take the median of Z_1, \dots, Z_k . Use it to estimate $\|x\|_0$. (See reading material for details.)
- This algorithm assumes we have access to ideal hash functions.
- Ideas for improvement:
 - Use real hash functions. Discretize the range. Possibly use k -wise independent hash family for appropriate k .
 - The minimum of $h(i_t)$ tends to be volatile: a single bad event ruins the estimate.

An Algorithm Assuming Ideal Hash

- Maintain k independent, ideal hash functions.
- For each hash function h_i , store $Z_i = \min_t \{h_i(i_t)\}$, the smallest address used throughout the stream.
- Take the median of Z_1, \dots, Z_k . Use it to estimate $\|x\|_0$. (See reading material for details.)
- This algorithm assumes we have access to ideal hash functions.
- Ideas for improvement:
 - Use real hash functions. Discretize the range. Possibly use k -wise independent hash family for appropriate k .
 - The minimum of $h(i_t)$ tends to be volatile: a single bad event ruins the estimate.
 - To make the estimate more stable, we may keep track of more than one smallest hash values.

KMV

The following KMV (k minimum values) algorithm is due to Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan (2002).

KMV

The following KMV (k minimum values) algorithm is due to Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan (2002).

- Sample a hash function h from a pairwise independent hash family mapping $[d]$ to $[D]$, for $D \in [d^3, 2d^3]$ that is a power of 2.

KMV

The following KMV (k minimum values) algorithm is due to Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan (2002).

- Sample a hash function h from a pairwise independent hash family mapping $[d]$ to $[D]$, for $D \in [d^3, 2d^3]$ that is a power of 2.
- Initialize S to \emptyset . Set $t = 12/\delta\epsilon^2$.

KMV

The following KMV (k minimum values) algorithm is due to Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan (2002).

- Sample a hash function h from a pairwise independent hash family mapping $[d]$ to $[D]$, for $D \in [d^3, 2d^3]$ that is a power of 2.
- Initialize S to \emptyset . Set $t = 12/\delta\epsilon^2$.
- When i_j arrives,
 - If $|S| < t$, then add $h(i_j)$ to S ;
 - Otherwise, only if $h(i_j) < y, \forall y \in S$, add $h(i_j)$ to S and remove the largest element of S .

KMV

The following KMV (k minimum values) algorithm is due to Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan (2002).

- Sample a hash function h from a pairwise independent hash family mapping $[d]$ to $[D]$, for $D \in [d^3, 2d^3]$ that is a power of 2.
- Initialize S to \emptyset . Set $t = 12/\delta\epsilon^2$.
- When i_j arrives,
 - If $|S| < t$, then add $h(i_j)$ to S ;
 - Otherwise, only if $h(i_j) < y, \forall y \in S$, add $h(i_j)$ to S and remove the largest element of S .
- For output at the end:
 - If $|S| < t$, return $|S|$.

KMV

The following KMV (k minimum values) algorithm is due to Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan (2002).

- Sample a hash function h from a pairwise independent hash family mapping $[d]$ to $[D]$, for $D \in [d^3, 2d^3]$ that is a power of 2.
- Initialize S to \emptyset . Set $t = 12/\delta\epsilon^2$.
- When i_j arrives,
 - If $|S| < t$, then add $h(i_j)$ to S ;
 - Otherwise, only if $h(i_j) < y, \forall y \in S$, add $h(i_j)$ to S and remove the largest element of S .
- For output at the end:
 - If $|S| < t$, return $|S|$.
 - Otherwise, let X be the largest element in S , return $\frac{tD}{X}$.

Analysis of KMV

Proposition

If Y is a Bernoulli random variable, then $\text{Var}[Y] \leq \mathbf{Pr}[Y = 1]$.

Analysis of KMV

Proposition

If Y is a Bernoulli random variable, then $\text{Var}[Y] \leq \mathbf{Pr}[Y = 1]$.

Proof.

$$\text{Var}[Y] = \mathbf{E}[Y^2] - \mathbf{E}[Y]^2 \leq \mathbf{E}[Y^2] = \mathbf{Pr}[Y = 1]. \quad \square$$

Analysis of KMV

Proposition

If Y is a Bernoulli random variable, then $\text{Var}[Y] \leq \mathbf{Pr}[Y = 1]$.

Proof.

$$\text{Var}[Y] = \mathbf{E}[Y^2] - \mathbf{E}[Y]^2 \leq \mathbf{E}[Y^2] = \mathbf{Pr}[Y = 1]. \quad \square$$

- Let's denote $\ell := \|x\|_0$, assume $\epsilon < \frac{1}{2}$, and $d > \frac{2}{\epsilon^2 \delta}$.
- First case of output: if $|S| < t$, what's the chance that $\|x\|_0 > |S|$?

Analysis of KMV

Proposition

If Y is a Bernoulli random variable, then $\text{Var}[Y] \leq \mathbf{Pr}[Y = 1]$.

Proof.

$$\text{Var}[Y] = \mathbf{E}[Y^2] - \mathbf{E}[Y]^2 \leq \mathbf{E}[Y^2] = \mathbf{Pr}[Y = 1]. \quad \square$$

- Let's denote $\ell := \|x\|_0$, assume $\epsilon < \frac{1}{2}$, and $d > \frac{2}{\epsilon^2 \delta}$.
- First case of output: if $|S| < t$, what's the chance that $\|x\|_0 > |S|$?
 - For any pair of indices, they are mapped to the same address with probability $\frac{1}{D}$.

Analysis of KMV

Proposition

If Y is a Bernoulli random variable, then $\text{Var}[Y] \leq \mathbf{Pr}[Y = 1]$.

Proof.

$$\text{Var}[Y] = \mathbf{E}[Y^2] - \mathbf{E}[Y]^2 \leq \mathbf{E}[Y^2] = \mathbf{Pr}[Y = 1]. \quad \square$$

- Let's denote $\ell := \|x\|_0$, assume $\epsilon < \frac{1}{2}$, and $d > \frac{2}{\epsilon^2 \delta}$.
- First case of output: if $|S| < t$, what's the chance that $\|x\|_0 > |S|$?
 - For any pair of indices, they are mapped to the same address with probability $\frac{1}{D}$.
 - There are $\binom{\ell}{2}$ pairs, so the probability that any clash happens is $\leq \binom{\ell}{2} \cdot \frac{1}{D} \leq \frac{1}{d}$. (Recall $D \geq d^3$.)

Analysis of KMV

Proposition

If Y is a Bernoulli random variable, then $\text{Var}[Y] \leq \mathbf{Pr}[Y = 1]$.

Proof.

$$\text{Var}[Y] = \mathbf{E}[Y^2] - \mathbf{E}[Y]^2 \leq \mathbf{E}[Y^2] = \mathbf{Pr}[Y = 1]. \quad \square$$

- Let's denote $\ell := \|x\|_0$, assume $\epsilon < \frac{1}{2}$, and $d > \frac{2}{\epsilon^2 \delta}$.
- First case of output: if $|S| < t$, what's the chance that $\|x\|_0 > |S|$?
 - For any pair of indices, they are mapped to the same address with probability $\frac{1}{D}$.
 - There are $\binom{\ell}{2}$ pairs, so the probability that any clash happens is $\leq \binom{\ell}{2} \cdot \frac{1}{D} \leq \frac{1}{d}$. (Recall $D \geq d^3$.)
 - So the output is exactly correct w.p. $1 - \frac{1}{d}$.

Analysis of KMV

Proposition

If Y is a Bernoulli random variable, then $\text{Var}[Y] \leq \mathbf{Pr}[Y = 1]$.

Proof.

$$\text{Var}[Y] = \mathbf{E}[Y^2] - \mathbf{E}[Y]^2 \leq \mathbf{E}[Y^2] = \mathbf{Pr}[Y = 1]. \quad \square$$

- Let's denote $\ell := \|x\|_0$, assume $\epsilon < \frac{1}{2}$, and $d > \frac{2}{\epsilon^2 \delta}$.
- First case of output: if $|S| < t$, what's the chance that $\|x\|_0 > |S|$?
 - For any pair of indices, they are mapped to the same address with probability $\frac{1}{D}$.
 - There are $\binom{\ell}{2}$ pairs, so the probability that any clash happens is $\leq \binom{\ell}{2} \cdot \frac{1}{D} \leq \frac{1}{d}$. (Recall $D \geq d^3$.)
 - So the output is exactly correct w.p. $1 - \frac{1}{d}$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.
- Recall: X is the largest element in S . We'll bound $\Pr\left[\left|\frac{tD}{X} - \ell\right| > \epsilon\ell\right]$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.
- Recall: X is the largest element in S . We'll bound $\Pr[|\frac{tD}{X} - \ell| > \epsilon\ell]$.
- Consider the event $\frac{tD}{X} > (1 + \epsilon)\ell$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.
- Recall: X is the largest element in S . We'll bound $\Pr[|\frac{tD}{X} - \ell| > \epsilon\ell]$.
- Consider the event $\frac{tD}{X} > (1 + \epsilon)\ell$.
- This happens only if more than t of the ℓ elements are hashed to addresses smaller than $X < \frac{tD}{(1+\epsilon)\ell} \leq \frac{(1-\epsilon/2)tD}{\ell}$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.
- Recall: X is the largest element in S . We'll bound $\Pr\left[\left|\frac{tD}{X} - \ell\right| > \epsilon\ell\right]$.
- Consider the event $\frac{tD}{X} > (1 + \epsilon)\ell$.
- This happens only if more than t of the ℓ elements are hashed to addresses smaller than $X < \frac{tD}{(1+\epsilon)\ell} \leq \frac{(1-\epsilon/2)tD}{\ell}$.
- W.l.o.g let the ℓ elements be $1, \dots, \ell$, and let Z_i be the indicator variable for the event $h(i) < \frac{(1-\epsilon/2)tD}{\ell}$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.
- Recall: X is the largest element in S . We'll bound $\Pr\left[\left|\frac{tD}{X} - \ell\right| > \epsilon\ell\right]$.
- Consider the event $\frac{tD}{X} > (1 + \epsilon)\ell$.
- This happens only if more than t of the ℓ elements are hashed to addresses smaller than $X < \frac{tD}{(1+\epsilon)\ell} \leq \frac{(1-\epsilon/2)tD}{\ell}$.
- W.l.o.g let the ℓ elements be $1, \dots, \ell$, and let Z_i be the indicator variable for the event $h(i) < \frac{(1-\epsilon/2)tD}{\ell}$.
- Then $\mathbf{E}[Z_i] = (1 - \epsilon/2)t/\ell$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.
- Recall: X is the largest element in S . We'll bound $\Pr\left[\left|\frac{tD}{X} - \ell\right| > \epsilon\ell\right]$.
- Consider the event $\frac{tD}{X} > (1 + \epsilon)\ell$.
- This happens only if more than t of the ℓ elements are hashed to addresses smaller than $X < \frac{tD}{(1+\epsilon)\ell} \leq \frac{(1-\epsilon/2)tD}{\ell}$.
- W.l.o.g let the ℓ elements be $1, \dots, \ell$, and let Z_i be the indicator variable for the event $h(i) < \frac{(1-\epsilon/2)tD}{\ell}$.
- Then $\mathbf{E}[Z_i] = (1 - \epsilon/2)t/\ell$.
- The bad event is $Z := \sum_{i=1}^{\ell} Z_i \geq t$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.
- Recall: X is the largest element in S . We'll bound $\Pr[|\frac{tD}{X} - \ell| > \epsilon\ell]$.
- Consider the event $\frac{tD}{X} > (1 + \epsilon)\ell$.
- This happens only if more than t of the ℓ elements are hashed to addresses smaller than $X < \frac{tD}{(1+\epsilon)\ell} \leq \frac{(1-\epsilon/2)tD}{\ell}$.
- W.l.o.g let the ℓ elements be $1, \dots, \ell$, and let Z_i be the indicator variable for the event $h(i) < \frac{(1-\epsilon/2)tD}{\ell}$.
- Then $\mathbf{E}[Z_i] = (1 - \epsilon/2)t/\ell$.
- The bad event is $Z := \sum_{i=1}^{\ell} Z_i \geq t$.
- $\text{Var}[Z_i] \leq \Pr[Z_i] = (1 - \frac{\epsilon}{2})t/\ell$.

Analysis: The interesting case

- The interesting case: $|S| \geq t$.
- Recall: X is the largest element in S . We'll bound $\Pr\left[\left|\frac{tD}{X} - \ell\right| > \epsilon\ell\right]$.
- Consider the event $\frac{tD}{X} > (1 + \epsilon)\ell$.
- This happens only if more than t of the ℓ elements are hashed to addresses smaller than $X < \frac{tD}{(1+\epsilon)\ell} \leq \frac{(1-\epsilon/2)tD}{\ell}$.
- W.l.o.g let the ℓ elements be $1, \dots, \ell$, and let Z_i be the indicator variable for the event $h(i) < \frac{(1-\epsilon/2)tD}{\ell}$.
- Then $\mathbf{E}[Z_i] = (1 - \epsilon/2)t/\ell$.
- The bad event is $Z := \sum_{i=1}^{\ell} Z_i \geq t$.
- $\text{Var}[Z_i] \leq \Pr[Z_i] = (1 - \epsilon/2)t/\ell$.
- By pairwise independence we have $\text{Var}[Z] = \sum_i \text{Var}[Z_i] \leq t$.

Analysis of KMV (Cont.)

- We have so far $\{\frac{tD}{X} > (1 + \epsilon)\ell\} \Rightarrow \{Z \geq t\}$, $\mathbf{E}[Z] \leq (1 - \frac{\epsilon}{2})t$, and $\text{Var}[Z] \leq t$.

Analysis of KMV (Cont.)

- We have so far $\{\frac{tD}{X} > (1 + \epsilon)\ell\} \Rightarrow \{Z \geq t\}$, $\mathbf{E}[Z] \leq (1 - \frac{\epsilon}{2})t$, and $\text{Var}[Z] \leq t$.
- By Chebysheve inequality, we have

$$\mathbf{Pr} \left[\frac{tD}{X} > (1 + \epsilon)\ell \right] \leq \mathbf{Pr} [Z \geq t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{4}{\epsilon^2 t} \leq \frac{\delta}{3}.$$

Analysis of KMV (Cont.)

- We have so far $\{\frac{tD}{X} > (1 + \epsilon)\ell\} \Rightarrow \{Z \geq t\}$, $\mathbf{E}[Z] \leq (1 - \frac{\epsilon}{2})t$, and $\text{Var}[Z] \leq t$.
- By Chebysheve inequality, we have

$$\Pr \left[\frac{tD}{X} > (1 + \epsilon)\ell \right] \leq \Pr [Z \geq t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{4}{\epsilon^2 t} \leq \frac{\delta}{3}.$$

- Almost symmetrically, the event $\{\frac{tD}{X} < (1 - \epsilon)\ell\}$ happens only if fewer than t of the ℓ elements are hashed to addresses smaller than $X > \frac{tD}{(1-\epsilon)\ell}$.

Analysis of KMV (Cont.)

- We have so far $\{\frac{tD}{X} > (1 + \epsilon)\ell\} \Rightarrow \{Z \geq t\}$, $\mathbf{E}[Z] \leq (1 - \frac{\epsilon}{2})t$, and $\text{Var}[Z] \leq t$.
- By Chebysheve inequality, we have

$$\Pr \left[\frac{tD}{X} > (1 + \epsilon)\ell \right] \leq \Pr [Z \geq t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{4}{\epsilon^2 t} \leq \frac{\delta}{3}.$$

- Almost symmetrically, the event $\{\frac{tD}{X} < (1 - \epsilon)\ell\}$ happens only if fewer than t of the ℓ elements are hashed to addresses smaller than $X > \frac{tD}{(1-\epsilon)\ell}$.
- Let Z_i be the indicator variable for the event $h(i) < \frac{tD}{(1-\epsilon)\ell}$.

Analysis of KMV (Cont.)

- We have so far $\{ \frac{tD}{X} > (1 + \epsilon)\ell \} \Rightarrow \{ Z \geq t \}$, $\mathbf{E}[Z] \leq (1 - \frac{\epsilon}{2})t$, and $\text{Var}[Z] \leq t$.
- By Chebysheve inequality, we have

$$\Pr \left[\frac{tD}{X} > (1 + \epsilon)\ell \right] \leq \Pr [Z \geq t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{4}{\epsilon^2 t} \leq \frac{\delta}{3}.$$

- Almost symmetrically, the event $\{ \frac{tD}{X} < (1 - \epsilon)\ell \}$ happens only if fewer than t of the ℓ elements are hashed to addresses smaller than $X > \frac{tD}{(1-\epsilon)\ell}$.
- Let Z_i be the indicator variable for the event $h(i) < \frac{tD}{(1-\epsilon)\ell}$.

$$\frac{t}{(1 - \epsilon)\ell} \geq \mathbf{E}[Z_i] \geq \frac{t}{(1 - \epsilon)\ell} - \frac{1}{D} \geq \frac{(1 + \epsilon)t}{\ell} - \frac{1}{D} \geq \frac{(1 + \epsilon/2)t}{\ell}.$$

Analysis of KMV

$$\text{Var}[Z_i] \leq \mathbf{E}[Z_i] \leq \frac{t}{(1-\epsilon)\ell} \leq \frac{2t}{\ell}.$$

Analysis of KMV

$$\text{Var}[Z_i] \leq \mathbf{E}[Z_i] \leq \frac{t}{(1-\epsilon)\ell} \leq \frac{2t}{\ell}.$$

Let Z be $\sum_{i=1}^{\ell} Z_i$, then $\mathbf{E}[Z] \geq (1 + \frac{\epsilon}{2})t$, $\text{Var}[Z] \leq 2t$.

Analysis of KMV

$$\text{Var}[Z_i] \leq \mathbf{E}[Z_i] \leq \frac{t}{(1-\epsilon)\ell} \leq \frac{2t}{\ell}.$$

Let Z be $\sum_{i=1}^{\ell} Z_i$, then $\mathbf{E}[Z] \geq (1 + \frac{\epsilon}{2})t$, $\text{Var}[Z] \leq 2t$.

By Chebyshev inequality,

$$\Pr\left[\frac{tD}{X} < (1-\epsilon)\ell\right] \leq \Pr[Z < t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{8}{\epsilon^2 t} \leq \frac{2\delta}{3}.$$

Analysis of KMV

$$\text{Var}[Z_i] \leq \mathbf{E}[Z_i] \leq \frac{t}{(1-\epsilon)\ell} \leq \frac{2t}{\ell}.$$

Let Z be $\sum_{i=1}^{\ell} Z_i$, then $\mathbf{E}[Z] \geq (1 + \frac{\epsilon}{2})t$, $\text{Var}[Z] \leq 2t$.

By Chebyshev inequality,

$$\Pr\left[\frac{tD}{X} < (1-\epsilon)\ell\right] \leq \Pr[Z < t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{8}{\epsilon^2 t} \leq \frac{2\delta}{3}.$$

Combining everything, we have that with probability at least $1 - \delta$,

$$\left|\frac{tD}{X} - \ell\right| \leq \epsilon\ell.$$

Analysis of KMV

$$\text{Var}[Z_i] \leq \mathbf{E}[Z_i] \leq \frac{t}{(1-\epsilon)\ell} \leq \frac{2t}{\ell}.$$

Let Z be $\sum_{i=1}^{\ell} Z_i$, then $\mathbf{E}[Z] \geq (1 + \frac{\epsilon}{2})t$, $\text{Var}[Z] \leq 2t$.

By Chebyshev inequality,

$$\Pr\left[\frac{tD}{X} < (1-\epsilon)\ell\right] \leq \Pr[Z < t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{8}{\epsilon^2 t} \leq \frac{2\delta}{3}.$$

Combining everything, we have that with probability at least $1 - \delta$,

$$\left|\frac{tD}{X} - \ell\right| \leq \epsilon\ell.$$

Space usage:

- Storing the hash takes space $O(\log D) = O(\log d)$.

Analysis of KMV

$$\text{Var}[Z_i] \leq \mathbf{E}[Z_i] \leq \frac{t}{(1-\epsilon)\ell} \leq \frac{2t}{\ell}.$$

Let Z be $\sum_{i=1}^{\ell} Z_i$, then $\mathbf{E}[Z] \geq (1 + \frac{\epsilon}{2})t$, $\text{Var}[Z] \leq 2t$.

By Chebyshev inequality,

$$\Pr\left[\frac{tD}{X} < (1-\epsilon)\ell\right] \leq \Pr[Z < t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{8}{\epsilon^2 t} \leq \frac{2\delta}{3}.$$

Combining everything, we have that with probability at least $1 - \delta$,

$$\left|\frac{tD}{X} - \ell\right| \leq \epsilon\ell.$$

Space usage:

- Storing the hash takes space $O(\log D) = O(\log d)$.
- Storing S takes space $tO(\log D) = O\left(\frac{\log d}{\epsilon^2 \delta}\right)$.

Analysis of KMV

$$\text{Var}[Z_i] \leq \mathbf{E}[Z_i] \leq \frac{t}{(1-\epsilon)\ell} \leq \frac{2t}{\ell}.$$

Let Z be $\sum_{i=1}^{\ell} Z_i$, then $\mathbf{E}[Z] \geq (1 + \frac{\epsilon}{2})t$, $\text{Var}[Z] \leq 2t$.

By Chebyshev inequality,

$$\Pr\left[\frac{tD}{X} < (1-\epsilon)\ell\right] \leq \Pr[Z < t] \leq \frac{\text{Var}[Z]}{(\epsilon t/2)^2} \leq \frac{8}{\epsilon^2 t} \leq \frac{2\delta}{3}.$$

Combining everything, we have that with probability at least $1 - \delta$,

$$\left|\frac{tD}{X} - \ell\right| \leq \epsilon\ell.$$

Space usage:

- Storing the hash takes space $O(\log D) = O(\log d)$.
- Storing S takes space $tO(\log D) = O(\frac{\log d}{\epsilon^2 \delta})$.
- The optimal algorithm uses space $O(\log d + \epsilon^{-2})!$