

MTBF (Mean Time Between Failure)

System with N components C_1, \dots, C_N that at each time t may fail indep.

$$p = P[C_i \text{ fails at } t] \text{ so } P[\exists C_i \text{ that fails at time } t] = 1 - (1-p)^N$$

Let $X = \#$ time units before the next failure occurs. Due to the fact that each C_i fails independently, we have

$$\mathbb{E}[X] = \frac{1}{1 - (1-p)^N} \xrightarrow[N \rightarrow \infty]{} \frac{1}{p}$$

WORD COUNT WITH 1 ROUND OF MAPREDUCE.

Input: k documents $D_1, \dots, D_k \Rightarrow D_i = (\text{name}, \text{list of words})$

$D_1: \{\text{orange, apple, pear}\}$

Output $\{(w, c(w)) : w \text{ is a word in one of } D_i \text{'s doc, } c(w) = \# \text{ occurrences}\}$ set of all w followed by $c(w)$

Round 1 \rightarrow MAP : $\forall D_i$, separately

$$D_i \rightarrow \{(w, c_i(w)) : w \text{ is a word in } D_i, c_i(w) = \# \text{ occur. of } w \text{ in } D_i\}$$

\rightarrow REDUCE : $\forall w$ separately

$L_w = \text{list of } c_i(w)$

$$\text{REDUCE}(w, L_w) \rightarrow w, \sum_{i=1}^k c_i(w) \text{ with } c_i(w) \in L_w$$

ANALYSIS $R=1 \Rightarrow M_L \Rightarrow$ Suppose that $N_i = \#$ of words in D_i

$$N_{\max} = \max \{D_i, i \in (0, k)\}$$

MapPhase : $O(N_{\max})$ local space

ReducePhase : $O(k)$ local space $\Rightarrow M_L = O(\max \{N_{\max}, k\}) = O(\underline{N_{\max} + k})$

MA : $\Rightarrow N = \# \text{ total number of words}$

\times Space occupied by input pairs $O(N)$

" intermediate pairs $O(N)$

" output pairs $O(N)$

$$\Rightarrow M_A = O(N)$$

CLASS COUNT PROBLEM

PROBLEM: Set of N objects $\langle o_i, 1 \leq i \leq N \rangle$, labeled with γ_i . We want to count how much objects belong to each class γ .

R=1 CLASS COUNT 1-ROUND (NO PARTITIONING)

MapPhase: \forall input pairs do $(i, (o_i, \gamma_i)) \rightarrow (\gamma_i, 1)$

Reduce Phase: Let $L_\gamma =$ list of $(\gamma_i, 1)$ with $\gamma_i = \gamma$

$$(\gamma, L_\gamma) \rightarrow (\gamma, c(\gamma) = |L_\gamma|)$$

Analysis: $R=1 \Rightarrow M_L \Rightarrow$ RedPhase: if N objects $\langle o_i \rangle$ has label γ_i equal, we have

MapPhase: $O(N) \Rightarrow O(1) \Rightarrow O(N)$ doesn't satisfy design goal

CLASS COUNT, 2-ROUND WITH PARTITIONING (deterministic)

Input: $(i, (o_i, \gamma_i))$ with $i \in [0, N]$ l partitions

Round 1 $\forall i$ do $(i, (o_i, \gamma_i)) \rightarrow (\underbrace{i \bmod l}_{\text{take the rest}}, \gamma_i) \leftarrow \text{MAPPHASE}$

Let $L_j =$ List of elements labeled with γ_j , assigned to $i \bmod l = j$

$$(j, L_j) \rightarrow \{(\gamma, c(j, \gamma)) : \forall \text{ label } \gamma \text{ in } j, \text{ count occurrences}\} \leftarrow \text{REDUCE}$$

list of labels with respective occ. in L_j

Round 2 Empty $\leftarrow \text{MapPhase}$

$\forall \gamma$, let $L_\gamma =$ list of $c(j, \gamma)$'s from outputs of R1

$$(\gamma, L_\gamma) \rightarrow (\gamma, c(\gamma) = \sum_{c(j, \gamma) \in L_\gamma} c(j, \gamma))$$

Analysis:

$R=2 \Rightarrow R_1 \approx \text{MapPhase: } O(1)$

RedPhase: $O(N/l)$

$R_2 \approx \text{Map: } O(l)$

$$\Rightarrow M_L = \max(N/l, l) \Rightarrow O(\sqrt{N})$$

RedPhase: $O(l)$

CLASS-COUNT, 2-ROUND with RANDOM PARTITIONING **METHOD TWO** 22AUG

$l = \text{numb. of part.}$

Round 1: \forall input $(i, (0, y_i)) \rightarrow (x, y_i)$ with $x = \text{number} \in [0, l]$ with unif. prob.

↓ equal to R_2 , R_1 Red Phase of det. partitioning

Analysis: what change is only the MapPhase of R_1 ; we have that m_x is the number of interv. pairs with key x .

So RedPhase $\rightarrow O(m)$ where $m = \max\{m_x\}$

m is surely $\geq \frac{N}{l}$ since when partition N objects in l groups there must exists a group with $\frac{N}{l}$ objects \Rightarrow hope: $m \leq \frac{N}{l}$

Theorem: Fix $\ell = \sqrt{N}$ and suppose that in round 1 keys are assigned to intermediate pairs independently and with uniform prob. from 0 to \sqrt{N} . Then, with probability $\geq 1 - \frac{1}{N^5} \Rightarrow m = O(\sqrt{N})$

$$\downarrow (0, \sqrt{N})$$

$$M_L = O(\sqrt{N})$$

Proof:

We use 2 things: 1) \rightarrow Union Bound: $P[\bigcup_{i=1}^r E_i] \leq \sum_{i=1}^r P[E_i]$ where E_i are events

2) \rightarrow Chernoff Bound:

X_1, \dots, X_n Bernoulli variables s.t. $P[X_i=1] = p \quad \forall i \in \{1, \dots, n\}$

So $X = \sum X_i$ is \Rightarrow Binomial(n, p) with $E[X] = \mu = n \cdot p$

We have that $\forall \delta_1 \geq 0$ and $\delta_2 \in (0, 1)$:

$$a) P[X \geq \delta_1 \mu] \leq 2^{-\delta_1 \mu}$$

$$b) P[X \leq (1 - \delta_2) \mu] \leq 2^{-(\delta_2)^2 \mu / 2}$$

So the probability is bounded!

so, in Re Map Phase of W.C algorithm with random partition, we assign a key $x \in (0, \ell)$ with $m_x = \#$ number of times an intermediate pairs has key x

Define $y_i = \begin{cases} 1 & \text{if } (x_i, o_i) \rightarrow (x, y_i) \\ 0 & \text{otherwise} \end{cases}$

so y_i is a Bernoulli variable with $P[y_i=1] = \frac{1}{\sqrt{N}} = p$

so $m_x = \sum_{i=0}^{N-1} y_i \Rightarrow m_x \text{ is a Binomial}(N, \frac{1}{\sqrt{N}})$ \rightarrow since for hypo $\ell = \sqrt{N}$ and the probability of assignment is uniform $\frac{1}{\ell}$

$$\text{with } E[m_x] = N \cdot \frac{1}{\sqrt{N}} = \sqrt{N}$$

Using the Chernoff Bound (a) we have $P[m_x \geq 8\sqrt{N}] \leq 2^{-6\sqrt{N}}$

$$\text{and for } N \geq 1 \text{ and } \sqrt{N} \geq \log_2 N \Rightarrow 2^{-6\sqrt{N}} \leq 2^{-6\log_2 N} = \left(\left(\frac{1}{2}\right)^{\log_2 N}\right)^6 = \frac{1}{N^6}$$

$$\text{so } \boxed{P[m_x \geq 8\sqrt{N}] \leq \frac{1}{N^6}}$$

Now we want to extend the bound to $m = \max \{m_x : 0 \leq x \leq \sqrt{N}\}$

so we define $E_x = "m_x \geq 6\sqrt{N}"$ and so $P[m \geq 6\sqrt{N}] = P[E_0 \cup E_1 \cup \dots \cup E_{\sqrt{N}-1}]$

$$\leq \sqrt{N} \cdot P[E_x] \leq \sqrt{N} \cdot \frac{1}{N^6} < \frac{1}{N^5} \quad \text{so } P[m \geq 6\sqrt{N}] \leq \frac{1}{N^5} \text{ and}$$

by union bound
 $\sum_{i=0}^{\sqrt{N}-1} P[E_i]$

$$P[m \leq 6\sqrt{N}] \geq 1 - \frac{1}{N^5}$$

↓
So, with $P[\cdot] \geq 1 - \frac{1}{N^5}$ we have $m \leq 6\sqrt{N}$

and then for $m = O(\sqrt{N})$

CONCLUSION: If we distribute N objects at random among l partitions, and N/l is sufficient large, with HIGH PROB. no partition receives an excessive number of objects.

Theorem:

Let S be the set of centers returned by running FFT on P . Then:

$$\bigoplus_{\text{center}} (P, S) \leq 2 \bigoplus_{\text{center}} (P, k) \quad \rightsquigarrow \begin{array}{l} \text{so the output of FFT} \\ \text{is "close" to the optimal solution} \\ \text{by a factor of 2} \end{array}$$

So FFT is a 2-approx. algorithm.

Proof:

$S = \{c_1, c_2, \dots, c_n\}$ is the set of centers selected by FFT, } APPROX.
where $c_i = i\text{-th center selected}$ } CENTERS

Define $q = \text{point of } P \text{ farthest from } S$, $d(q, S) \geq d(x, S)$

It's important to do an observation: Using FFT, c_2 is the farthest point from c_1 ,
 c_3 is the farthest point from $\{c_2, c_1\} \Rightarrow d(c_2, c_1) \geq d(c_3, c_1) \geq d(c_3, \{c_1, c_2\})$
and so, if we iterate, $d(c_k, \{c_1, c_2, \dots, c_{k-1}\}) \geq d(q, S)$ $\xrightarrow{\min\{d(c_3, c_1), d(c_3, c_2)\}}$
logic since it would be chosen as a center

Now, consider the set $\{c_1, c_2, \dots, c_k, c_{k+1} = q\} = S \cup \{q\}$

If we prove that $d(q, S) \leq d(c_i, c_j) \quad \forall 1 \leq i < j \leq k+1$ the rest of the proof
is easy

Fix i, j arbitrarily $1 \leq i < j \leq k+1 \Rightarrow d(q, S) \leq d(q, \{c_1, \dots, c_{j-1}\})$ "less possibilities" and we
 $\leq d(c_j, \{c_1, \dots, c_{j-1}\})$ could get at most the
 $\leq d(c_j, c_i)$ result for S

\hookrightarrow since FFT takes
the point choosing
as the one at max
distance

\hookrightarrow since c_j
is selected as the max
distance and $i < j$, hence $c_i \in \{c_1, \dots, c_{j-1}\}$,
the $d(c_j, \{c_1, \dots, c_{j-1}\})$ is the min
distance between c_j and the set and
being c_i part of the set, the
equivalence holds

Let $S^* = \{c_1^*, c_2^*, \dots, c_k^*\}$ OPTIMAL CENTERS

$$\text{So } \phi_{\text{center}}^{\text{opt}}(P, k) = \phi_{\text{center}}^{\text{opt}}(P, S^*) \leq \phi_{\text{center}}^{\text{opt}}(P, S)$$

$$\bullet \forall x \in P, d(x, S^*) \leq \phi_{\text{center}}^{\text{opt}}(P, k) \Rightarrow \text{since } \phi_{\text{center}}^{\text{opt}}(P, k) = \max_{x \in P} d(x, S^*)$$

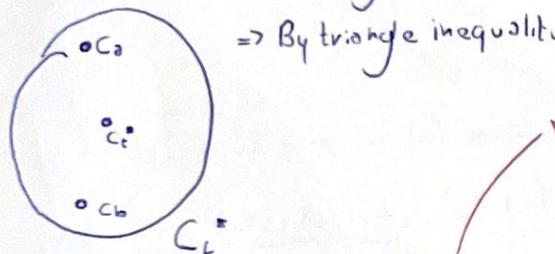
For $t=1..k$ let $C_t^* = \{x \in P : c_t^* \text{ is the center closest to } x\} \Rightarrow P = \bigcup_{t=1}^k C_t^*$

We have defined before that the set $\{c_1, c_2, \dots, c_q, c_{q+1}=q\} \subseteq P$

and for the PIGEONHOLE PRINCIPLE must $\exists 2$ points of $\{c_1, \dots, c_{q+1}\}$,

called c_a, c_b that belong to the same cluster C_t^*

$$\Rightarrow \text{By triangle inequality} \Rightarrow d(c_a, c_b) \leq d(c_a, c_t^*) + d(c_b, c_t^*) \leq 2\phi_{\text{center}}^{\text{opt}}(P, k)$$

$$= d(c_a, S^*)$$


so $\forall x \in P$ we have that $d(x, S) \leq d(q, S) \Rightarrow$ by the choice of q

$$\leq d(c_a, c_b) \Rightarrow \text{proof before}$$

$$\leq 2\phi_{\text{center}}^{\text{opt}}(P, k)$$

Lemma: Let T be the union of coresets T_i computed by MR-FFT on P_i .

$\forall x \in P$, we have $d(x, T) \leq 2 \cdot \underline{\phi_{\text{center}}^{\text{OPT}}(P, k)}$ \Rightarrow This implies that T is a good repres. of P , in the sense that each $x \in P$ has a "close-by" representative in T

Proof:

Recall that $P = P_1 \cup P_2 \cup \dots \cup P_l$ and $T = T_1 \cup T_2 \cup \dots \cup T_r$ where $T_i \leftarrow \text{FFT}(P_i, k)$

↓ So ∀ index $\in (1, l)$ we define $q_r = \text{dist} \text{ point of } P_r \text{ from } T_r \Rightarrow d(q_r, T_r) \geq d(x, T_r) \quad \forall x \in P_r$

CLAIM: $d(q_r, T_r) \leq 2 \phi_{\text{center}}^{\text{OPT}}(P, k) \quad \forall r$ \Rightarrow if it's true means \Rightarrow that it's valid $\forall x \in P$ $d(x, T) \leq d(x, T_r) \leq d(q_r, T_r) \leq 2 \phi_{\text{center}}^{\text{OPT}}(P, k)$

We need to prove the claim:

we fix r arbitrarily and we repeat the same argument in the analysis of FFT

$S^* = \{c_1^*, \dots, c_n^*\}$ OPTIMAL CENTERS

C_t^* = optimal clusters around C_t^* with $P = \bigcup_{i=1}^k C_t^*$

$T_r = \{c_1, \dots, c_n\}$ output of FFT and consider the set $\{c_1, \dots, c_n, c_{n+1} = q_r\}$ that are k_{r+1} points $\in P_r$ and so $\in P$

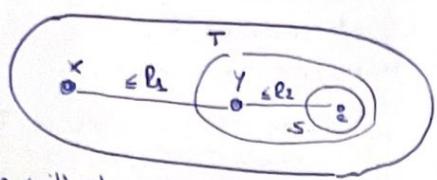
We have that $d(q_r, T_r) \leq d(c_i, c_j) \quad \forall 1 \leq i < j \leq k_{r+1}$ previous proof

•] 2 points c_a, c_b with $1 \leq a < b \leq k_{r+1}$ which \in to some C_t^* for some t

$$d(q_r, T_r) \leq d(c_a, c_b) \leq 2 \phi_{\text{center}}^{\text{OPT}}(P, k)$$

\downarrow
previous proof

Why Composable coreset technique work well with k-center and FPT



$x = \text{arbitrary point of } P$

$y = \text{closest point of } T \text{ to } x$

$c = \text{closest center of } S \text{ to } y$

$$\begin{aligned} \text{we will show that } R_1, R_2 &\leq 2\phi_{\text{k-center}}^{\text{OPT}}(P, k) \Rightarrow d(x, S) \leq d(x, c) \leq d(x, y) + d(y, c) \\ &\leq r_1 + r_2 \leq 4\phi_{\text{k-center}}^{\text{OPT}}(P, k) \end{aligned}$$

ANALYSIS OF MR-FFT

$R = 2$

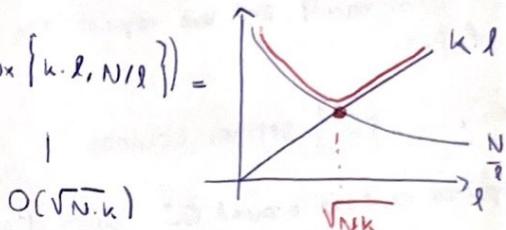
Round 1: MapPhase \Rightarrow 1 point $p \rightarrow (l, p)$ where l is the number of partitions $\Rightarrow O(1)$

ReducePhase $\Rightarrow O(N/l)$

Round 2: \propto MapPhase

$$\Rightarrow M_2 = O(\max \{k \cdot l, N/l\}) =$$

ReducePhase: $O(k \cdot l)$



$O(\sqrt{N \cdot k})$

Theorem: Let S be the set of k centers returned by running MR-FFT on P .

Then

$$\phi_{\text{center}}(P, S) \leq 4 \cdot \phi_{\text{center}}^{\text{OPT}}(P, k)$$

so MR-FFT is a 4-approx algorithm \rightarrow so obtained

Proof: From the previous lemma, we know that $\forall x \in P \exists y \in T$ s.t.

$$d(x, y) \leq 2 \phi_{\text{center}}^{\text{OPT}}(P, k)$$

Recall that, as we can see from Pseudocode of MRFPT, S is extracted from T running FFT(T, k).

Let \bar{y} the point of T farthest from S and observe that $T \subseteq P$.

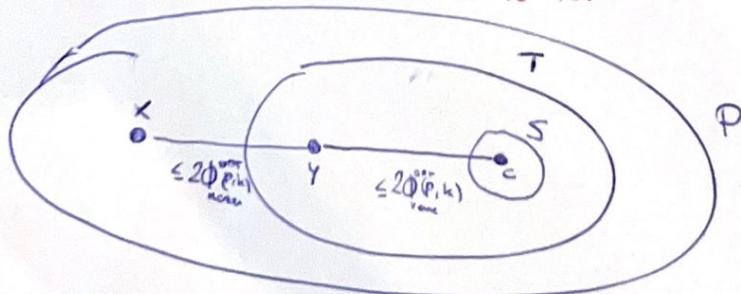
The same argument used in the previous lemma with the CLAIM can be repeated here

$$d(\bar{y}, S) \leq 2 \phi_{\text{center}}^{\text{OPT}}(P, k)$$

So how we have that $\bullet \forall x \in P \exists y \in T$ s.t. $d(x, y) \leq 2 \phi_{\text{center}}^{\text{OPT}}(P, k) \rightarrow$ lemma

$\bullet \forall y \in T, d(y, S) \leq d(\bar{y}, S) \leq 2 \phi_{\text{center}}^{\text{OPT}}(P, k)$

these two things implies that $\forall x \in P \exists$ "a good representative" of this point in T and \forall point of T the distance of the closest center is not too far



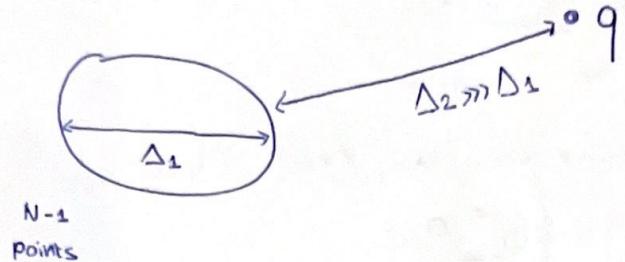
$$d(x, S) \leq d(x, c)$$

$$\leq d(x, y) + d(y, c)$$

$$\leq 4 \phi_{\text{center}}^{\text{OPT}}(P, k)$$

Example with $k=2$

Suppose that P is done like this:



so that points of P are "close" within a cluster and q is an outlier

We have that if T is selected randomly from P , with $|T| = \sqrt{N \cdot k}$ and $|P| = N$ with N very large

$$P[q \in T] = \sqrt{N \cdot k} \cdot \frac{1}{N} = \sqrt{\frac{k}{N}} \xrightarrow[N \rightarrow +\infty]{} 0$$

So if T doesn't contain q , and we search for 2 centers in T ,

for any selection $S \subseteq T$ we have that $\phi_{\text{center}}^{\text{over}}(P, S) \approx \Delta_1$ while

$$\phi_{\text{center}}^{\text{over}}(P, S) \approx \Delta_2$$

↑ the case where q is a center

↓ since q is at one cluster \uparrow

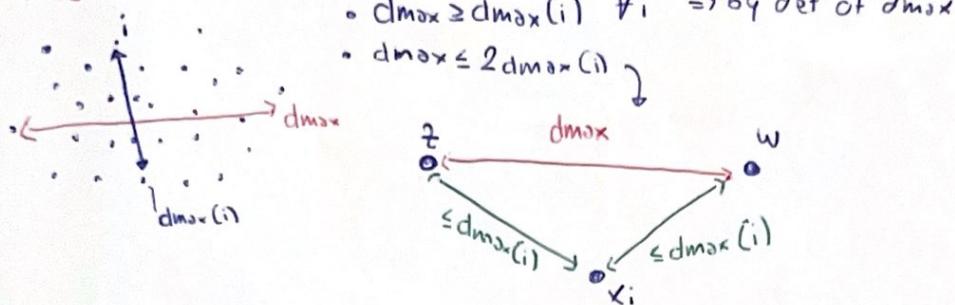
$$\phi(P, S^*) = \max_{x \in P} d(x, S^*)$$

Lemma:

For an arbitrary $x_i \in P$, we define $d_{\max}(i) = \max \{d(x_i, x_j) : 0 \leq j < N\}$

For any $0 \leq i < N$, we have that $d_{\max} \in [d_{\max}(i), 2d_{\max}(i)]$

Proof



By triang. ineq. $d_{\max} \leq d(z, x_i) + d(x_i, w) \leq 2d_{\max}(i)$

Is D_T a good approx. of D_{\max} ?

$T = \{c_1, \dots, c_k\}$, $q = \text{point of } P \text{ farthest from } T$

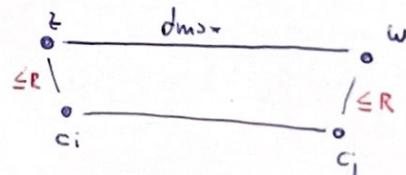
$$R = d(q, T)$$

$d_{\max} = d(z, w) \equiv \text{true diameter}$

Let $c_i = \text{center of } T \text{ closest to } z$

$c_j = \text{ " " to } w$

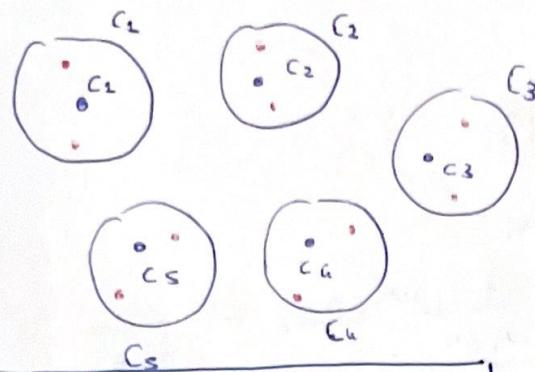
$$\begin{aligned} d_{\max} &= d(z, w) \leq d(z, c_i) + d(c_i, c_j) + d(c_j, w) \\ &\leq 2R + d(c_i, c_j) \quad \rightarrow \text{since } R = d(q, T) \geq d(x, T) \quad \forall x \in P \\ &\leq 2R + d_T \quad \rightarrow \text{since } d_T = \max_{(x, y) \in T} d(x, y) \\ \Rightarrow d_T &\leq d_{\max} \leq d_T + 2R \end{aligned}$$



As k grows, the value of R becomes smaller until it's negligible w.r.t d_{\max} .
For low d.m. space can be shown that $\log O(1)$ suffices to obtain $R \leq \epsilon \cdot d_{\max}$

Coreset-based algorithm for diversity maximization

Example with $h=5, k=3$



- Run FFT to extract T^h centers from P and assign cluster

- Select k points from each cluster

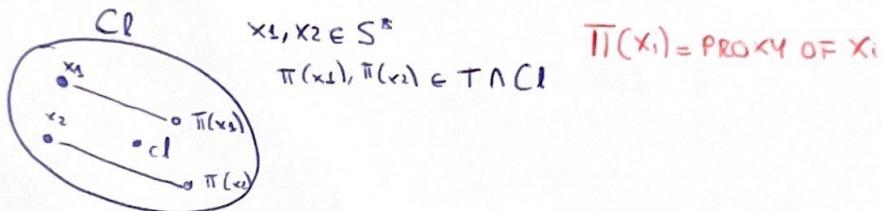
$$T = \{c_1 \dots c_5\} \cup \{\text{additional points for each } C_l\}$$

Why T is a good coresset?

Let S^* = optimal solution to diversity maximization for P .

IDEA: For cluster C_l ($1 \leq l \leq h$) create an injective mapping between $S^* \cap C_l$ and $T \cap C_l \Rightarrow$ this is always possible since $|T \cap C_l| = k$ and $|S^* \cap C_l| \leq |S^*| = k$

For example



Let R be the radius of the clustering i.e. $d(x, \{c_1 \dots c_h\}) \leq R \quad \forall x \in P$

For $i=1,2$

$$d(x_i, \pi(x_i)) \leq d(x_i, c_l) + d(c_l, \pi(x_i)) \leq 2R$$

\Rightarrow Means that point of S^* has a "nearby" proxy in T !

We now argue that T contains a good solution to div. maxim. for P

Consider the set $S = \{\pi(x) : x \in S^*\} \subset T \quad |S| = k$

$$\begin{aligned} \text{div}(S) &= \sum_{x_i, x_j \in S^*} d(\pi(x_i), \pi(x_j)) \geq \sum_{x_i, x_j \in S^*} d(x_i, x_j) - d(x_i, \pi(x_i)) - d(x_j, \pi(x_j)) \geq \sum_{\substack{x_i, x_j \\ |x_i-x_j| \leq 2R}} (d(x_i, x_j) - 4R) \\ &= \text{div}^{\text{opt}}(P, k) - \binom{k}{2} 4R \end{aligned}$$

↓
↳ Since it's done
k couples of points
and k combinations

If h is large enough we can make $4R \in \Phi_{\text{center}}^{\text{OPT}}(P, k) \rightarrow$ for fact

In this case $\text{div}(S) \geq \text{div}^{\text{OPT}}(P, k) - 4R \binom{k}{2}$

$$\geq \text{div}^{\text{OPT}}(P, k) - \epsilon \Phi_{\text{center}}^{\text{OPT}}(P, k) \binom{k}{2}$$

$$\geq \text{div}^{\text{OPT}}(P, k) - \epsilon \text{div}^{\text{OPT}}(P, k)$$

$$\Rightarrow \text{div}(S) \geq (1-\epsilon) \text{div}^{\text{OPT}}(P, k) = \frac{1}{1+\epsilon'} \text{div}^{\text{OPT}}(P, k) \quad \text{with } \epsilon' \text{ close to } \epsilon$$

$$\text{div}^{\text{OPT}}(T, k) \geq \text{div}(S) \geq \frac{1}{1+\epsilon'} \text{div}^{\text{OPT}}(P, k) \Rightarrow \checkmark$$

Example of application of Boyer-Moore algorithm

$\Sigma = x_1 x_2 \dots x_9$	\Rightarrow	Step	Cand	Count
A A A C C B B A A		0	NULL	0
	↑	1	A	1
	2	A	2	
	3	A	3	
	4	A	2	
	5	A	1	
	6	A	0	
	7	B	1	
	8	B	0	
	9	A	1	

\Rightarrow at the end: $\text{cand} = A$ ~~not true majority~~

Theorem: Given a stream Σ with a majority element m , the Boyer-Moore algorithm returns m using:

- working memory of size $O(1)$
- 1 pass
- $O(n)$ time per element

Proof:

• working memory $O(1)$
 • 1 pass
 • $O(n)$ time } straightforward (simple and direct proof)

CORRECTNESS: we must show that if a majority element m exists, then at the end of the stream $\text{cand} = m$

For $t = [0, n]$, with $n = |\Sigma|$, let cand_t , count_t be the values of cand and count after $x_1 \dots x_t$ ($t=0$ means initialization value)

It's easy to show that after processing $x_1 \dots x_t$, we have that these t elements can be partitioned into:
 * count occurrences of cand
 * $(t - \text{count})/2$ pairs (e_1, e_2) with $e_1 \neq e_2$

For the above example, $t=5 \Rightarrow x_1 \dots x_5 = A A A C C \sim \text{cand} = A, \text{count} = 1$
 $\Rightarrow 1$ occurrence of A and 2 pairs $(A, C), (A, C)$

We now show that $\text{cond}_m = m$.

By contradiction suppose that $\text{cond}_m \neq m$. Because of the invariants stated before, we have that Σ contains

- * $\text{Count}_n \geq 0$ occurrences of cond_n

- * $(\text{hcount}_n)/2$ pairs (e_1, e_2) with $e_1 \neq e_2$

Therefore, if $m \neq \text{cond}_n$, it can occur at most $(n - h\text{count}_n)/2 \leq \frac{n}{2}$ since
 $\text{Count}_n \geq 0 \Rightarrow m$ cannot be the majority \Rightarrow contradiction

$\frac{n-h\text{count}}{2} + h\text{count}$, It can appear
only in couples (e_1, e_2)

Theorem:

Let $\Sigma = x_1, x_2, \dots$. For any time $t \geq m$, the set S maintained by the reservoir sampling algorithm is an m -sample of $\Sigma_t = x_1, \dots, x_t$.

Proof:

Let $S_t = S$ after processing x_1, \dots, x_t . We now show by induction on $t \geq m$ that S_t is an m -sample of x_1, \dots, x_t , that is, for each $i \in [1, t]$, we must have that $P[x_i \in S_t] = m/t$.

Base of the induction: $t = m$, we can choose only these elements ✓

Inductive step: Suppose that the property holds for $t-1 \geq m$

$$\text{so } \forall i \in [1, t-1] \Rightarrow P[x_i \in S_{t-1}] = \frac{m}{t-1}$$

Now we evaluate $P[x_i \in S_t] \forall i \in [1, t]$:

$$\bullet P[x_t \in S_t] = \frac{m}{t} \text{ by construction}$$

- Consider $x_i, i < t$, and define the event $A = "x_t \text{ is added to } S \text{ at time } t"$

otherwise

$$\text{we have that } P[A] = \frac{m}{t} \quad P[\sim A] = 1 - \frac{m}{t}$$

since the algorithm choose with prob $\frac{m}{t}$ if evict an x_i in favor of x_t

By the law of total probabilities

$$P[x_i \in S_t] = P[x_i \in S_t \cap A] + P[x_i \in S_t \cap \sim A]$$

$$= P[x_i \in S_t | A] \cdot P[A] \sim \frac{m}{t}$$

$$P[x_i \in S_t \cap \sim A] = P[x_i \in S_t | \sim A] \cdot P[\sim A]$$

$$= \underbrace{P[x_i \in S_{t-1} | A]}_{\frac{m}{t-1}} \cdot \underbrace{P[x_i \text{ not evicted in step } t | A]}_{1 - \frac{1}{m}}$$

$$\left| \frac{m}{t-1} \cdot \left(1 - \frac{1}{m}\right) \cdot \frac{m}{t} \right| \rightarrow$$

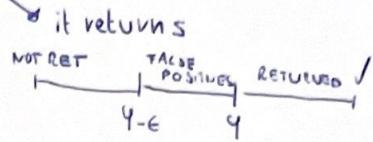
$$= \boxed{\frac{m}{t}}$$

$$P[x_i \in S_{t-1}] = \frac{m}{t-1}$$

$$+ \boxed{\frac{m}{t-1} \cdot \left(1 - \frac{m}{t}\right)}$$

Theorem: Sticky sampling solves the ϵ -AFI problem correctly with $P[\cdot] \geq 1-\delta$ and requires

- working memory of size $O(\ln(1/\delta\gamma)/\epsilon)$ in expectation
- 1 pass
- $O(1)$ expected time per element



Proof: Recall that $\Sigma = x_1 \dots x_n$

- 2) each input is seen only one time, so at the end the total of items = pass
- 3) whence we have find the item (working memory) we need to increment a count

WORKING MEMORY: • It's proportional to the size of S and, if S is a properly set, the size of S is proportional to the number of items it contains

- Each x_i contributes an extra unit of space to S with prob $\frac{r}{n}$ since it's added with probability $\frac{r}{n}$ only if not already in S

$$\Rightarrow \text{Expected working memory } O(n \cdot \frac{r}{n}) = O(r) \text{ and } r = \frac{\ln(\frac{1}{\delta\gamma})}{\epsilon}$$

CORRECTNESS: We need to show that with $P[\cdot] \geq 1-\delta$, the set of items returned is s.t.

- A) all true freq. items are in the set
- B) no item occurring $< (4-\epsilon) \cdot n$ times in Σ are in the set
 - ↳ ensured by the fact that we return only items whose # of occurrences is $\geq (4-\epsilon) \cdot n$
- Let's verify A)

Consider a frequent item $a \Rightarrow P[a \text{ is not returned in the output}] \leq P[\text{none of the}$

$$\Rightarrow \leq \left(1 - \frac{r}{n}\right)^{E \cdot n} \leq \left(1 - \frac{r}{n}\right)^{E \cdot n} = \left(1 - \frac{r}{n}\right)^{\frac{n \cdot r \cdot \epsilon n}{n}} = \left(\left(1 - \frac{1}{n/n}\right)^{n/n}\right)^{E \cdot r} \text{ first } E \cdot n \text{ occur. of } a \text{ in } \Sigma \text{ are sampled} \rightarrow$$

$$\leq \left(\frac{1}{e}\right)^{E \cdot r} \text{ where we use the fact that } \left(1 - \frac{1}{x}\right)^x \leq \frac{1}{e} \quad \forall x \geq 1$$

Let $a_1 \dots a_k$ be the k true freq. items and recall that $k \leq \frac{1}{\epsilon}$

$$P[\text{some } a_i \text{ not returned}] \leq \sum_{i=1}^k P[a_i \text{ not returned}] \leq k \cdot \left(\frac{1}{e}\right)^{E \cdot r} \leq$$

$$\leq \frac{1}{\epsilon} \left(\frac{1}{e}\right)^{E \cdot \ln(\dots)/\epsilon} \stackrel{U.B.}{=} \frac{1}{\epsilon} \left(\frac{1}{e}\right)^{\ln\left(\frac{1}{\delta\gamma}\right)} = \delta \Rightarrow P[\text{all items are returned}] \geq 1 - \delta$$

Why? because
S returns $\geq (4-\epsilon) \cdot n$ items.
So, if $E \cdot n$ occ. are missed, we have $\leq \epsilon$ occ. \Rightarrow not ret.

Example

$U = \text{alphabet } |U|=26 \Rightarrow \text{we need to represent } L=25 \text{ in binary: 5 bits!}$

$\Sigma = A, D, A, A, C, B, F, F, B, A, E, C$

x	$h(x)$	$\text{tr}(h(x))$	C (6 bits)	$R=3 \Rightarrow \tilde{F}_0 = 2^{R-3}$ while $F_0 = 6$
A	01001	0	100000	
D	11100	2	101000	
C	11000	3	101100	
B	01110	1	111100	
F	10100	2	"	
E	10011	0	"	

Why does it work? (intuition)

For simplicity, assume $|U| \geq \text{power of 2}$

For $x \in \Sigma$, what is the prob that $h(x)$ has at least j trailing zeros?

$$\Pr[\text{tr}(h(x)) \geq j] = \frac{1}{2^j}$$

⋮

$$\Pr[\text{tr}(h(x)) \geq j] = \frac{1}{2^j}$$

Intuitively, interval $[0, |U|-1]$ contains $\frac{|U|}{2^j}$ integers whose binary configuration is $\dots 000\dots 0$ (j trailing zeros)

Configurations with $\geq j$ trailing zeros. So, a random integer have $\geq \frac{1}{2^j}$ probability to be an integer with $\geq j$ zero (trailing)

The algorithm will generate F_0 random integers in $[0, |U|-1]$ and the expected number of these integers with $\geq j$ trailing zeros is $\frac{F_0}{2^j} \Rightarrow$ due to indep. B(F_0, p) $\rightarrow E[X] = \tilde{F}_0 \cdot p$

So we expect to generate integers with $\geq j$ trailing zero's only if $\frac{F_0}{2^j} \geq 1$

Therefore, the maximum number of trailing 0's obtained when processing the stream Σ (i.e. R) is expected to be s.t

$$\frac{F_0}{2^R} \geq 1 \text{ but, } \frac{F_0}{2^R} \rightarrow 2^{R-j} \tilde{F}_0$$

$\underbrace{2^R}_{\text{3 trailing}} \underbrace{\tilde{F}_0}_{\text{the element with 2 trailing 0's with prob. }} \rightarrow 2^{R-j} \tilde{F}_0$

example:

$$n=15, d=3, w=3$$

$$\Sigma = A, B, C, B, D, A, C, D, A, B, D, C, A, A, B$$

U, f_U	h_0	h_1	h_2
A, 5	0	1	1
B, 4	1	2	1
C, 3	0	0	2
D, 3	1	1	2

array C		
SA	4B	
3C	3D	
3C	SA	4B
	3D	
	SA	3C
	4B	3D

$$\tilde{f}_A = \min\{8, 8, 9\} = 8 > f_A = 5$$

$$\tilde{f}_B = \min\{7, 6, 9\} = 6 = f_B$$

$$\tilde{f}_C = \min\{8, 3, 6\} = 3 = f_C$$

$$\tilde{f}_D = \min\{9, 8, 6\} = 6 > f_D = 3$$

d
1
2

↔ w 1 2

Theorem: Consider a $d \times w$ count min sketch for a stream Σ of length n , where

$$d = \log_2\left(\frac{1}{\delta}\right) \text{ and } w = \frac{2}{\epsilon}, \text{ for some } \epsilon, \delta \in (0, 1)$$

for any given $u \in U$, occurring in Σ ,

$$\tilde{f}_u - f_u \leq \epsilon \cdot n$$

with prob $1 - \delta$

Proof: Fix an arbitrary $u \in U$ and $j \in [0, d-1]$

- $C[j, h_j[u]]$ receives, in expectation, a fraction w of the total mass $n = |\Sigma|$
- \Rightarrow receives a value $\frac{n}{w}$ and includes for sure $f_u \Rightarrow E[C[j, h_j[u]] - f_u] \leq \frac{n}{w} = \frac{\epsilon n}{2}$
- \Rightarrow By Markov inequality

$E[\cdot] = \sum u \cdot \frac{1}{w} = \frac{\sum u}{w} \Rightarrow \frac{n}{w}$

↓
since we insert it
the table surely f_u
times u

$$P[C[j, h_j[u]] - f_u \leq \epsilon \cdot n] \geq \frac{1}{2} \Rightarrow P[C[j, h_j[u]] - f_u > \epsilon \cdot n] < \frac{1}{2}$$

Since \tilde{f}_u is the minimum among those rows ($d = \log_2\left(\frac{1}{\delta}\right)$), we can say that

$\tilde{f}_u - f_u > \epsilon \cdot n$ only if all rows $> \epsilon \cdot n$ and this happen with prob $< \left(\frac{1}{2}\right)^d = \delta$

$$\text{So } P[\tilde{f}_u - f_u \leq \epsilon \cdot n] \geq 1 - \delta$$

Example:

$$n=15, d=3, w=3$$

$\Sigma = A, B, C, B, D, A, C, D, A, B, D, C, A, A, B$

u, v	h_0	g_0	h_1	g_1	h_2	g_2	h_3	g_3
A, S	0	1	1	1	1	1	1	1
B, 4	1	-1	2	1	1	1	-1	
C, 3	0	-1	0	-1	2	1	1	
D, 3	1	-1	1	1	2	1	1	

array C		
$5A - 3C$	$-4B - 3D$	
$-3C$	$5A + 3D$	$4B$
	$5A - 4B$	$3C + 3D$

$$\begin{aligned} \tilde{f}_{A,D} &= 2 & \tilde{f}_{A,B} &= 8 & \tilde{f}_{A,C} &= 1 \Rightarrow \tilde{f}_A = \text{median}\{2, 8, 1\} = 2 < f_A = 5 \\ \tilde{f}_B &= \text{median}\{7, 4, -1\} = 4 = f_B \\ \tilde{f}_C &= \text{median}\{-2, 3, 6\} = 3 = f_C \\ \tilde{f}_D &= \text{median}\{7, 8, 5\} = 7 > f_D \end{aligned}$$

Theorem

Proof: We prove only (A)

Fix $u \in U$ and $j \in [0, d-1]$ arbitrarily. We show that $\mathbb{E}[f_{u,j}] = f_u$

$\forall a \in U, a \neq u$, define $\hookrightarrow g_j(a) \cdot C[j, h_j(a)]$

$$y_a = \begin{cases} f_a & \text{if } h_j(a) = h_j(u) \text{ AND } g_j(a) = g_j(u) \Rightarrow \text{are in the same cell for row } i \text{ and have the same sign} \\ -f_a & \text{if } h_j(a) = h_j(u) \text{ AND } g_j(a) \neq g_j(u) \\ 0 & \text{i.e. } h_j(a) \neq h_j(u) \end{cases}$$

CRUCIAL OBSERVATION: $y_a = \text{contribution of } a \neq u \text{ to } \tilde{f}_{u,j}$

$$\tilde{f}_{u,j} = f_u + \sum_{a \in U, a \neq u} y_a$$

Now for $a \neq u$

$$P[y_a = f_a] = P[y_a = -f_a] = \frac{1}{w} \cdot \frac{1}{2} \quad / \quad g_j(a) \text{ random } \in \{-1, 1\}$$

$$P[y_a = 0] = 1 - 2 \cdot \frac{1}{w} \cdot \frac{1}{2} = 1 - \frac{1}{w} \quad / \quad \begin{matrix} h_j(a) = h_j(u) \\ \downarrow \\ 1-2 \text{ possibilities above} \end{matrix} \quad \Rightarrow \text{IN}$$

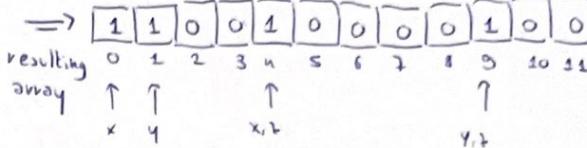
$$\mathbb{E}[\tilde{f}_{U,j}] = \mathbb{E}[f_U + \sum_{\omega \in U, \omega \neq j} y_\omega] = f_U + \mathbb{E}[\sum y_\omega] = f_U + \mathbb{E}[y_\omega] = f_U$$
$$f_\omega\left(\frac{1}{\omega}, \frac{1}{2}\right) - f_\omega\left(\frac{1}{\omega}, \frac{1}{2}\right) = 0$$

example

$$S = \{x_1, 4, 2\} \quad n=12 \quad k=2$$

	X	4	2
h ₀	0	1	4
h ₂	4	9	9

resulting array



Check(X): $h_0(x) = 0 \quad h_1(x) = 4 \quad$ Yes, true positive

Check(T): $h_0(T) = 4 \quad h_2(T) = 7 \quad$ No, true negative

Check(P): $h_0(P) = 1 \quad$ ~~not 1~~ $h_2(P) = 9 \quad$ Yes, ~~true positive~~ False]

Theorem:

Proof: Let $S = \{e_1, \dots, e_m\}$

Under the assumption that the k hash functions are independent, provide indices uniformly in $[0, n-1]$ and $h_j(e)$ and $h_j(e')$ are independent, $\forall j \in [0, k-1]$ and $e, e' \in S$.

The indices of the cells of A which contain 1's can be seen as $k \cdot m$ independent variables uniformly distributed in $[0, n-1]$

\Rightarrow for any index $l \in [0, n-1]$, we have that

$$P[A[l]=0] = \left(1 - \frac{1}{n}\right)^{k \cdot m} = \left(1 - \frac{1}{n}\right)^{n \cdot km} = \left(\left(1 - \frac{1}{n}\right)^n\right)^{km} \approx \left(\frac{1}{e}\right)^{\frac{km}{n}}$$

we define $p = e^{-\frac{km}{n}}$ and we make the SIMPLIFYING ASSUMPTION that A contains exactly $p \cdot n$ 0's (the theorem says true even without this assumption) true if n is large enough

Consider $x \in \Sigma$ s.t. $x \notin S$ and let $l_j = h_j(x) \quad 0 \leq j \leq k$

Since $h_j(x)$ is uniformly distributed in $[0, n-1]$

$$P[A[l_j]=1] = 1 - p^m = 1 - p \quad \text{i.e. prob of not hitting a 0}$$

$$P[A[l_j]=0]$$

$$P[A[l_j]=1 \quad \forall 0 \leq j < k] = \prod_{j=0}^{k-1} P[A[l_j]=1] = (1-p)^k = (1 - e^{-\frac{km}{n}})^k$$

Proof:

* PROBABILISTIC CORRECTNESS

Case 1: $B_r(q) \cap P \neq \emptyset$. In this case, a legal output is any point $p' \in P$ s.t.

$$d(q, p') \leq c.r$$

Now consider an arbitrary point $p \in B_r(q) \cap P$ (exist)

$$\begin{aligned} \text{Then: } \Pr(\text{answer correct}) &= \Pr[\text{answer} \neq \text{null}] \geq \Pr[p \text{ is mapped in the} \\ &\quad \text{same region by} \\ &\quad h] \\ &= p_1 \end{aligned}$$

\downarrow
by property dc H

Case 2: $B_r(q) \cap P = \emptyset$. The answer will be surely correct both if a point p at distance $\leq c.r$ is returned and if null is RETURNED.

* CONSTRUCTION AND SPACE TIME: Straightforward

* QUERY TIME: Suppose that the bucket $T[h(q)]$ is implemented as a list.

The scan of $T[h(q)]$ ends as soon as a point at distance at most $\leq c.r$ from q is found or the end of the list is reached.

Therefore, at most $x+1$ points of $T[h(q)]$ are checked where

x is the number of points in the bucket at distance $> c.r$ from q .

Since a far point (i.e. at distance $> c.r$) has prob. at most p_2 of being in the bucket $T[h(q)]$ and since there are at most n far points in P , then, the expected number $E[x] \leq n \cdot p_2$

all the
points