

Algorithm engineering

ANALYSIS WITH PSEUDOCODES AND THEOREMS WITH PROOFS
GIUSEPPE LOMBARDI

Table of contents

Introduction	5
2-level memory model (disk model)	5
B-tree	6
Atomic items sorting	8
Merge-based sorting paradigm	8
Disk striping	11
Permuting	12
Upper bounds	13
Lower bound on sorting	14
Distribution-based sorting	16
Random sampling	22
Disk model (known sequence)	22
Streaming model (known sequence)	24
Streaming model (unknown sequence)	25
Randomized data structures	27
Treap	27
Skip List	31
Set intersection	34
Intro	34
Mutual partitioning	36
Doubling searching	37
2-level storage approach	39
String sorting	41
Intro	41
Radix sort (MSD first)	43
Radix sort (LSD first)	46
Multi-key Quicksort	49
Dictionary problem	52
Intro	52
Direct address table	53
Hash Table (hashing with chaining)	54
Universal hashing	57
Perfect hashing, minimal, ordered	61
Cuckoo hashing	63
Bloom filters	66

Spectral bloom filters.....	68
Prefix search	70
Intro.....	70
Array of string pointers	71
Contiguous allocation of strings	73
Front coding with bucketing.....	75
Interpolation search.....	77
Compacted trie	79
Patricia trie	82
Substring search.....	84
Intro.....	84
Suffix array.....	85
Text mining.....	87
Integer coding.....	88
Intro.....	88
Unary code	90
Fixed-length binary code.....	91
Gamma-code (Elias)	92
Delta-code (Elias)	94
Rice code	96
Variable-byte code.....	97
(s,c)-dense code	99
PForDelta code	100
Interpolative code.....	101
Elias-Fano code	102
Compressed data structure	105
Rank and Select.....	105
Succint solution for Rank (via Rank data structure).....	106
Compressed Rank/Select based on Elias-Fano coding.....	109
Succint representation of binary tree (LOUDS-like).....	112
Statistical coding	114
Intro.....	114
Huffman coding	115
Canonical Huffman coding	119
Arithmetic coding.....	121
Dictionary-based compressors	125

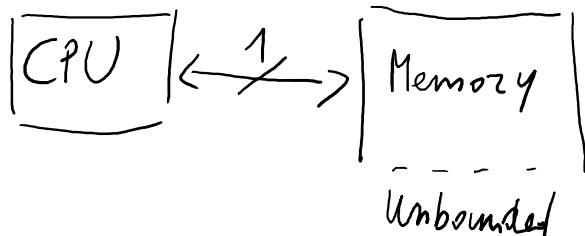
Intro.....	125
LZ77	126
LZss.....	128
Gzip (fast LZ77)	130
LZ78.....	132
Burrows-Wheeler transform	134
Intro.....	134
Forward transform.....	135
Backward transform.....	136
Move-To-Front transform (MTF)	138
Run-Length Encoding transform (RLE)	140
BZIP	141

2-level memory model (disk model)

venerdì 27 gennaio 2023 15:32

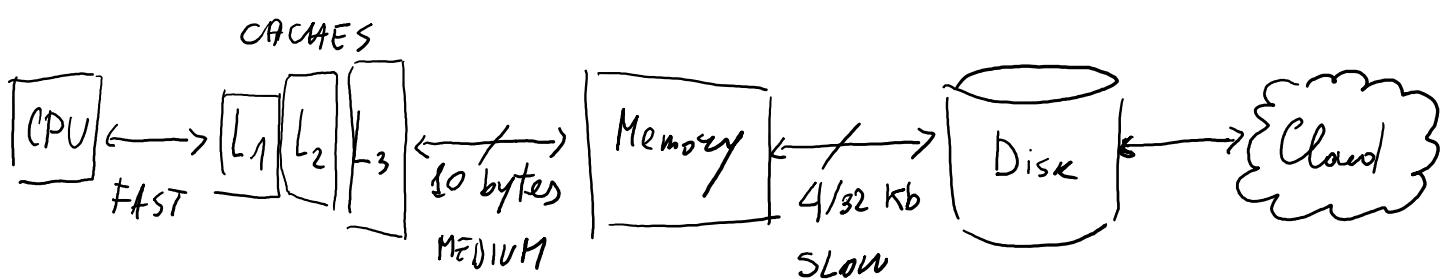
I/Os good estimator for time complexity

Von Neumann Model (RAM model)

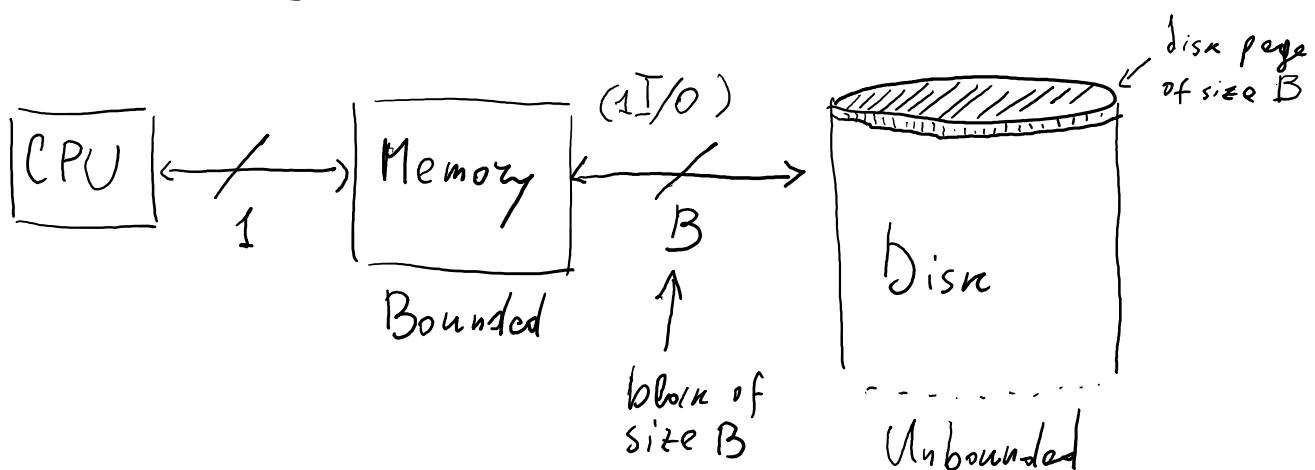


$T(n) = \# \text{ steps for input size } n$
(asymptotic and worst case)

Modern machine (hierarchy of memory levels)



2-layered Memory model



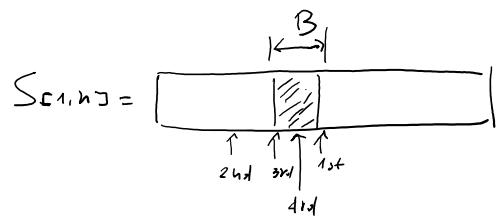
$$T(h) = \# \text{ I/Os}$$

B-tree

venerdì 3 marzo 2023 19:33

Binary Search

$$T(n) = \mathcal{O}(\log_2 n)$$

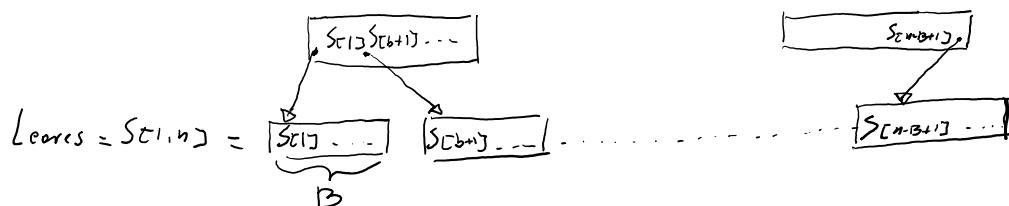


We expect $\mathcal{O}(\log_2 n) - \underbrace{\mathcal{O}(\log_2 B)}_{\text{binary search}} \text{ I/Os} = \mathcal{O}(\log_2 \frac{n}{B}) \text{ I/Os}$
in subseq. of size $\leq B$

An improvement:



B^+ -tree



$$\begin{aligned} \text{FAN-OUT} &= \frac{B}{\text{size(key)} + \text{size(pointer)}} = \mathcal{O}(B) \\ \# \text{LEAVES} &= \frac{B}{B} \end{aligned} \quad \left. \right\} \Rightarrow \text{height} = \log_B \frac{n}{B}$$

SEARCH(B^+ -tree):

node \leftarrow ROOT(B^+ -tree);

while ($\text{node} \neq \emptyset$) {

$k \leftarrow \text{BINARYSEARCH}(\text{node});$ // 1 I/O (some page)

 node $\leftarrow \text{CHILDREN}(\text{node}, k);$ // goes down to k-th children

}

- Extra space:
- Time: $(\log_B \frac{n}{B}) \cdot \log_2 B = \frac{\log_2 \frac{n}{B}}{\log_2 B} \cdot \log_2 B = \log_2 \frac{n}{B}$
- I/Os: $\log_B \frac{n}{B}$ (= height of tree)

Usually: $B = 32 \text{ kb} = 2^5 \cdot 2^{10} \text{ byte} = 2^{15} \text{ byte}$

key/pointer = 2³ byte

$$\log_B \frac{n}{B} \text{ I/Os} = \frac{\log_2 \frac{n}{B}}{\log_2 B} = \frac{\log_2 \frac{n}{B}}{15} \quad \text{Assuming } h = 1 \text{ terabyte} = 2^{40} \text{ byte}$$

$\approx 2^{40}/2^{15}$

$$\frac{ly_2}{15} = \frac{25}{15} \text{ ok } !!$$

Merge-based sorting paradigm

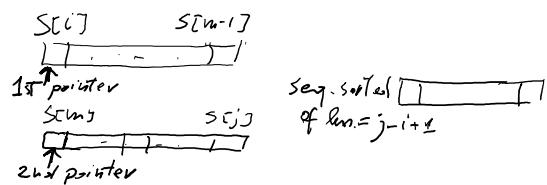
lunedì 13 febbraio 2023 20:04

$$S[1..n] = S_1, \dots, S_n$$

MERGE-SORT (S, i, j)

```

if  $i < j$  {
     $m = (i+j)/2;$ 
    MERGESORT ( $S, i, m-1$ );
    MERGESORT ( $S, m, j$ );
    MERGE ( $S, i, m, j$ );
}
  
```



- Extra space: $O(n)$ (aux. array for merging)

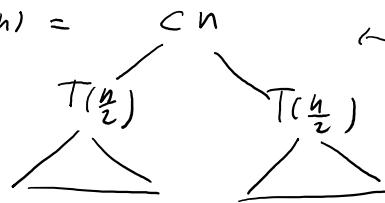
- Time: Merge sort (S, n)

$$T(n) = 2T(n/2) + O(n) \stackrel{(*)}{\Rightarrow} O(n \lg n)$$

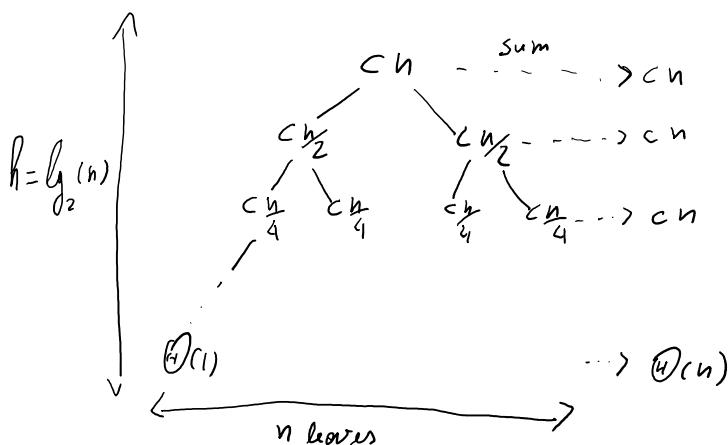
↑
Merging

$$(*) T(n) = \begin{cases} \Theta(1) & \text{if } n=1 \\ 2T(n/2) + \Theta(n) & \text{if } n>1 \end{cases}$$

Recursion tree: $T(n) =$



~ the cost of this algorithm is the sum of the cost of all the node in the recursive tree

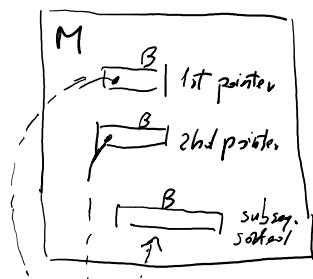


$$\text{Total cost} = \underbrace{cn}_{R_n} \cdot cn + \Theta(n)$$

$$= cn \lg n + \Theta(n) = \Theta(n \lg n)$$

Disk Model ($n > M$)

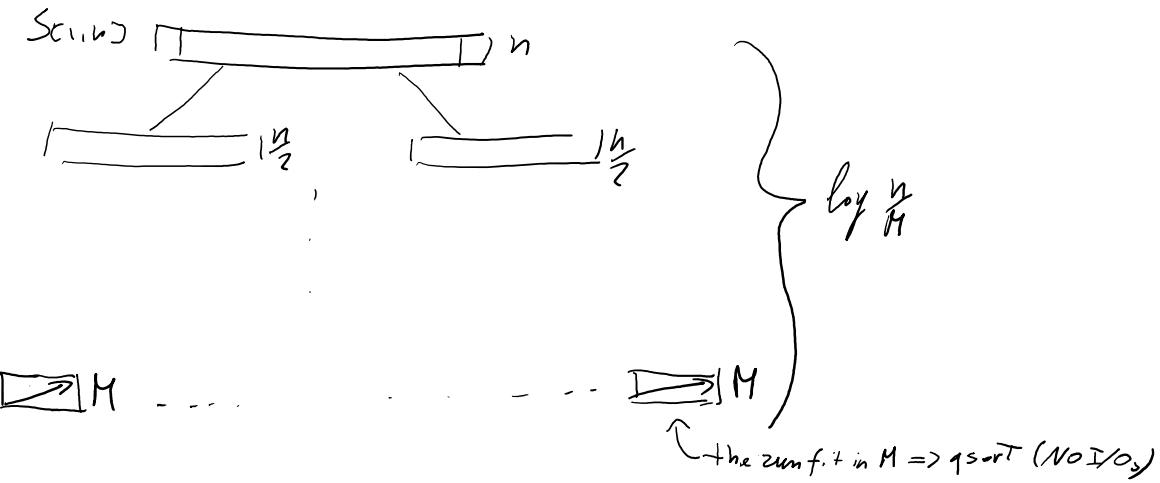
Suppose $M \geq 3B$



$\Theta(n/B)$ I/Os to merge
2 subseq. of total
size = n

$$S = [B | B | \dots | B |]$$

$$S[1..n] \quad [\quad] \quad n$$



$\bullet \text{I/Os} : O\left(\frac{n}{B} \lg_2 \frac{n}{M}\right) \quad (= O\left(\frac{n}{B} \lg_2 n\right) - O\left(\frac{n}{B} \lg_2 M\right))$

$\uparrow \text{I/Os for merging}$ $\uparrow \text{levels}$ $\uparrow \text{I/Os sorted}$

(1st RUN)
SNOW PLOW (S, M):

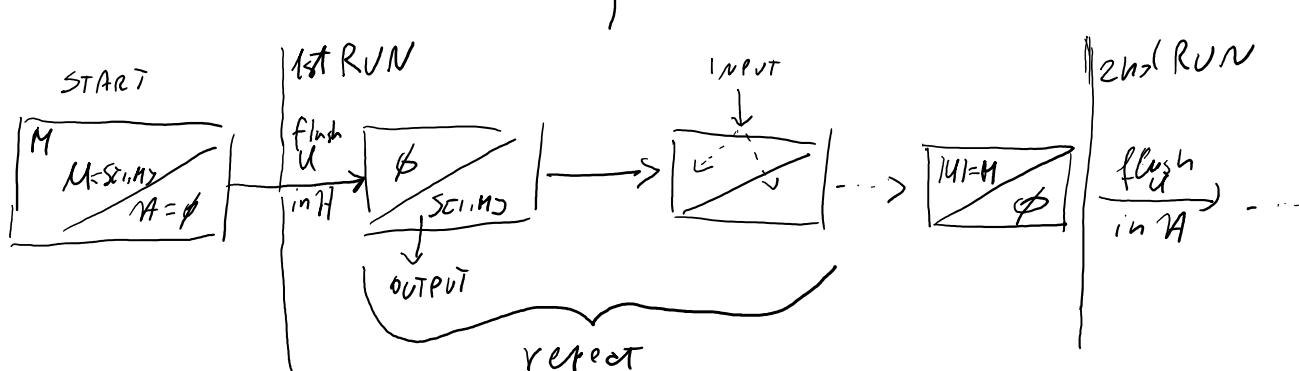
```

 $M \leftarrow S[1, M];$ 
 $A \leftarrow \text{MIN-HEAP}(M);$ 
 $M \leftarrow \emptyset;$ 
 $\text{while } (M \neq \emptyset) \{$ 
     $\text{min} \leftarrow \text{POP}(A);$ 
     $\text{OUTPUT}(\text{min});$ 
     $\text{next} \leftarrow \text{NEXT}(S);$ 
     $\text{if } (\text{next} < \text{min})$ 
         $M \leftarrow \text{ADD}(M, \text{next});$ 
     $\text{else}$ 
         $A \leftarrow \text{WSERT}(A, \text{next});$ 
     $\}$ 

```

Sci, n, $U = \text{set}$, $A = \text{min-heap}$

- Extra space: No (array-based impl. for heap)
 - I/Os: $O\left(\frac{n}{B} \lg_2 \frac{n}{(2M)}\right)$ (see Prop)
- \uparrow orig. size of runs.



Prop Average length of runs is $\gamma = 2M$, supposing random distribution

Proof γ items of a run go in U or A . But in U they will not be removed in a run.

$$\gamma = M + (\gamma - M) \quad , \text{Random dist.} \Rightarrow p(\text{next} < \min) = \frac{1}{2}$$

$\uparrow \quad \uparrow$
inserted in U inserted in \mathcal{B}

$$\Rightarrow E[\text{items per run}] = \frac{1}{2} \cdot T = M$$

$$\Rightarrow E[\text{size run}] = \bar{x} = 2M \quad \square$$

Cor $p(\text{next} < \min) = \alpha \quad (0 < \alpha < 1) \Rightarrow \gamma = \frac{M}{\alpha}$

Suppose $M \gg 3B$, and $K = \left(\frac{M}{B}-1\right) \gg 2$ block erasables
aux. space

START: k ordered subseq. $R_{1[1]}, R_k$ (n distinct runs)

\mathcal{H} min-heap containing $\langle R_{1[1]}, 1 \rangle, \dots, \langle R_k[1], k \rangle$

↑ ↑
smallest item run

MULTI-WAY MERGESORT(R_1, \dots, R_n):

$x_1 \leftarrow 1; \dots; x_n \leftarrow 1;$

$\mathcal{H} \leftarrow \text{MIN-HEAP}((R_1[x_1], 1), \dots, (R_n[x_n], n))$,
while ($\mathcal{H} \neq \emptyset$) {

$(\min, \text{run}) \leftarrow \text{POP}(\mathcal{H})$;

$\text{OUTPUT}(\min)$;

$x_{\text{run}}++$;

 if $x_{\text{run}} \leq \text{LEN}(R_{\text{run}})$

$H \leftarrow \text{INSERT}(R_{\text{run}}[x_{\text{run}}], \text{run})$;

}

• Extra space: NO

• Time / # Comparisons: $O(n \log n)$

• I/Os: $O\left(\frac{n}{B} \log \frac{\binom{M}{B} M}{n}\right) = O\left(\frac{n}{B} \log \frac{n}{B} \frac{n}{B}\right)$

seg (runs)
merged at time

$$g_{\frac{M}{B}}(M) = g_{\frac{M}{B}}\left(\frac{M}{B} \times B\right)$$

$$= 1 + g_{\frac{n}{B}}(B)$$

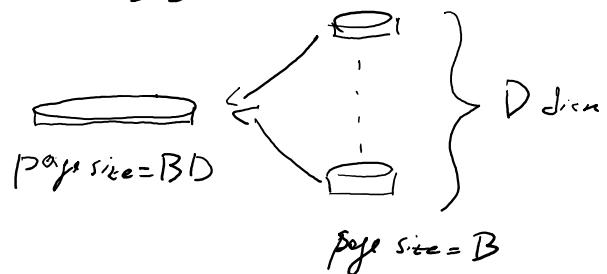
$$= \Theta(g_{\frac{n}{B}}(B))$$

Disk striping

mercoledì 1 marzo 2023 15:46

Using $D > 1$ disks to increase the bandwidth of the I/O. Simplest approach is to consider one big disk of page size $B' = DB$

(DB must be $\leq M$)



Using the multi-way mergesort the I/O complexity is

$$\mathcal{O}\left(\frac{n}{DB} \lg_{\frac{M}{DB}} \frac{n}{M}\right) > \Omega\left(\frac{n}{BD} \lg_{\frac{M}{B}} \frac{n}{M}\right)$$

↑
MULTI-WAY MERGESORT ↑
LOWER BOUND

Analysing the ratio r :

$$r = \frac{\frac{n}{BD} \lg_{\frac{M}{BD}} \left(\frac{n}{M}\right)}{\frac{n}{BD} \lg_{\frac{M}{B}} \left(\frac{n}{M}\right)} = \frac{\frac{\lg_2 \left(\frac{n}{M}\right)}{\lg_2 \left(\frac{M}{BD}\right)}}{\frac{\lg_2 \left(\frac{n}{M}\right)}{\lg_2 \left(\frac{M}{B}\right)}} = \frac{\lg_2 \left(\frac{M}{B}\right)}{\lg_2 \left(\frac{M}{BD}\right)} = \frac{\lg_2 \left(\frac{M}{B}\right)}{\lg_2 \left(\frac{M}{B}\right) - \lg_2(D)}$$

$$= \frac{1}{1 - \frac{\lg_2(D)}{\lg_2 \left(\frac{M}{B}\right)}} = \frac{1}{1 - \frac{\lg_M(D)}{\lg_B(D)}}$$

Cor Multi-way Mergesort is not optimal with $D > 1$ and the ratio r is $r > 1$, r is high if $D \geq \frac{M}{B}$ ($1 < D \leq \frac{M}{B}$)

Permuting

mercoledì 1 marzo 2023 11:28

$$\pi_{[1,n]} = \pi^{(\overbrace{1}^{\pi(1)})}, \dots, \pi^{(n)}$$

$$S = S_{[1,n]} \quad \langle S, \pi \rangle \longrightarrow S[\pi^{(1)}], \dots, S[\pi^{(n)}]$$

- RAM MODEL:

`PERMUTING_RAM(S, π):`

```

 $S' = S;$  | time:  $O(n)$ 
for  $i=1, \dots, n$ 
 $S'[i] = S[\pi[i]],$ 
return  $S'$ 

```

- DISK MODEL:

`PERMUTING1_DISK(S, π)` mimics `PERMUTING_RAM(S, π)`

• I/Os: $O(n)$

`PERMUTING2_DISK(S, π):`

```

 $P \leftarrow \langle 1, \pi(1) \rangle, \dots, \langle n, \pi(n) \rangle;$ 
 $P \leftarrow \text{SORT}(P, 2);$  // sort on 2nd component
for  $i=1, \dots, n$ 
 $P[0, 2] \leftarrow S[i];$  //  $P = \langle n_1, S[1] \rangle, \dots, \langle n_n, S[n] \rangle$ 
 $P \leftarrow \text{SORT}(P, 1);$  // sort on 1st component
return  $P[:, 2];$ 

```

$$\begin{aligned}
 & \cdot \text{I/Os: } 2 \underbrace{\frac{n}{B}}_{\text{2 scans}} + 2 \underbrace{\frac{n}{B} \lg \frac{n}{B}}_{\text{2 sorts}} \frac{n}{M} \\
 & = O\left(\frac{n}{B} \lg \frac{n}{B} \frac{n}{M}\right)
 \end{aligned}$$

Thm permuting n items takes $O(\min\{n, \frac{n}{B} \lg \frac{n}{M}\})$ I/Os on disk model

Cuz On Disk Model I/O complexity of sorting and of permuting are equal iff $\frac{n}{B} \lg \frac{n}{M} < n \Leftrightarrow B > \lg \frac{n}{M}$, so always true in practice

Upper bounds

mercoledì 1 marzo 2023 12:44

	RAM model	Disk Model	
Permuting	$O(n)$ scanning	$O(\min\{n, \frac{n}{B} \lg \frac{n}{M}\})$	SECRETARY OF PERMUTING
Sorting	$O(n \lg n)$ MERGE SORT	$O\left(\frac{n}{B} \lg \frac{n}{M}\right)$	MULTI-WAY MERGESORT

In case $B > 1$

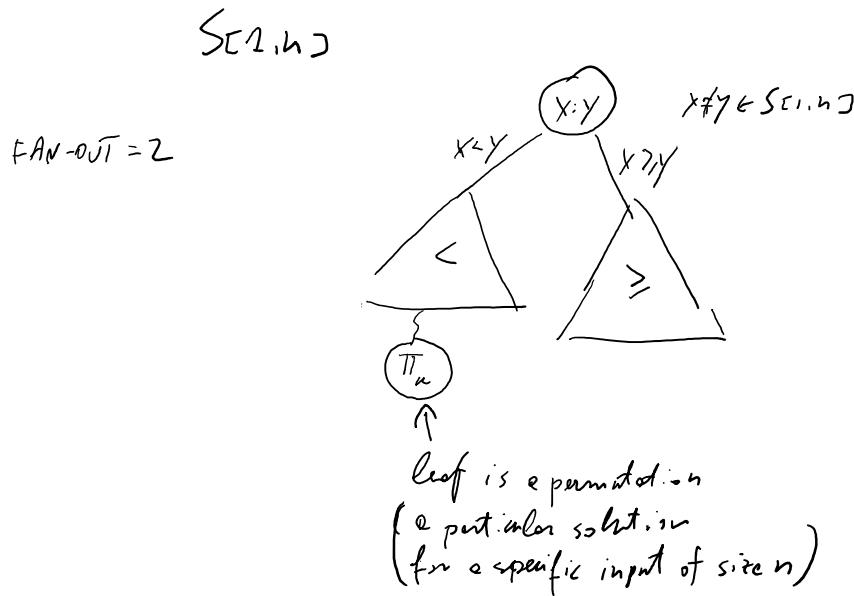
	Disk model	
Sorting	$O\left(\frac{n}{BD} \lg \frac{n}{M}\right)$	MULTI-WAY MERGESORT with DISK STRIPPING TECHNIQUE

Lower bound on sorting

mercoledì 1 marzo 2023 15:00

- RAM model (Comparison-based sorting)

A binary tree for each input size: each node represents a comparison between two items



$$We must have \ 2^h \geq n! \quad (\Rightarrow h \geq \lg_2(n!) \approx n \lg_2(n))$$

\uparrow all possible permutations (leaves) \uparrow Stirling's approximation

Prop time complexity lower bound (RAM model) for comparison based sorting is $\Omega(n \lg n)$

- Disk Model

A node in the comparison tree is an I/O (so reads B items from disk)

$$FAN-OUT = \# \text{ permutations generated by the I/O}$$

$$FAN-OUT = \begin{cases} \binom{M}{B} B! & \text{for new items} \\ \binom{M}{B} & \text{for old item (just fetched)} \end{cases}$$

$$\left| \begin{array}{l} \# \text{ leaves} = n! \\ \uparrow \\ \text{all possible permutations} \end{array} \right.$$

The #I/Os occurring to new items must be $\frac{n}{B}$. So assuming

$$t = \frac{n}{B} + \left(t - \frac{n}{B}\right)$$

↑
all I/Os ↑
"new" I/Os ↑
"old" I/Os

We must have:

$$\left[\binom{M}{B} B! \right]^{\frac{n}{B}} \cdot \binom{M}{B}^{t - \frac{n}{B}} \geq n!$$

$$\lg \left(\left[\binom{M}{B} B! \right]^{\frac{n}{B}} \cdot \binom{M}{B}^{t - \frac{n}{B}} \right) \geq \lg(n!)$$

$$\frac{n}{B} \lg \left(\binom{M}{B} B! \right) + \left(t - \frac{n}{B} \right) \lg \left(\binom{M}{B} \right) \geq n \lg n$$

$$\frac{n}{B} \left(B \lg \left(\frac{M}{B} \right) + B \lg B \right) + \left(t - \frac{n}{B} \right) B \lg \left(\frac{M}{B} \right) \geq n \lg n$$

$$n \lg B + t B \lg \left(\frac{M}{B} \right) \geq n \lg n$$

$$t \geq \frac{n \lg n - n \lg B}{B \lg \left(\frac{M}{B} \right)} = \frac{n \lg \left(\frac{n}{B} \right)}{B \lg \left(\frac{M}{B} \right)} = \frac{n}{B} \lg_{\frac{M}{B}} \left(\frac{n}{B} \right)$$

Cor 1 Lower bound for comparison-based sorting in disk-model is $\Omega \left(\frac{n}{B} \lg_{\frac{M}{B}} \left(\frac{n}{B} \right) \right)$ I/Os

$\forall D > 1$ disks $\Omega \left(\frac{n}{BD} \lg_{\frac{M}{BD}} \left(\frac{n}{BD} \right) \right)$ I/Os (same requirement with I/O over BD items)

Cor 2 Multi-way Mergesort is time and I/Os optimal if $D=1$.
($D>1$ see Disk striping)

Distribution-based sorting

mercoledì 1 marzo 2023 16:20

```

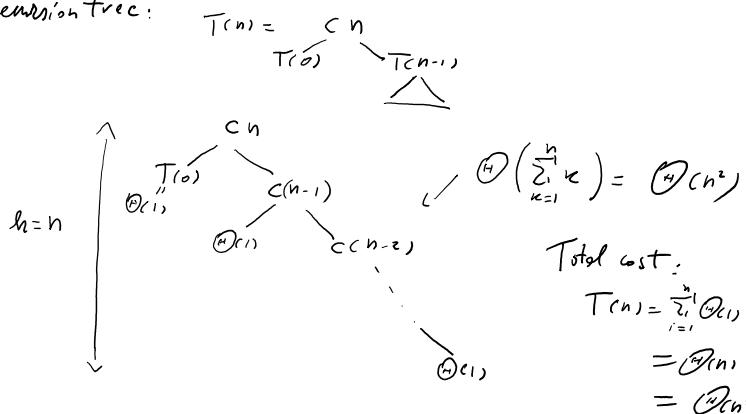
QUICKSORT(S, i, j) {
    if (i < j) {
        r ← pick the position of "a good pivot";
        S ← SWAP(S, r, i); // swap S[i] with S[r]
        p ← PARTITION(S, i, j); // new position of the pivot
        QViCKSORT(S', i, p-1);
        QViCKSORT(S', p+1, j);
    }
}
    
```

$\boxed{< | p | \geq}$

- Extra space: NO
- Time: $O(n \lg n)$ for balanced partitioning
 $O(n^2)$ for unbalanced //

Worst case: $T(n) = T(o) + T(n-1) + c_n$
 (unbalanced)

recursion tree:



Best case:

$$T(n) = 2T\left(\frac{n}{2}\right) + \Theta(n) = \Theta(n \lg n) \quad (\text{see Merge Sort})$$

Average Case:

$$\text{For } k=0, \dots, n-1 \quad X_k := \begin{cases} 1 & \text{if PARTITION generates } k:n-k-1 \text{ split} \\ 0 & \text{otherwise} \end{cases}$$

We assume $E[X_k] = 1 \cdot p(X_k=1) + 0 \cdot p(X_k=0) = p(X_k=1) = \frac{1}{n}$ (all splits equally likely)

$$T(n) = \begin{cases} T(o) + T(n-1) + \Theta(n) & \text{if } 0:n-1 \text{ split (if } X_0=1) \\ T(n-1) + T(o) + \Theta(n) & \text{if } n-1:0 \text{ split (if } X_{n-1}=1) \end{cases}$$

$$T(n) = \sum_{k=0}^{n-1} X_k \times (T(k) + T(n-k-1) + \Theta(n))$$

$$\begin{aligned} E[T(n)] &= E\left[\sum_{k=0}^{n-1} X_k \times (T(k) + T(n-k-1) + \Theta(n))\right] \\ &= \sum_{k=0}^{n-1} E[X_k] \times E[T(k) + T(n-k-1) + \Theta(n)] \\ &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(k)] + \frac{1}{n} \sum_{k=0}^{n-1} E[T(n-k-1)] + \frac{1}{n} \sum_{k=0}^{n-1} \Theta(n) \\ &= \frac{2}{n} \sum_{k=0}^{n-1} E[T(k)] + \Theta(n) = \frac{2}{n} \sum_{k=0}^{n-1} E[T(k)] + \Theta(n) \end{aligned}$$

$$= \frac{2}{n} \sum_{k=2}^{n-1} E[T_{(k)}] + O(n) = \frac{2}{n} \sum_{k=2}^{n-1} E[T_{(k)}] + O(n)$$

\uparrow
 $k=2, 4$
dominated by $O(n)$

Now, prove by induction that $E[T_{(n)}] \leq \alpha n \lg_2 n$ for const $\alpha > 0$
because $n=2$: choose α large enough $\underline{\alpha}$

$$\stackrel{n \rightarrow n+1}{E[T_{(n)}]} = \frac{2}{n} \sum_{k=1}^{n-1} E[T_{(k)}] + O(n) \stackrel{\substack{\text{use } \alpha \text{ chosen} \\ \text{the max } \alpha}}{\leq} \frac{2}{n} \sum_{k=2}^{n-1} \alpha k \lg_2 k + O(n)$$

\downarrow
inductive hypothesis

$$\begin{aligned} \sum_{k=2}^{n-1} \alpha k \lg_2 k &= \frac{1}{2} n^2 \lg_2 n - \frac{1}{8} n^2 \\ &\stackrel{\substack{\longrightarrow \\ \text{if } \alpha \text{ chosen} \\ \text{s.t. } \frac{\alpha n}{4} \text{ dominates } O(n)}}{\leq} \alpha n \lg_2 n - \left(\frac{\alpha n}{4} - O(n) \right) \end{aligned}$$

THREE-WAY-PARTITIONING-QUICKSORT(S, i, j):

{
 $i < j$:
 $(l, r) \leftarrow \text{THREE-WAY-PARTITION}(S, i, j);$
 $\text{THREE-WAY-PARTITIONING-QUICKSORT}(S, i, l-1);$
 $\text{THREE-WAY-PARTITIONING-QUICKSORT}(S, r, j)$

}

THREE-WAY-PARTITION(S, i, j)

$p \leftarrow S[i]; // \text{pivot element}$
 $l \leftarrow i;$

$r \leftarrow i+1;$

for ($c = r, \dots, j$) {

if ($S[c] == p$) {
 $\text{SWAP}(S[l], S[r]);$
 $r++;$

}

else if ($S[c] < p$) {

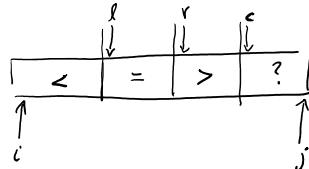
$\text{CYCLE}(S[l], S[r], S[c]);$

$r++;$

$l++;$

}

}



- Extra space: NO
- Time: $O(n \log n)$ on avg. (*)

(*) Thm Random selection of the pivot \Rightarrow the algorithm converges

(*) Thm Random selection of the pivot \Rightarrow the algorithm compares no more than $2n \ln n$ items on avg.

Proof Let be
 $RV \rightarrow X_{u,v} = \begin{cases} 1 & \text{if } S_{u,v} \text{ and } S_{v,u} \text{ are compared by PARTITION} \\ 0 & \text{otherwise} \end{cases}$

$$P_{u,v} = P(X_{u,v} = 1)$$

$$\xrightarrow{\substack{\text{total #} \\ \text{comparisons}}} E\left[\sum_{u,v} X_{u,v}\right] = \sum_u^1 \sum_{v>u}^1 E[X_{u,v}] = \sum_u^1 \sum_{v>u}^1 (1 \times p_{u,v} + 0 \times (1 - p_{u,v})) = \sum_u^1 \sum_{v>u}^1 p_{u,v}$$

- We have 3 cases:
- 1) $\begin{array}{l} \text{pivot} \in S_{u,v}, S_{v,u} \\ \text{OR} \\ S_{u,v}, S_{v,u} \subset \text{pivot} \end{array} \} \begin{array}{l} \text{same recursive} \\ \text{call} \end{array} \} \begin{array}{l} \text{no validity} \\ \text{for computation of } p_{u,v} \end{array}$
 - 2) $\begin{array}{l} S_{u,v} \subset \text{pivot} \subset S_{v,u} \\ \text{OR} \\ S_{v,u} \subset \text{pivot} \subset S_{u,v} \end{array} \} \begin{array}{l} \text{different recursive} \\ \text{call} \end{array} \} \rightarrow u, v \text{ NOT compared}$
 - 3) $\begin{array}{l} \text{pivot} = S_{u,v} \\ \text{OR} \\ \text{pivot} = S_{v,u} \end{array} \} \rightarrow u, v \text{ compared}$

Considering cases 2) and 3) and a bijective map $\phi: S \xrightarrow{\phi} S'_{\text{sorted}(S)}$
 $S_{u,v} \xrightarrow{\phi} S'_{u',v'}$

$$P_{u,v} = \frac{2}{v' - u' + 1} \quad \left(= \frac{\# \text{cases 2)}}{\# \text{cases 3}} \right)$$

$$\begin{aligned} \Rightarrow E\left[\sum_{u,v} X_{u,v}\right] &= \sum_{u=1}^n \sum_{v=u+1}^n p_{u,v} = \sum_{u=1}^n \sum_{v'=u+1}^n \frac{2}{v' - u' + 1} \\ &= 2 \sum_{u=1}^n \sum_{k=2}^{n-u+1} \frac{1}{k} \leq 2 \sum_{u=1}^n \underbrace{\sum_{k=2}^n}_{\sum_{k=1}^n \frac{1}{k} \leq 1 + \ln n} \frac{1}{k} \leq 2n \ln n \end{aligned}$$

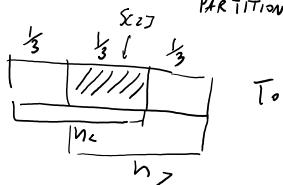
```
RANDSELECT(S, k) // Select k-th smallest item
r ← RAND(1, n);
(S<, S>) ← PARTITION(S, r); // item at position r is pivot
n< ← |S<|;
n> ← |S| - (|S<| + |S>|);
if (k ≤ n<)
    return RANDSELECT(S<, k);
else if (k ≤ n< + n>)
    return S>; // S> is in S<
else
    return RANDSELECT(S>, k - n< - n>); // skipping the items in S< US<
```

<ul style="list-style-type: none"> Extra space: NO Time: O(n) in avg. I/Os: $O(\frac{n}{\theta})$ 	(*)
---	-----

(*) Thm Selecting k -th smallest item in an unsorted sequence of size n takes $\Theta(n)$ avg time in RAM model and $O(\frac{n}{B}) I/Os$ in Disk model

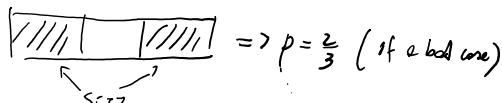
Proof let be $p_c = \text{prob. of go to left}$
 $p_r = " " " " \text{right}$ $T(n) = O(n) + p_c T(n_c) + p_r T(n_r)$

good case: $n_c, n_r \leq \frac{2}{3}n$:



To ensure this $S_{2,7}$ must be in $[\frac{2}{3}n, \frac{4}{3}n]$
 $\Rightarrow p = \frac{1}{3}$ (of a good case)

bad case: $n_c \text{ or } n_r > \frac{2}{3}n$:



$$\Rightarrow E[T(n)] \leq O(n) + \underbrace{\frac{1}{3} \cdot E[T(\frac{2}{3}n)]}_{\substack{\text{partition} \\ \text{good case}}} + \underbrace{\frac{2}{3} \cdot E[T(\frac{4}{3}n)]}_{\substack{\text{partition} \\ \text{bad case}}}$$

$$\frac{1}{3}E[T(n)] \leq O(n) + \frac{1}{3}E[T(\frac{2}{3}n)]$$

$$E[T(n)] \leq O(n) + E[T(\frac{2}{3}n)] = O(n)$$

on Disc model:

$$E[T(n)] \leq O(\frac{n}{B}) + E[T(\frac{2}{3}n)] = O(\frac{n}{B})$$

contradiction
of $S_{2,7} \approx S_3$ is a single pass

$$\begin{aligned} cn & \text{ cost} \\ \frac{cn}{3} &= \left[\sum_{k=1}^{\infty} \left(\frac{2}{3} \right)^k \right] \cdot cn \\ &\leq \left[\sum_{k=1}^{+\infty} \left(\frac{2}{3} \right)^k \right] \cdot cn \\ &= \frac{1}{1-\frac{2}{3}} \cdot cn = 3cn \\ &= O(n) \end{aligned}$$

Quick Sort is in place but we must consider the management of the recursive calls
#total recursive calls = $\begin{cases} O(\lg n) \text{ best case} \\ O(n) \text{ worst case} \end{cases}$ (For each rec. call $O(c_1)$ space cost)

So $O(n)$ extra working space. Solution:

BOUNDEDQS(S, i, j):

while ($j - i > n_0$) {

$r \leftarrow$ position of a good pivot;

$S \leftarrow SWAP(S, r, i);$

$p \leftarrow PARTITION(S, i, j);$

if ($p \leq \frac{i+j}{2}$) {

BOUNDEDQS($S, i, p-1$);

$i = p+1;$

} else {

BOUNDEDQS($S, p+1, j$);

$j = p-1;$

• Extra space: $O(\lg n)$ AVL

• Time: $O(n \lg n)$ average.

"elimination
of tail recursion"

BOUNDEDQS($S, p+1, j$);
 $j = p-1;$
 }
 }
 }
 INSERTION SORT(S, i, j);

(*) Thm On RAM model BOUNDED JS uses $O(\log n)$ extra working space.

Proof The recursive call is exerted on a subarray whose size is no more than half of the starting array. So the total recursive call = $O(g^{\lfloor \frac{n}{2} \rfloor})$

MULT-WAY QUICKSORT

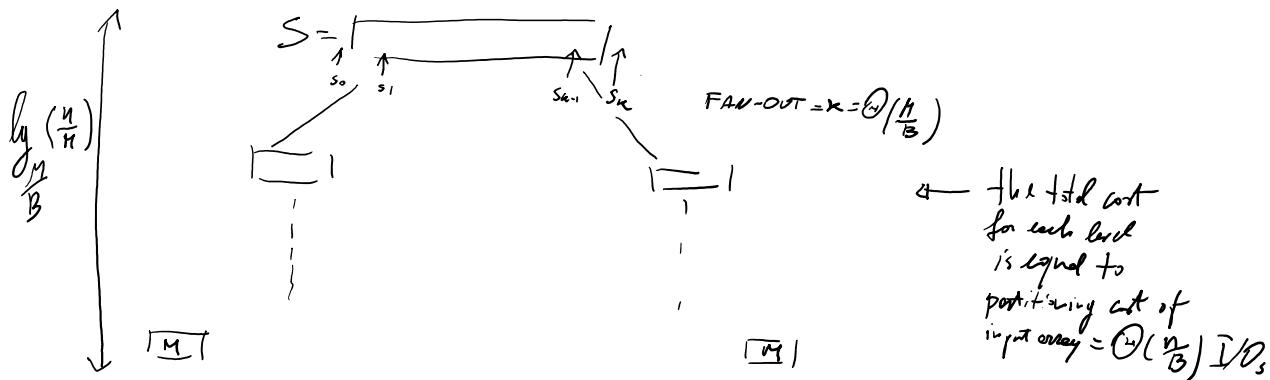
$k-1$ pivot s_1, \dots, s_{k-1} s.t. $\kappa = \Theta\left(\frac{M}{B}\right)$ $s_0 = 0, s_k = +\infty$
~~dummy pivot~~

$$B_i = \{S_{i,j} \mid s_{i-1} < S_{i,j} \leq s_i\}$$

We suppose good choice of points $\Rightarrow |B_i| = \Theta\left(\frac{n}{k}\right)$

The partition can be implemented in $O(\frac{4}{3})$ using 1 input block (to read input seq.) and 1 output block (to write its n partition)

Using recursive tree:



GOOD SELECTION of $\kappa-1$ PIVOTS via OVERSAMPLING ($\alpha \geq 0$)
 $A \leftarrow \text{SAMPLING}(S, (\alpha+1)\kappa-1)$; \uparrow
int Time: $O(\kappa \lg^2 \kappa)$ (**)
 $A \leftarrow \text{SORT}(A)$; tuning $\alpha = \Theta(\lg \kappa)$

$A \leftarrow \text{SORT}(A)$;

for $i=1, \dots, k-1$

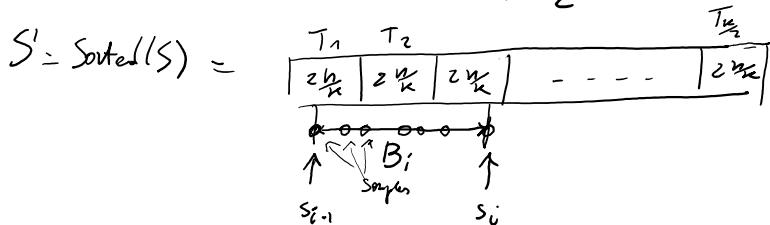
OUTPUT ($s_i = A[(\alpha+1)i]$);

(*) If $\alpha = \Theta(\lg n)$ over sampling time $O(k \lg^2 n)$. Prof time $\text{sort } A = O(\alpha \lg(\alpha n)) = O(\lg n \cdot \alpha \cdot \lg(\lg n \cdot n))$
 $= O(\lg n \cdot n \cdot [\lg(\lg n) + \lg n]) = O(n \cdot \lg^2 n)$

Lemma Let $k \geq 2$ and $\alpha+1 = 12 \lg k \Rightarrow |B_i| < \frac{4n}{k}$ with probability $\geq \frac{1}{2}$

Proof By contradiction:

We suppose $P(\exists B_i \mid |B_i| \geq \frac{4n}{k}) > \frac{1}{2}$



$$\begin{aligned} P(\exists B_i \mid |B_i| \geq \frac{4n}{k}) &\leq P(\exists T_j \mid T_j \subset B_i) = P(\exists T_j \mid T_j \text{ contains } \leq (\alpha+1) \text{ samples}) \\ &= \sum_{j=1}^{k-1} P(T_j \text{ contains } \leq (\alpha+1) \text{ samples}) \\ &\stackrel{\text{Union Bound}}{\leq} \frac{k}{2} P(T_j \text{ contains } \leq (\alpha+1) \text{ samples}) \end{aligned}$$

(*) $P(\text{a sample is in } T_j) = \frac{2n/k}{n} = \frac{2}{k}$, let $X = \# \text{samples in } T_j$
 $P(X_i)$ $\xrightarrow{\text{uniform sampling}}$ $X_i = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ sample is in } T_j \\ 0 & \text{otherwise} \end{cases}$

$$\begin{aligned} E[X] &= E[\sum_{i=1}^{\# \text{sample}} X_i] = \sum_{i=1}^{\# \text{sample}} E[X_i] = \sum_{i=1}^{\# \text{sample}} p(X_i=1) = \underbrace{[(\alpha+1)k - 1]}_{\# \text{sample}} \frac{2}{k} \\ \text{So } E[X] &\approx 2(\alpha+1) - \frac{2}{k} \geq 2(\alpha+1) - 1 \geq \frac{3}{2}(\alpha+1) \end{aligned}$$

$$\text{So } \alpha+1 \geq (1 - \frac{1}{3}) E[X] \quad (*)$$

$$\Rightarrow P(X < \alpha+1) \leq P(X < (1 - \frac{1}{3}) E[X]) \leq e^{-\frac{(1/3)^2 \cdot E[X]}{2}} = e^{-\frac{E[X]}{18}} \leq e^{-\frac{[\frac{3}{2}(\alpha+1)]/18}{2}} \quad (*)$$

$$= e^{-\frac{(\alpha+1)/12}{2}} = e^{-\frac{\alpha+1}{24}} = \frac{1}{k}$$

$\uparrow \alpha = 12 \lg k - 1$

Finally:

$$P(\exists B_i \mid |B_i| \geq \frac{4n}{k}) \leq \frac{k}{2} \cdot \frac{1}{k} = \frac{1}{2}$$

Disk model (known seq.)

lunedì 13 febbraio 2023 20:11

$S_{(1,h)} = s_1, \dots, s_n \quad (m < n)$
stored on disk
Aux $S'_{(1,n)} =$ 

```

 $S'_{(1,n)} \leftarrow S_{(1,n)}$ 
for  $s=0,1,\dots,m\}$ 
     $p \leftarrow \text{RANDOM}(1, n-s)$ ;
     $\text{SELECT}(S'[p])$ ;
     $\text{SWAP}(S'[p], S'[n-s])$ ;
}

```

- Extra space: $O(n \log n)$ bits $\leftarrow S'$ array of positions
- Time:
- I/Os: $\Theta(m) \text{ I/Os} \leftarrow$ access in all random way & swap items

Dictionary

```

 $\emptyset \leftarrow \emptyset$ 
while ( $|\emptyset| < m$ ) {
     $p \leftarrow \text{RANDOM}(1, h)$ ;
    if  $p \notin \emptyset$ 
         $\text{INSERT}(p, \emptyset)$ ;
}

```

Fact $P(p \in \emptyset) < \frac{1}{2}$ and Oois re-sampling strategy.

Dim $P(p \in \emptyset) = |\emptyset| \leq \frac{m}{h} < \frac{1}{2}$
s.p.g. $m < \frac{h}{2}$

$$E[\text{#resampling}] < \sum_{\ell=1}^{\infty} \left(\frac{1}{2}\right)^{\ell} \cdot \ell = 2$$

- Extra space: $O(m)$
- Time: $O(m)$ in avg.
- I/Os: $O(\min\{m, \frac{m}{B}\}) \text{ I/Os}$

Using Hash table with chaining

Sorting

```

 $\emptyset \leftarrow \emptyset$ 
while ( $|\emptyset| < m$ ) {
     $X \leftarrow \text{RANDOM}(1, n), m$ ;
     $X \leftarrow \text{SORT}(X)$ ;
     $\emptyset \leftarrow \text{REMOVEDUPLICATES}(\emptyset)$ ;
}

```

- Extra space: $O(m)$
- Time: $O(m \log m)$ avg.
- I/Os: $O(\min\{m, \frac{m}{B}\}) \text{ I/Os}$
if $m \leq B$ small enough
to map sampled pos. in I.M.

Prop $P(\text{No duplicates in } \emptyset) \leq e^{-m(m-1)/2n}$

Dim Birthday problem:

$$\begin{aligned}
 P(\text{No duplicate in } S) &= \frac{m! \binom{n}{m}}{n^m} = \frac{n(n-1)\dots(n-m+1)}{n^m} = 1 \times \left(\frac{n-1}{n}\right) \times \dots \times \left(\frac{n-m+1}{n}\right) \\
 &= 1 \times \left(1 - \frac{1}{n}\right) \times \dots \times \left(1 - \frac{m-1}{n}\right) \leq e^0 \times e^{-\frac{1}{n}} \times \dots \times e^{-\frac{(m-1)}{n}} = e^{-\sum_{i=1}^{m-1} \frac{i}{n}}
 \end{aligned}$$

Corollary If $m = \sqrt{n}$ $P(\text{Duplicate in } S) \approx 1 - \frac{1}{\sqrt{e}} \approx 0.9$ ($= P(\text{re-sampling})$)

Dim $P(\text{No dup. in } S) \leq e^{-\frac{1}{2}} = \frac{1}{\sqrt{e}}$

Streaming model (known seq.)

lunedì 13 febbraio 2023 21:57

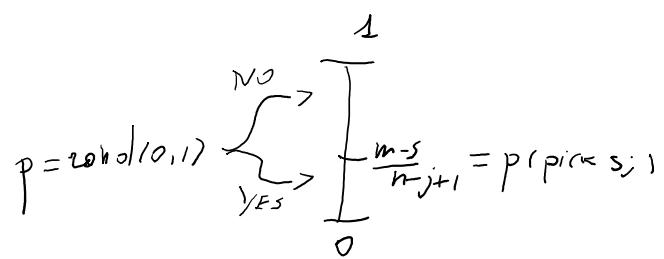
$$\{s_1, \dots, s_n\} = s_1, \dots, s_n, m < n$$

SCANNING & SELECTING (S, m):

```

 $s \leftarrow \emptyset$  // # selected
for( $j=1$ ;  $j \leq n$  &  $s \neq m$ ;  $j++$ ) {
     $p \leftarrow \text{RAND}(0,1)$ ;
    if  $p < \frac{m-s}{n-j+1}$  {
        select( $s[j]$ );
         $s++$ ;
    }
}

```



- Extra space: $O(m)$
- Time: $O(n)$
- I/O: $O\left(\frac{n}{B}\right)$ I/Os

CASE $m=1$: $p(\text{pick } s_j) = \frac{1}{n-j+1} \Rightarrow p(\text{sample } s_j) = \frac{1}{n}$

Proof By induction:

Suppose $p(\text{sample } s_i, i \in \{1, \dots, j-1\}) = \frac{1}{n}$

$$\begin{aligned}
 p(\text{sample } s_j, j \in \{1, \dots, n\}) &= p(\text{not sample } s_1, \dots, s_{j-1}) \times p(\text{pick } s_j) = \\
 &= 1 - p(\text{sample } s_1 \text{ or } \dots \text{ or } s_{j-1}) \times p(\text{pick } s_j) = \\
 &= \left(1 - \frac{j-1}{n}\right) \times \frac{1}{n-j+1} = \frac{1}{n}
 \end{aligned}$$

□

CASE $m \geq 1$:

$$p(\text{pick } s_j) = \frac{m-s}{n-j+1}$$

$$\begin{aligned}
 \text{Prof: } p(\text{pick } s_j) &= \frac{\# \text{ elem. to pick}}{\# \text{ remain. pos. to scan}} = \frac{|\{s' \subseteq S[j, n], |s'| = m-s, s_j \in s'\}|}{|\{s \subseteq S[j, n], |s| = m-s\}|}
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{\binom{n-j}{m-s-1}}{\binom{n-j+1}{m-s}} = \frac{\frac{(n-j)!}{(m-s-1)!}}{\frac{(n-j+1)!}{(m-s)!}} = \frac{\frac{(n-j)!}{(m-s-1)!(n-j-m+s+1)!}}{\frac{(n-j+1)!}{(m-s)!(n-j+1-m+s)!}} = \frac{m-s}{n-j+1}
 \end{aligned}$$

□

Streaming model (unknown seq.)

lunedì 27 febbraio 2023 17:29

$$S = s_1, \dots, s_n, \dots$$

$$\text{min-Heap } H = \left\{ \langle r_i, \text{item} \rangle \mid \begin{array}{c} \uparrow \\ \text{random} \\ \in [0,1] \end{array} \right\} \subseteq S$$

$$\text{START: } H = \{ \langle -\infty, \phi \rangle, \dots, \langle -\infty, \phi \rangle \}, |H| = m$$

```

for j = 1, ..., n {
    r_j ← RAND(0,1);
    y ← MIN(H);
    if (r_j > y) {
        POP_MIN(H);
        INSERT(<r_j, s[j]>, H);
    }
}
return H;

```

- Extra space: $O(m)$
- Time: $O(n \lg m)$ ($O(\lg m)$ to insert item in H)
- I/Os: $O(\frac{n}{B})$ I/Os

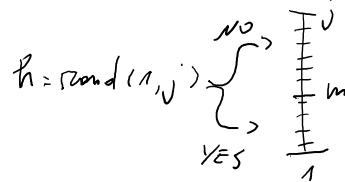
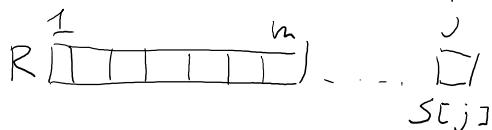
$$\text{START: } R[1, m] = S[1, m]$$

RESERVOIR SAMPLING:

```

for j = m+1, ..., n, ... {
    h ← RAND(1, j); // rand integer
    if (h ≤ m)
        R[h] ← S[j];
}
return R;

```



- Extra space: $\Theta(m)$
- Time: $O(n)$
- I/Os: $O(\frac{n}{B})$ I/Os

$$\text{Prop } P(S_j \in R \text{ after } n \text{ steps}) = \frac{m}{n} \quad (\forall 1 \leq j \leq n)$$

$$\text{Proof Case } n=m: P(S_j \in R) = 1 = \frac{m}{n}$$

$$\text{Start: } R[1, m] = S[1, m]$$

$$\text{Case } n > m: P(S_j \in R \text{ after } n \text{ steps}) = \prod_{i=1}^{n-1} (S_i \text{ go inside } R \text{ at } i\text{-th step with prob. } = \frac{m}{n})$$

$$P(S_j \in R \text{ after } n-1 \text{ steps}) \times [P(S_n \text{ not picked}) + P(S_n \text{ picked}) \times P(S_j \text{ not removed})] =$$

$$\begin{aligned}
 & \frac{m}{n-1} \times \left[\left(1 - \frac{m}{n}\right) + \frac{m}{n} \times \frac{m-1}{m} \right] = \\
 & \text{inductive Step} \quad P(S_n \text{ picked}) \quad P(\text{remove } S_i \neq S_j \text{ from } R) \\
 & = \sum_{S_i \neq S_j} P(\text{remove } S_i \text{ from } R) = \sum_{i \neq j} \frac{1}{m} = \frac{m-1}{m}
 \end{aligned}$$

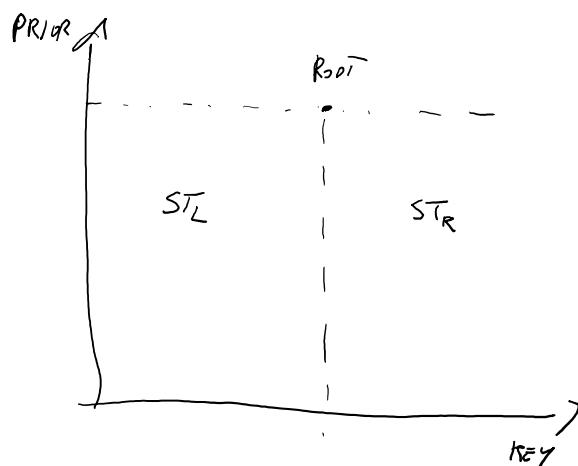
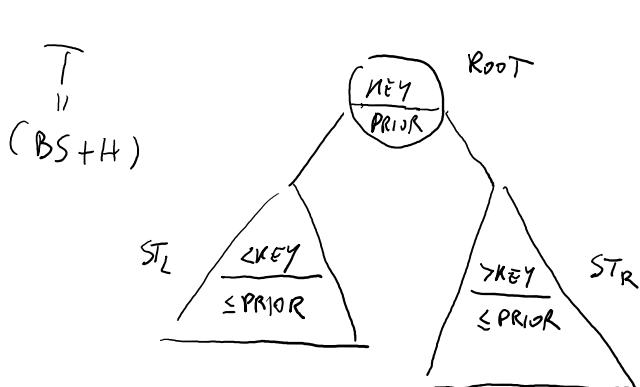
$$= \frac{m}{n-1} \times \left[1 + \left(\frac{m-1}{m} - 1 \right) \times \frac{m}{n} \right] = \frac{m}{n-1} \times \left[1 - \frac{1}{n} \right] = \frac{m}{n-1} \times \frac{n-1}{n} = \frac{m}{n}$$

Treap

lunedì 13 febbraio 2023 20:09

Def Treap T is a binary tree based data structure. Each node $u \in T$ is given by a search key and a priority $u = (\text{search key}, \text{priority})$. The treap is simultaneously a binary searchtree for the search key, and a max or min heap for the priority of the node.

Node = (Search key , pr. vity)
 ↑
 ↑
 BIN. SEARCH
 TREAP MW/MAX
 HEAP



Operations

- Search($\text{key} = k$)
- Insert ($\text{node} = z$)
- Delete ($\text{node} = z$)
- Split ($T, \text{key} = k$) = return 2 treaps T_L, T_R such that each key (T_L) $<$ key (T_R)
- Merge (T_L, T_R) = return a treap $T = T_L \cup T_R$ ($\text{key}(T_L) < \text{key}(T_R)$)

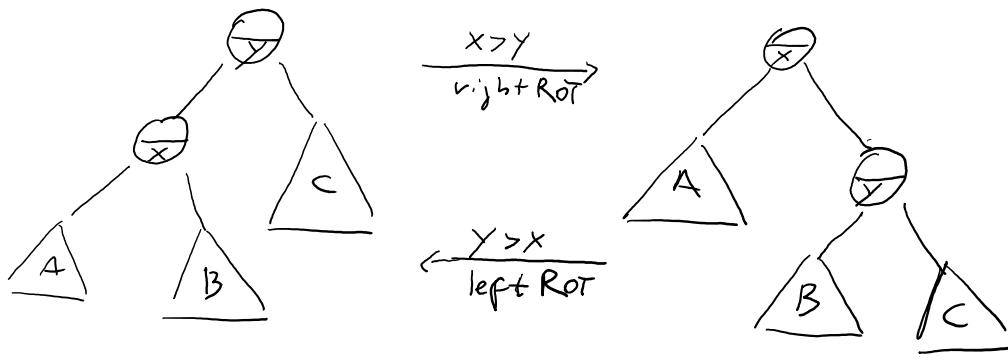
SEARCH($T, \text{key} = k$)

return SEARCH(BST(T), k); // search k as in bin. searchtree

INSERT($T, \text{node} = z$)

```

    INSERT(BST( $T$ ), key(node));
    while (!HEAP( $T$ )) { // the priority doesn't satisfy the heap
        ROTATION( $T, z$ );
    }
    return  $T$ ;
  
```



$\text{DELETE}(T, \text{node} = z)$

```

while ( $z \neq \text{LEAF}(T)$ ) {
    ROTATION( $T, z$ ); // rotate with children
    }
    CUT( $T, z$ );
return  $T$ ;

```

• Time = $O(n)$ worst case (*),
 $O(\lg n)$ expected case (*).

$\text{SPLIT}(T, \text{key} = k)$

```

INSERT( $T, \text{node} = (k, -\infty / +\infty)$ ); // return a tree p
 $T_L \leftarrow \text{LEFT-SUBTREAP}(T, \text{root})$ . with root =  $(k, \pm \infty)$ 
 $T_R \leftarrow \text{RIGHT-SUBTREAP}(T, \text{root})$ .
return  $(T_L, T_R)$ .

```

$\text{MERGE}(T_L, T_R)$

```

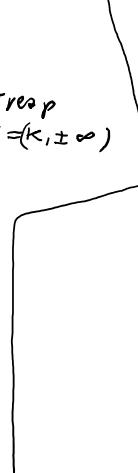
dummy node  $\leftarrow (? , -\infty / +\infty)$ ;
 $T \leftarrow \text{TREAP}(\text{Root} = \text{dummy node},$ 
    LEFT-SUBTREAP =  $T_L$ ,
    RIGHT-SUBTREAP =  $T_R$ ).

```

```

while ( $\text{dummy node} \neq \text{LEAF}(T)$ )
    ROTATION( $T, \text{dummy node}$ );
CUT( $T, \text{dummy node}$ );
return  $T$ ;

```



Thm (*) All the operation on tree take time $O(n)$ time

Proof The time need by the operation is $\Theta(\text{depth}(\text{node}))$ where node is the node that must be inserted, deleted, searched or used by split or merge.

$\Theta(\text{depth}(\text{node})) = O(h) = O(n)$ (the tree is not necessary balanced)

\uparrow
node = leaf worst case

71

→ up ↑ ↓ worst case
 node = leaf worst case
 (the key is not necessary necessary)

□

Random prioritization

$$\text{PRIOR}_i = \text{PRIOR}(\text{node}_i) \quad \text{PRIOR}_1, \dots, \text{PRIOR}_n \text{ i.i.d. RV, } \text{PRIOR}_i = \text{rand}(0,1)$$

Thm (*): Assumption min-tree, x_n generic node of T . $E[\text{depth}(x_n)] = O(\lg n)$

Proof x_n s.t. $x_n \in \text{NODES}(T) = \{x_1, \dots, x_n\}$

$$\text{KEY}(x_1) < \dots < \underset{\substack{\uparrow \\ k\text{-th smallest}}}{\text{KEY}(x_n)} < \dots < \text{KEY}(x_n)$$

- I.V. $A_n^i := \begin{cases} 1 & \text{if } x_i \text{ is predecessor of } x_n \\ 0 & \text{otherwise} \end{cases}$

$$\text{depth}(x_n) = \sum_{i=1}^n A_n^i, \quad E[\text{depth}(x_n)] = \sum_{i=1}^n P(A_n^i = 1)$$

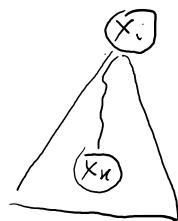
A_n^i is I.V.

- Def $X_{(i,n)} = X_{(n,i)} := \begin{cases} \{x_1, \dots, x_n\} & \text{if } i < n \\ \{x_n, \dots, x_i\} & \text{if } i > n \end{cases}$

Lemma Now we prove that: $A_n^i = 1 \Leftrightarrow \text{PRIOR}(x_i) = \min \{\text{PRIOR}(X_{(i,n)})\}$

→ we have 4 possible case of tree

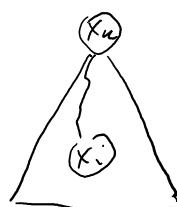
1) x_i is root



$$\text{PRIOR}(x_i) = \min \text{PRIOR}(X_{(i,n)})$$

$$A_n^i = 1 \quad \text{on}$$

2) x_n is root



$$\text{PRIOR}(x_i) \neq \min \text{PRIOR}(X_{(i,n)}) \\ (= \text{PRIOR}(x_n))$$

$$A_n^i = 0 \quad \text{on}$$

3) x_j root



or



$$\text{PRIOR}(x_i)$$

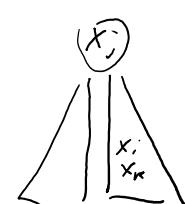
$$\neq \min \text{PRIOR}(X_{(i,n)}) \\ (= \text{PRIOR}(x_j))$$

$$A_n^i = 0 \quad \text{on}$$

4) x_j root



or



Recursively on the subtree

$$A_x^i = 0 \text{ or } (= \text{PRIOR}(X_j))$$

7

$\text{PRIOR}_1, \dots, \text{PRIOR}_n$ i.i.d.

$$\Rightarrow P(\text{PRIOR}_i = \min(\text{PRIOR}(X_{(i,k)}))) \stackrel{(*)}{=} \frac{1}{|k-i|+1} \quad \text{uniform distribution}$$

Anal

$$\begin{aligned} p(A_n^i = 1) &\stackrel{\text{Lemma}}{=} [i \neq n] \cdot p(\text{PRIOR}_i = \min(\text{PRIOR}(X_{(i,k)}))) = \\ &= [i \neq n] \cdot \frac{1}{|k-i|+1} = \begin{cases} \frac{1}{k-i+1} & \text{if } i < k \\ 0 & \text{if } i = k \\ \frac{1}{i-n+1} & \text{if } i > k \end{cases} \end{aligned}$$

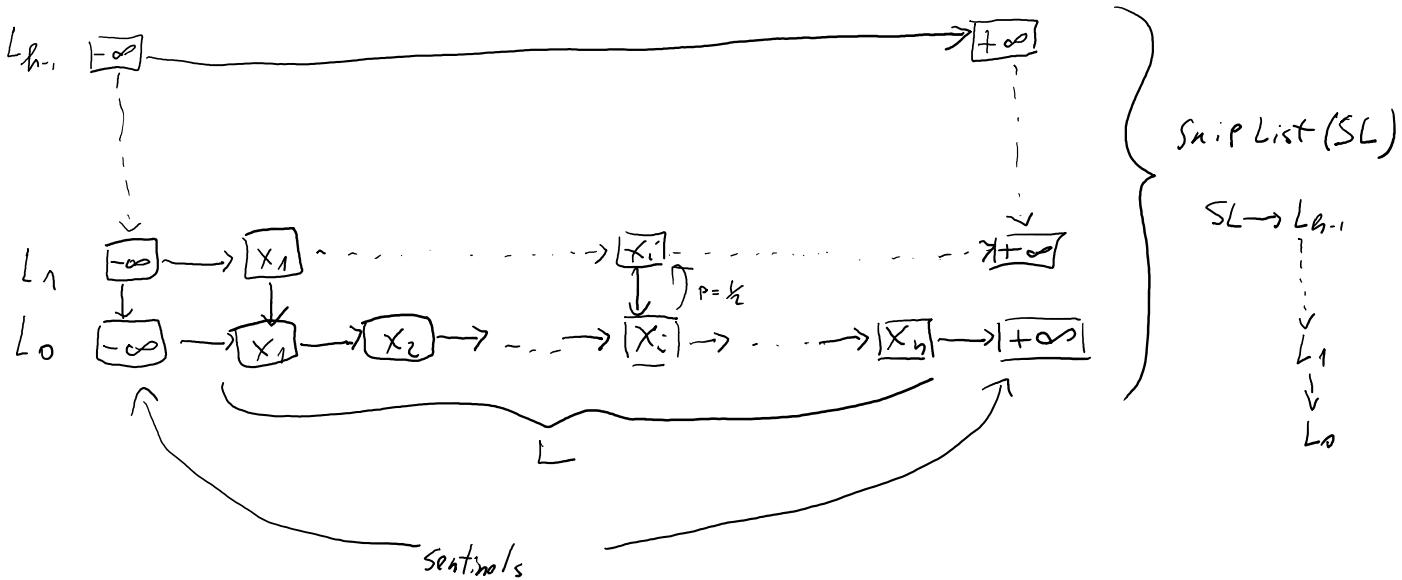
Therefore,

$$\begin{aligned} E[\text{depth}(x_n)] &= \sum_{i=1}^n p(A_n^i = 1) = \sum_{i < n} \frac{1}{k-i+1} + \sum_{i > n} \frac{1}{i-n+1} \\ &= \sum_{j=2}^n \frac{1}{j} + \sum_{j=2}^{n-k+1} \frac{1}{j} = \sum_{j=1}^n \frac{1}{j} - 1 + \sum_{j=1}^{n-k+1} \frac{1}{j} - 1 < \ln(n) + \ln(n-k+1) - 2 \\ &< 2\ln(n) - 2 = O(\ln n) \end{aligned}$$

Skip List

giovedì 23 marzo 2023 00:54

list $L = x_1 < \dots < x_n$



BUILD_SKIPLIST(L)

```

 $L_0 \leftarrow L;$ 
 $L_0 \leftarrow \text{ADD\_SENTINELS}(L_0);$  // add -∞ and +∞
 $SL \leftarrow L_0;$  // initialise Skip List
 $TO\_COPY \leftarrow \text{ITEM}(L_0);$  //  $L_0$  without sentinels
 $\text{while } (TO\_COPY \neq \emptyset) \{$ 
     $L_i \leftarrow \emptyset;$ 
     $\text{for each } item \text{ in } TO\_COPY \{$ 
         $p \leftarrow \text{RAND}(0,1);$ 
         $\text{if } p < 0.5$ 
             $\text{ADD}(\text{DUPLICATE}(item), L_i);$ 
    }
     $L_i \leftarrow \text{ADD\_SENTINELS}(L_i);$ 
     $SL \leftarrow \text{ADD}(SL, L_i);$ 
     $TO\_COPY \leftarrow \text{ITEM}(L_i);$ 
     $i++;$ 
}

```





$\text{SEARCH}(v, SL)$

```

 $v \in SL$ ; // point first item, fLh
while ( $v \neq \text{NULL}$  AND  $\text{key}(v) \neq k$ )
    if  $\text{key}(\text{RIGHT}(v)) < k$ 
         $v \leftarrow \text{RIGHT}(v)$ ;
    else
         $v \leftarrow \text{DOWN}(v)$ ;
}
return v;

```

Time: $O(\lg n)$ with high prob. (e)

Analysis:

Def 1 $h :=$ levels of SL

Def 2 Given $-\infty < x < \infty$, $C(x) :=$ #occ. of x in $L_1 \cup \dots \cup L_{n-1}$ (not L_n)
 $=$ #time x is copied

Lemma $\forall \text{key } x \text{ in } SL \quad E[C(x)] = 1$

$$\begin{aligned}
 &\text{Pf Recursively } E[C(x)] = p(\text{not copy } x) \cdot 0 + p(\text{copy } x) \cdot [E[C(x)] + 1] \\
 \Leftrightarrow E[C(x)] &= \frac{1}{2} \cdot 0 + \frac{1}{2} E[C(x)] + \frac{1}{2} \\
 \Leftrightarrow E[C(x)] &= 1
 \end{aligned}$$

Or equivalently $C(x) =$ geometric distribution $C(x) =$ #failures (copy) before success (not copy)

$$\begin{aligned}
 E[C(x)] &= \sum_{k=0}^{+\infty} k \cdot p(C(x)=k) = \sum_{k=0}^{+\infty} k \cdot (1-p)^k \cdot p & p = p(\text{not copy } x) \\
 &= \frac{1-p}{p} = \frac{1-\frac{1}{2}}{\frac{1}{2}} = 1 & 1-p = p(\text{copy } x)
 \end{aligned}$$

see Ge. Dist.

Thm $p(h \geq c \cdot \lg n) \leq \frac{1}{c}$, ($h = O(\lg n)$ with high prob.)

Thm $p(h \geq c \cdot \lg n) \leq \frac{1}{n^{c-1}}$ ($h = O(\lg n)$ with high prob)

Proof $p(h \geq l) = p((cx_1) \geq l \vee \dots \vee (cx_n) \geq l)$

$$\stackrel{\text{Union Bound}}{\leq} \sum_{i=1}^n p((cx_i) \geq l) = h \cdot p((cx) \geq l) \stackrel{\text{i.i.d.}}{\leq} n \cdot \frac{1}{2^l}$$

at least
l copy

$$\therefore \underbrace{p(h \geq l)}_{(*)} \leq \frac{n}{2^l} \Rightarrow p(h \geq c \lg n) \leq \frac{n}{2^{c \lg n}} \leq \frac{n}{n^c} = \frac{1}{n^{c-1}}$$

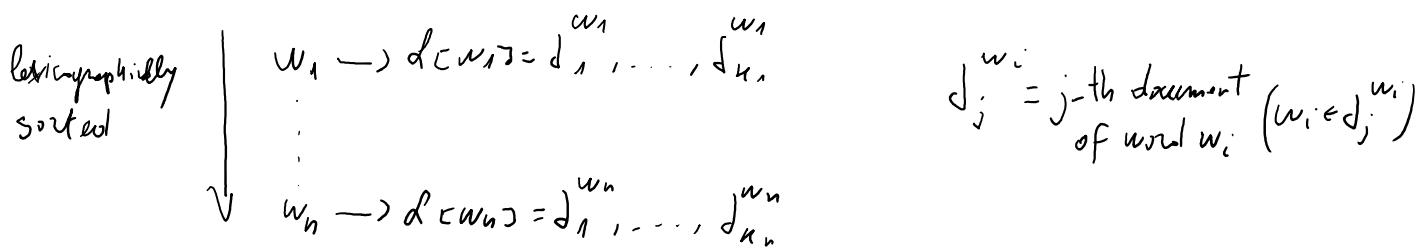
Thm $E[h] \leq \lg n + 1$ (in avg. at most 1 levels more than ideal)

$$\begin{aligned} \text{Proof } E[h] &= \sum_{l=0}^{+\infty} l p(h=l) = \sum_{l=1}^{\infty} p(h \geq l) \\ &\stackrel{\text{using IV}}{=} \sum_{l=1}^{\lg n} p(h \geq l) + \sum_{l=\lg n+1}^{+\infty} p(h \geq l) \\ &= \underbrace{\sum_{l=1}^{\lg n} p(h \geq l)}_{\text{using (*)}} + \underbrace{\sum_{i=1}^{\lg n} \frac{n}{2^{i \lg n}}}_{\substack{\text{substitute } l=i+\lg n}} = \lg n + \sum_{i=1}^{\lg n} \frac{1}{2^i} = \lg n + \left(\frac{1}{1-\frac{1}{2}} - 1\right) = \lg n + 1 \end{aligned}$$

Thm Search exec in skip list (copying with prob. $\frac{1}{2}$) takes $O(\lg n)$ with high probability

Proof Search takes $\Theta(h)$ time. $\Theta(h)$ can be bigger also more than $\Theta(n)$. But by thm before h is $O(\lg n)$ with high prob. $\Rightarrow O(\lg n)$ with high probability.

Inverted list



Query(w_1, w_2) \longrightarrow $\text{Intersection}(d(w_1), d(w_2))$

$\overset{\text{A}}{\parallel} \quad \overset{\text{B}}{\parallel} \leftarrow \text{set}$

We suppose A and B sorted

$$\begin{cases} \text{if not sorted intersection time} = \begin{cases} O(m \times n) \\ O(m \lg m + h \lg n + \dots) \end{cases} \\ \text{sort(A)} \quad \text{sort(B)} \quad \text{intersection(A, B)} \end{cases}$$

$$A = a_1, a_2, \dots, a_n \quad a_i < a_{i+1},$$

$$B = b_1, b_2, \dots, b_m \quad b_i < b_{i+1},$$

Cost $m \sim n$

MERGE-BASED INTERSECTION (A, B):

```

 $i \leftarrow 1; j \leftarrow 1;$ 
 $a \leftarrow A[i]; b \leftarrow B[j];$ 
while ( $i < n$  OR  $j < m$ ) {
    if ( $a = b$ ) {
        OUTPUT(a);
        if ( $i < n$ )  $a \leftarrow A[i+1];$ 
        if ( $j < m$ )  $b \leftarrow B[j+1];$ 
    }
    else if ( $a < b$  AND  $i < n$ )
         $a \leftarrow A[i+1];$ 
    else if ( $j < m$ )
         $b \leftarrow B[j+1];$ 
}

```

- Extra space: NO
- Time: $O(m+n)$ Optimal if $n = \Theta(m)$ ($O(m)$)
- I/Os: $O(\frac{m}{B} + \frac{n}{B})$ Optimal if $n = \Theta(m)$ ($O(\frac{m}{B})$)
- (+) cache friendly
- (-) bad if $m \ll n$
- (+) optimal if $m = n$

Case in con

BINARY SEARCH INTERSECTION (A, B)

for $i = 1, \dots, n$

$\text{BINARY_SEARCH}(A, B[i])$; // binary search b_i in A

- extra space: NO
- Time: $O(n \lg n)$
- I/O: $O(n_1 \lg \frac{n}{n_2})$
- (-) Not cache friendly
- (+) optimal if $n \ll n_1$

Lower Bound (General bound)

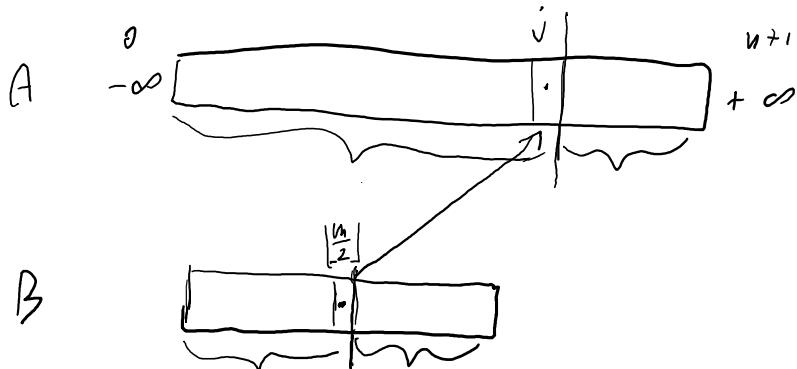
Supposing $B \subseteq A$ ($m < n$)



$$\Omega(\lg \binom{n}{m}) = \Omega(m \lg \frac{n}{m})$$

Mutual partitioning

lunedì 6 marzo 2023 11:52



MUTUAL PARTITIONING(A, B)

$n \leftarrow \text{LENGTH}(A); m \leftarrow \text{LENGTH}(B);$

$p \leftarrow B[\lfloor \frac{m}{2} \rfloor]; // \text{if } n < m \text{ exchange the role of } A \text{ and } B$

$i \leftarrow \text{BINARYSEARCH}(p, A); // a_i < b_{\lfloor \frac{m}{2} \rfloor} < a_{i+1}$

$\text{if } p = A[i]$

$\text{OUTPUT}(p);$

MUTUAL PARTITIONING(A[1:i], B[1, $\lfloor \frac{m}{2} \rfloor - i$]); $// a_i \text{ is discarded if just outputted}$
 $, , (A[i+1:n], B[\lfloor \frac{m}{2} \rfloor + 1, m]),$

- Extra space: NO
- Time: $O(m(1 + \lg \frac{n}{m}))$
- (+) combine the last of the merge and bin. search.
- (-) no cache friendly

fact time is $O(m(1 + \lg \frac{n}{m}))$

Proof $T(n, m) = O(\log n) + T(n_1, \frac{m}{2}) + T(n_2, \frac{m}{2}) \leq O(\lg n) + 2T(\frac{n}{2}, \frac{m}{2})$

$\uparrow \quad \uparrow$
BS balanced core (worst case)

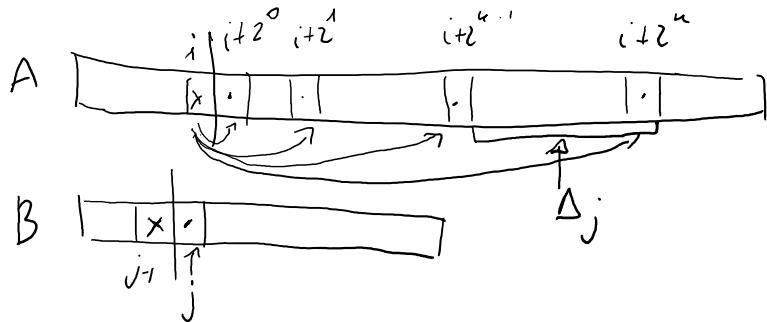
and $T(n, m) = 1$ if $n, m \leq 1$

$\Rightarrow T(n, m) = O(m(1 + \lg \frac{n}{m}))$ for any $m \leq n$

\uparrow
not proved

Doubling searching

lunedì 6 marzo 2023 12:37



Doubling search intersection (A, B)

$m \leftarrow |B|; n \leftarrow |A| // m < n$ otherwise exchange
 $i \leftarrow 0.$

for $j = 1, \dots, m \{$
 $k \leftarrow 0;$

while ($i + 2^k \leq n$ AND $B[i:j] > A[i+2^k:j]$)
 $k \leftarrow k + 1;$

$i' \leftarrow \text{BINARY SEARCH}(B[i:j], A[i+2^{k-1}+1, \min\{i+2^k, n\}]); // BS is A[i:i']$

if ($A[i'] = B[i:j]$)
 $\text{OUTPUT}(B[i:j]);$

$i \leftarrow i';$

}

- * extra space: NO
- * Time: $O(m(1 + \lg \frac{n}{m}))$ (*)
- (+) Optimal in comparison model
- (+) NO recursions
- (-) jumping

(*) fact Time is $O(m(1 + \lg \frac{n}{m}))$. Optimal for comparison model.

Proof Let i_j the position of b_j in A (from BS).

$$A[i_{j-1}+2^{k-1}] < b_j \leq A[i_{j-1}+2^k] \Rightarrow 2^{k-1} < i_j - i_{j-1} \leq 2^k \quad (1)$$

$$(2) \Delta_j = \min\{2^{k-1}, n\} \quad (\text{size of subarray of } A \text{ for BS}).$$

$$\Rightarrow \Delta_j \leq 2^{k-1} \leq i_j - i_{j-1} \Rightarrow \Delta_j < i_j - i_{j-1} \text{ and } \sum_{j=1}^n \Delta_j < \sum_{j=1}^n (i_j - i_{j-1}) = i_n - i_0 \leq n$$

↑ telescopic sum

$$T(n) = O\left(\sum_{j=1}^m (1 + \underbrace{\lg \Delta_j}_{BS})\right) = O(m + \sum_{j=1}^m (\lg \Delta_j)) = O(m(1 + \frac{1}{m} \sum_{j=1}^m (\lg \Delta_j))) =$$

(*) → $= O(m(1 + \lg (\frac{\sum_{j=1}^m \Delta_j}{m}))) =$

while
statement 10)

$$\begin{aligned}
 (*) &\longrightarrow = O\left(m\left(1 + \lg\left(\sum_{j=1}^m \frac{x_j}{m}\right)\right) = \right. \\
 &= O\left(m\left(1 + \lg\frac{n}{m}\right)\right)
 \end{aligned}$$

10

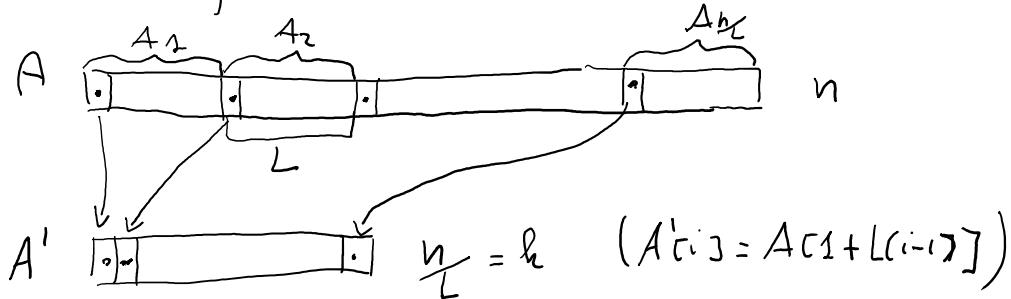
(*) Jensen inequality : $\begin{cases} \varphi\left(\frac{\sum_{j=1}^m x_j}{m}\right) \leq \frac{1}{m} \sum_{j=1}^m \varphi(x_j) & \text{if } \varphi \text{ is convex} \\ \varphi\left(\frac{\sum_{j=1}^m x_j}{m}\right) \geq \frac{1}{m} \sum_{j=1}^m \varphi(x_j) & \text{if } \varphi \text{ is concave} \end{cases}$ (*)

also for
 all convex
 linear combination
 of x_1, \dots, x_m

2-level storage approach

lunedì 6 marzo 2023 16:39

Let L size of block



2LEVEL INTERSECTION(A, A', B)

$(B_1, \dots, B_h) \leftarrow \text{MERGE}(A', B)$; // $B_i = \{B_{rj} \mid A'[i] \leq B_{rj} < A'[i+1]\}$

for $i=1, \dots, h$ {

if $B_i \neq \emptyset$

$\xrightarrow{\exists}$ MERGE-BASED INTERSECTION(A_i, B_i);

$\left. \begin{array}{l} \text{• Extra space: } O\left(\frac{n}{L}\right) \quad (A') \\ \text{• Time: } O\left(\frac{n}{L} + mL\right) \\ \text{• I/Os: } O\left(\frac{n}{LB} + \frac{mL}{B} + m\right) \end{array} \right\} \text{(*)}$

^(*) Fact Time is $O\left(\frac{n}{L} + mL\right)$ and #I/Os are $O\left(\frac{n}{LB} + \frac{mL}{B} + m\right)$ I/Os

Proof

$$T(n) = \underbrace{O\left(\frac{n}{L} + m\right)}_{\text{MERGE}(A', B)} + \sum_{B_i \neq \emptyset} O\left(\underbrace{|A_i| + |B_i|}_{\text{MERGE-BASED INT.}(A_i, B_i)}\right) = O\left(\frac{n}{L} + m\right) + O\left(\sum_{B_i \neq \emptyset} |A_i| + \sum_{B_i \neq \emptyset} |B_i|\right)$$

$$\stackrel{(*)}{=} O\left(\frac{n}{L} + m + mL + m\right)$$

at most SB blocks
touch m distinct A_i

$$= O\left(\frac{n}{L} + mL\right)$$

I/Os: $O\left(\underbrace{\frac{n}{LB} + \frac{m}{B}}_{\dots} + m\right)$

$\cup \left(\underbrace{A' \cup B}_{\text{MERGE}(A', B)} + m \uparrow \right)$
 at least 1 I/O (at most # of I/Os)
 for each $\text{MERGE BASED INT.}(A_i, B_i)$ (are no more than m)



General

lunedì 13 febbraio 2023 20:05

$$S = s_1, \dots, s_n$$

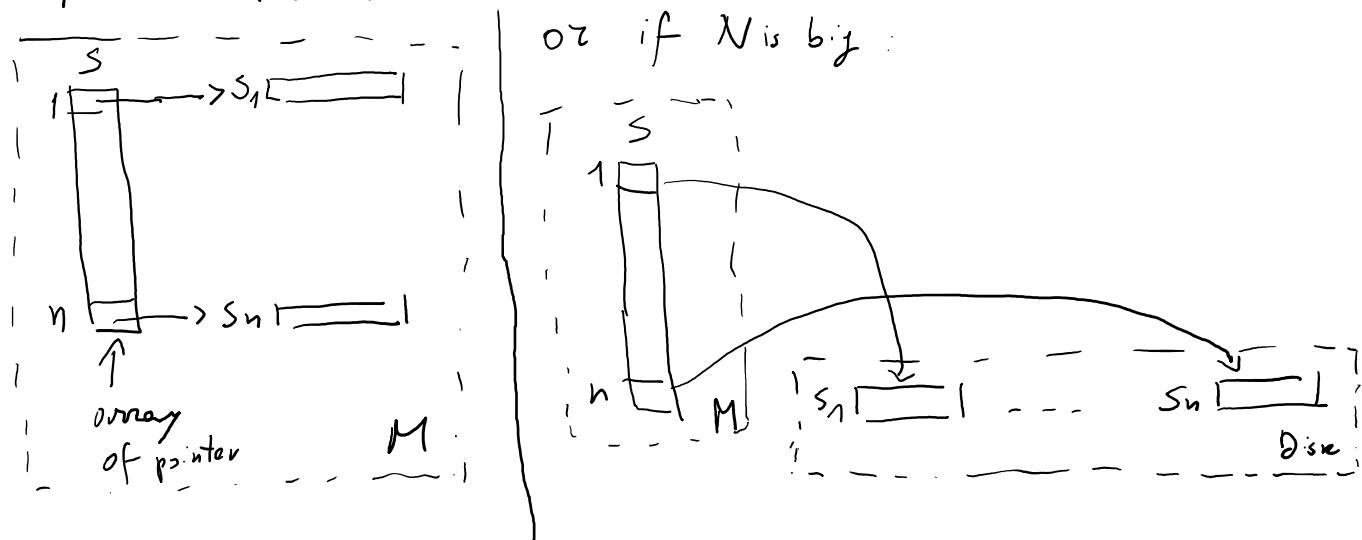
$n = \# \text{strings}$

$$N = \text{total length} = \sum_{i=1}^n |s_i|$$

$$L = \text{avg. length} = \frac{N}{n}$$

Σ alphabet, $\sigma = |\Sigma|$

Comparison based SORTER :



- $O(L \lg n) = O(N \lg n)$ time

$O(L)$
for each
string comparison

Comp. based sortor
for n items

- Every time comparing two string s_i and s_j , we must access to the string by the two pointer $\Rightarrow O(n \lg n)$ I/Os

Lower Bound

Lemma d_S = shortest prefix of S that distinguish it from others.

The lower bound of the string sorting problem is $\Omega(d + n \lg n)$ comparisons, where $d = \sum_{S \in S} d_S$ is the distinguishing prefix of S .

Potological case

ρ $\rho_{..}$

Petologial case

s_1, \dots, s_n stored in bit

$$|s_i| = l + \lg n, d_s = l + \lg n$$

$$J = \sum_{s \in S} d_s = n(l + \lg n)$$

$$s_1 = \overbrace{10\ldots0}^l \overbrace{10\ldots01}^{\lg n}$$

$$s_2 = \overbrace{10\ldots0}^l \overbrace{10\ldots10}^{\lg n}$$

⋮

$$s_n = \overbrace{1\ldots1}^l \overbrace{1\ldots1}^{\lg n}$$

$$\text{Lower bound} = \Omega(J + nl\lg n) = \Omega(n(l + \lg n) + nl\lg n) = \Omega(N + nl\lg n) = \Omega(N)$$

$$\text{Quicksort/Mergesort} = O(n(l + \lg n) \cdot \lg n) = O(N \lg n)$$

differs
 $\lg n$
by optimal

Radix sort (MSD-first)

martedì 7 marzo 2023 15:12

String or sequence of characters $s \in S$, $s = a_1 \dots a_{|s|}$, $a_i \in \{0, 1, \dots, \alpha-1\}$

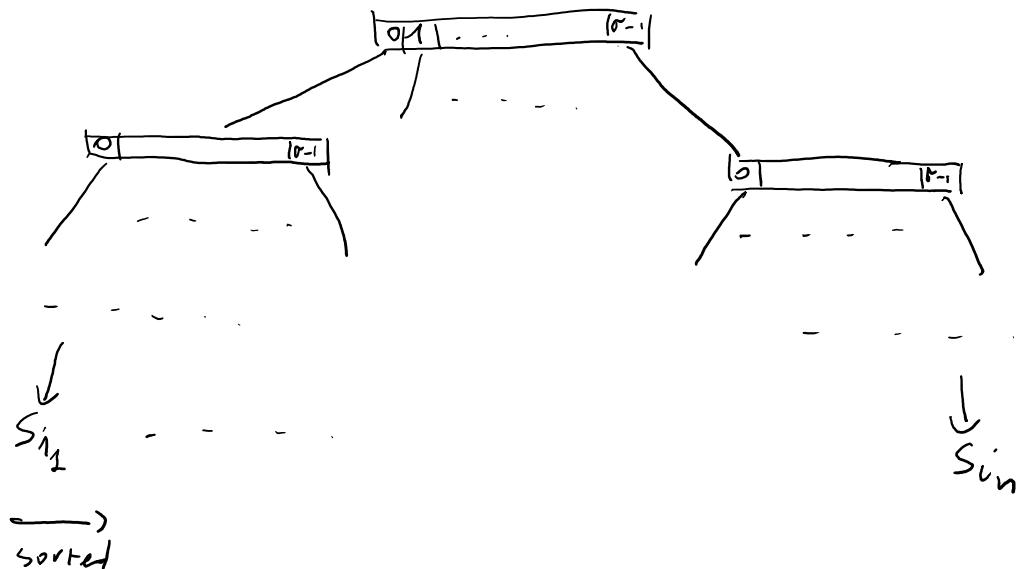
Naming process (assuming just given) $\mathcal{O}(N \lg \alpha)$

Space for S $\Theta(N \lg \alpha)$ bits (N = # chars, $\lg \alpha$ bits for each character)

MSD-first

Distribution-based approach \rightarrow α -ary tree (trie)

$$S^l = \{s_1, \dots, s_n\}$$



- Extra space: $\mathcal{O}(\# \text{internal node} \cdot \text{size(node)})$

$$\# \text{ internal node} = \# \text{branching node} + \# \text{leaves node} = \mathcal{O}(n) + \mathcal{O}(N) = \mathcal{O}(N)$$

Optimization: Drop unary nodes without branching node as descendent node
 $\Rightarrow \# \text{internal node} = \mathcal{O}(d)$

```

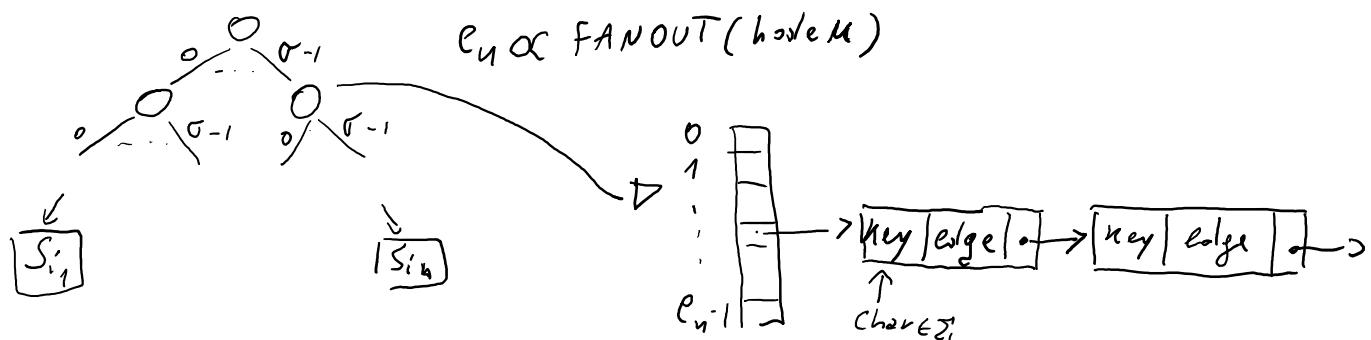
MSD-FIRST-ARRAY(S):
    tries ← φ; // empty tries
    for i = 1, ..., n
        tries ← INSERT( $s_i$ , tries);
        PREORDERVISIT(tries);
    
```

• Extra space: $O(d\alpha)$
 \uparrow
 $\# \text{nodes}$ size(node)
 } (*)
 • Time: $O(d\alpha)$

(*) Thm Extra space and time are $O(d\alpha)$ (using array)

Proof- To create a new internal node in the tree we require $O(\alpha)$ time and extra space (allocating the array). The number of internal nodes are $O(d)$ $\Rightarrow O(d\alpha)$ \square

Using hash table (with chaining) instead of array:



MSD-FIRST-HASHTABLE(S):
 // as MSD-FIRST-ARRAY(S)

• Extra space: $O(d)$
 • Time: $O(d\lg\alpha)$ in avg. } (*)
 } (avg)

(*) Thm Using hash with chaining time is $O(d\lg\alpha)$ in avg and extra space is $O(d)$.

Proof time in avg = time (construct trie) + time (sort edges) + time (pre-order visit)

$$\text{(1)} = \text{time(insert all strings)} = \# \text{nodes} \cdot \text{time(except the node)}$$

$$\begin{aligned} \mathcal{V} &= \text{time(insert all strings)} \approx \# \text{node} \cdot \text{time(exam the node)} \\ &\approx d \cdot O(1) \text{ in avg.} \end{aligned}$$

$$\text{time(exam the node)} \approx \text{avg. length of chain} = O\left(\frac{\# \text{keys}}{c_n}\right) = O\left(\frac{n\sigma}{n\sigma}\right) = O(1)$$

$$\textcircled{2} \quad \text{time}\left(\underset{\substack{\text{sort pointers of -} \\ \text{hash tables}}}{\sum_i c_n \cdot \text{key}}\right) = O\left(\sum_i c_n \cdot \log \sigma\right) = O\left(\sum_i c_n \log \sigma\right) = O(d \log \sigma)$$

$\sum_i c_n \propto \# \text{internal node}$

$$\begin{aligned} \textcircled{3} &= \text{time(scan trie rightward)} = \text{time(exam node)} \cdot \# \text{node} = O(1) \cdot d = O(d) \\ \Rightarrow \text{time} &= \textcircled{1} + \textcircled{2} + \textcircled{3} = O(d \log \sigma) \text{ in avg} \end{aligned}$$

$$\text{Extra space} = \text{space(list. node)} \cdot \# \text{internal node} = O(1) \cdot d = O(d)$$

OSS Not comparable with the LB on comparison based approach

Prop Using compacted trie $\# \text{internal node} = \# \text{branching node} = O(n)$

Radix sort (LSD-first)

martedì 7 marzo 2023 17:29

$L \sim \frac{N}{n}$ org. lenght $\forall i: S_i \in S, S_i := S[i:L]$ (otherwise padding)

$LSD\text{-first}(S) // S = S[1:n][1:L], S[i, :] = S_i$

for $i=L, \dots, 1 \{$

$A \leftarrow \text{CONCATENATE}(S_1[i], \dots, S_n[i]), // least significant digits (by now)$

$(index, A) \leftarrow \text{COUNTING SORT}(A);$

$S \leftarrow S[index, :];$

}

$\left. \begin{array}{l} \bullet \text{Spec: } O(N) \\ \bullet \text{Time: } O(N+L\sigma) \end{array} \right\} (*)$

\leftarrow cache misses

$\text{COUNTING SORT}(A)$

$\min \leftarrow \text{MIN}(A), \max \leftarrow \text{MAX}(A);$

$C \leftarrow \text{ZEROS}(1, \max - \min); // C = 0, \dots, 0$

for $i=1, \dots, n$

$C[A[i]-\min+1]++; // \# occurrences$

for $i=1, \dots, \max - \min$

$\text{OUTPUT}(\min + (i-1), C[i]); // C[i] times$

$\left. \begin{array}{l} \bullet \text{Time: } O(n + \underbrace{\max(A) - \min(A)}_{\min, \max, \#occ}) \\ \qquad \qquad \qquad \uparrow \\ \qquad \qquad \qquad \text{OUTPUT} \end{array} \right\}$

$(*)$
Lemma $\overline{Z}_1 = \{0, 1, \dots, \sigma-1\} \Rightarrow$ Correct, Spec is $O(N)$ and time $O(N+L\sigma)$

Proof $\text{Spec} = \text{Spec}(S) = O(N)$

$\text{Time} = L \times \text{time(counting sort)} = L \times (n + \sigma) = Ln + L\sigma = N + L\sigma$

Correctness (we prove the stability of Counting Sort): $\alpha, \beta \in S' \quad \alpha < \beta$ and

$\alpha = j \times \alpha_1, \beta = j \times \beta_1$ with $j = \text{lcp}(\alpha, \beta)$, $x < y$ and $|\alpha_1| = |\beta_1|$

So from right to left (LSD-first) when compare x and y (after $|\alpha_1| = |\beta_1|$ step) we have $\alpha < \beta$. For the next steps α and β hold the same positions \square

Optimization: $\forall S_i = B[1:b]$ binary \rightarrow 

Optimization: $\forall S_i = \text{Bit}[1..b]$ binary $S_i =$

$S_i =$
 ↑ r r r
 groups

Then Sort n strings of b bits each (using CS to sort r bits at time)
 with LSD-first \Rightarrow time $= \Theta\left(\frac{b}{r}(n+2^r)\right)$, space $= O(nb) = O(N \lg \sigma)$

Proof Each string S_i of $g = \frac{b}{r}$ "groups" of r bits each.

$$\begin{aligned} \Rightarrow \text{time} &= \# \text{ digits} * \text{time(CountingSort)} = \# \text{ groups} * \text{time(CS)} \\ &= O\left(\frac{b}{r}(n+2^r)\right) \end{aligned}$$

↑
max int with r bits

$$\text{Space} = O(bn) = O\left(\underbrace{N}_{\# \text{ digits}} \underbrace{\lg \sigma}_{\# \text{ bits per digit}}\right) \text{ bits}$$

④

Then Best choice of r is $r^* = \Theta(\lg n)$ bits, so that time $= O\left(\frac{bn}{\lg n}\right)$ and
 Space $= O(bn)$ bits. Extra space is $\Theta(n)$.

Proof Counting sort time $= O(n+2^r)$ so it is useless choose $r < \lg n$
 for $r \in [\lg n, n]$ $\frac{b}{r}(n+2^r) \uparrow \Rightarrow r^* = \lg n$. Time and space substituting r^*
 Extra space for Counting Sort $\Rightarrow \Theta(1(1)) = \Theta(2^{r^*}) = \Theta(n)$ ②

For time $= O\left(\frac{N \lg \sigma}{\lg n}\right)$, space $= O(N \lg \sigma)$ bits

Proof Recall $bn = \underbrace{N \lg \sigma}_{\# \text{ digits}} \underbrace{\lg \sigma}_{\# \text{ bits per digit}}$ ③

Cor MSD -first better of LSD -first iff $d = O\left(\frac{N}{\log n}\right)$

Proof	Time
MSD -first	$O(d \log r)$
LSD -first	$O\left(\frac{N \log r}{\log n}\right)$

□

Multi-key Quicksort

martedì 7 marzo 2023 20:23

Assume $R = \{s_1, \dots, s_n\}$ prefix free ($t_i \neq j \quad s_i \neq s_j \alpha$)
 (If not consider $R' = \{s_i'\} = \{s_i : \$\}$, with $\$ \notin \Sigma$.)

MULTIKEYQS(R, i): // start $i=1$

if $|R| \leq 1$
 return R ;

else {

$p \leftarrow \text{CHOOSEPIVOT}(R)$;

$(R_L, R_=, R_>) \leftarrow \text{THREE-WAY PARTITION}(R, p, i)$;

$A \leftarrow \text{MULTIKEYQS}(R_L, i)$;

$B \leftarrow \text{MULTIKEYQS}(R_=, i+1)$;

$C \leftarrow \text{MULTIKEYQS}(R_>, i)$.

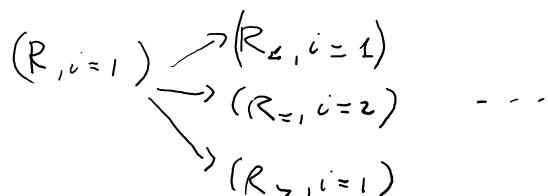
return $\text{CONCATENATE}(A, B, C)$;

}

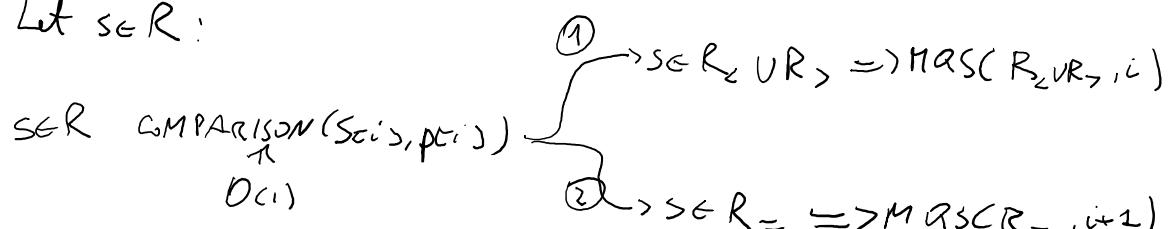
• Extra space: NO
 • Time: $O(d + n \lg n)$ on avg (*)

^(*) Thm It's correct, perform $O(d + n \lg n)$ comparison on avg. If good pivot selection strategy
 $\Rightarrow O(d + n \lg n)$ time in worst case.

Prf • Correctness: invariant of MULTIKEYQS(R, i): "all strings are sorted up to $(i-1)$ -th character".



• Time on avg.: Let $s \in R$:



if good pivot selection (balanced partition) $|R_L \cup R_>| \leq d/n, d < 1$:

① can occur at most $\frac{d}{n}$ times ($|R| = O(\lg n)$ times) } $\Rightarrow O(d_s + \lg n)$ # comparison for each s
 ② can occur at most $1/s$ times (is exploiting pref.) }
 $\Rightarrow \sum_{s \in R} O(d_s + \lg n) = O(d + n \lg n)$ comparisons.

$$\Rightarrow \sum_{s \in R} O(d_s + \lg n) = O(d + \lg n) \text{ comparisons.}$$

Good select pivot strategy: select the median pivot between 2 to 3 randomly selected pivots
(via oversampling). □

Multicay-AS or Ternary search Tree (TST)



// Start: SEARCH(S,T,i)
SEARCH(T) // T is ternary-search tree (pointer to root root)

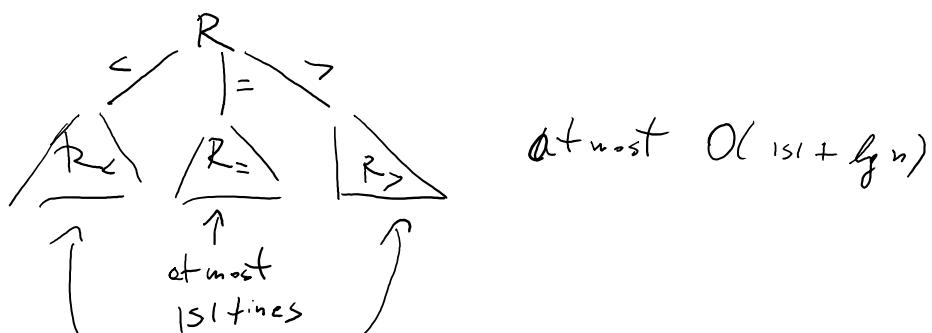
```

if (s == STRING(T))
    return true;
if (T == φ or i > |S|)
    return false;
if (S[i] == CHAR(T))
    SEARCH(S, MIDLCHILD(T), i+1);
else if (S[i] < CHAR(T))
    SEARCH(S, LEFTCHILD(T), i);
else
    SEARCH(S, RIGHTCHILD(T), i)
  
```

time: $O(|S| + \lg n)$ (*)

(*) Thm If ternary search tree is perfectly balanced time of search is $O(|S| + \lg n)$

Proof $|R_L \cup R_R| \leq \alpha n$, $\alpha < 1$ $n = |R|$



at most
151 times

at most
 $\log n$ times

□

General

lunedì 13 febbraio 2023 20:05

$$\mathcal{D} = \{\text{objects}\}, \mathcal{U} = \{\text{all keys}\} \quad n = |\mathcal{U}|, m = |\mathcal{D}|$$

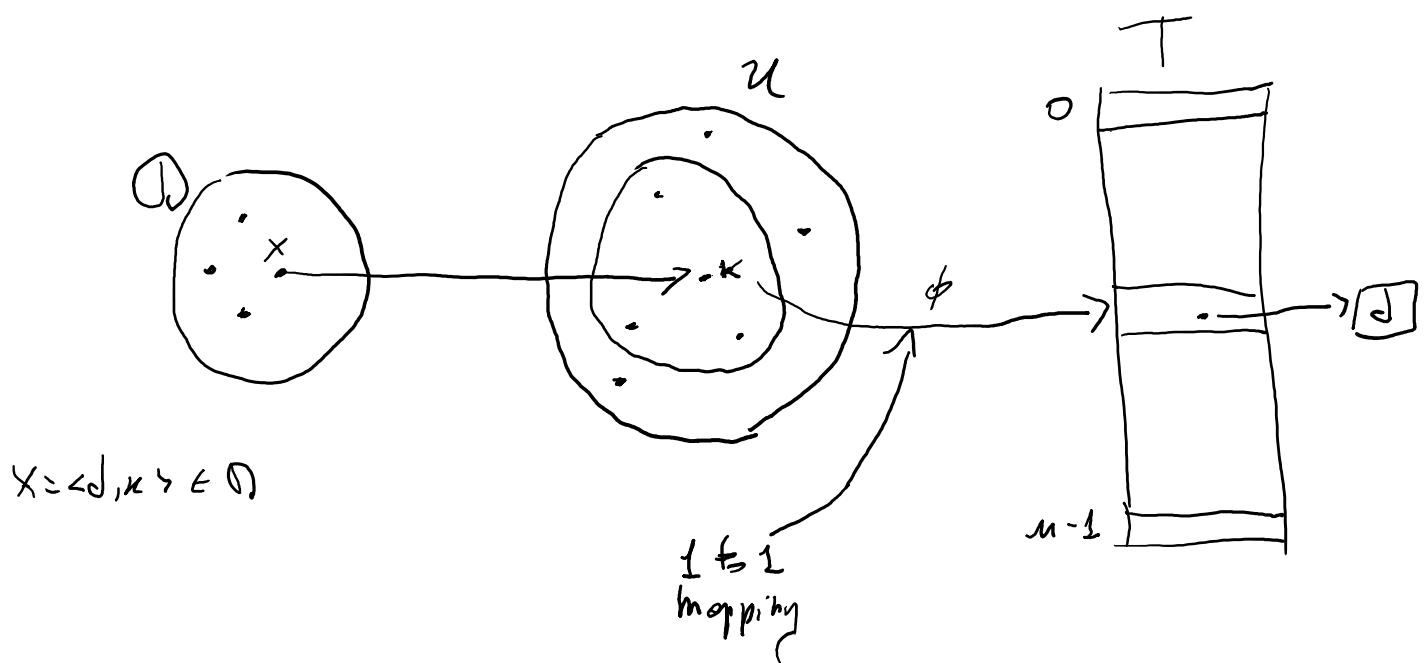
$$\begin{array}{ccc} \mathcal{D} & \xrightarrow{\quad} & S \subseteq \mathcal{U} \quad \text{bijective} \\ x & \longmapsto & \text{key}(x) = k \end{array}$$

object $x \in \mathcal{D}, x = \langle d, k \rangle$
 $d \uparrow$ satellite
 $k \uparrow$ date

ops: $\begin{cases} \text{Search}(k) \\ \text{Insert}(k) \\ \text{Delete}(k) \end{cases}$

Direct address table

venerdì 10 marzo 2023 17:15



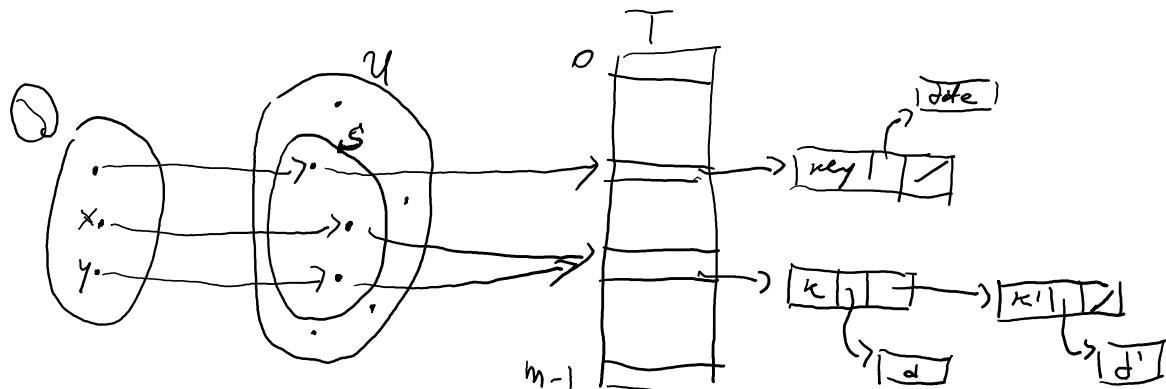
	Time	Space
Search(k)	$O(1)$	
Insert(k)	$O(1)$	$\Theta(n)$
Delete(k)	$O(1)$	

(+) Fast
(-) Space if $n \gg m$

Hash Table (hashing with chaining)

venerdì 10 marzo 2023 17:24

hash function $h: U \rightarrow \{0, 1, \dots, m-1\}$ $|U| = M \gg m$



Def Collision If $k' \neq k \in U$ and $h(k) = h(k')$

Def Simple uniform hashing $h(k) := k \pmod m = [k]_m$

	Time	Space
Search(k)	on avg	
Insert(k)	$\Theta(1 + \frac{n}{m})$ (*) "load factor"	$O((m+n) \lg n + n \lg n)$ bits (*)
Delete(k)		

(*) Using $h(k) = [k]_m$ } hash-table of size $= m$ such that time of Search(k) is $\Theta(1 + \frac{n}{m})$ on avg.

Proof 1st case) $k \notin S$:

$$\text{time on avg.} = \text{time}(\text{compute } h(k)) + E[\text{length}(\text{list } T[h(k)])]$$

$O(1)$

$$E[\text{length}(T[h(k)])] = \sum_{x \in U} (h(\text{key}(x)) = h(k)) = \sum_{x \in U} \frac{1}{m} = 1 \times \frac{1}{m} = \frac{1}{m} = \alpha$$

2nd case) $x \in S$ ($x = \langle d, k \rangle \in D$)

Let $\ell(x) = T[\text{hash}(x)]$, $E[\text{length } \ell(x)] = h_i \approx \frac{i-1}{m}$, if x is the i -th element inserted in T

$$\begin{aligned} \Rightarrow \text{time on } \text{array} &= \frac{1}{n} \sum_{i=1}^n \left(\underset{\substack{\# \text{ items} \\ \text{with } x \text{ inserted as } i\text{-th element}}}{\sum} \right) = \frac{1}{n} \sum_{i=1}^n \left(i + h_i \right) = \frac{1}{n} \sum_{i=1}^n \left(1 + \frac{i-1}{m} \right) \\ &= 1 + \frac{1}{n} \sum_{i=1}^{n-1} \underbrace{\frac{i}{m}}_{\frac{n(n-1)}{m}} = O\left(1 + \frac{n}{m}\right) \end{aligned}$$

□

(*) Thin hash-table with chaining occupies $(m+n) \lg_2 n + n \lg_2 u$ bits

Proof

$$\text{Space} = \text{space(pointers)} + \text{space(item keys)} \quad (\text{in bits})$$

$$\begin{aligned} &= (\underset{m}{\# \text{ pointers of } T} + \underset{O(n)}{\# \text{ pointers of chains}}) \cdot \underset{\lg_2 n}{\text{size(pointer)}} + (\underset{n}{\# \text{ item keys}}) \cdot \underset{\lg_2 u}{\text{size(key)}} \\ &\approx \frac{m}{2} \cdot \frac{n}{m} \cdot \lg_2 n \quad (\text{index one out of } n \text{ items}) + n \cdot \lg_2 u \end{aligned}$$

$$\approx (m+n) \lg_2 n + n \lg_2 u$$

□

Dynamical hash-table:

$$1. \text{ Start: } |D| = n_0 \Rightarrow m \leftarrow 2n_0$$

Rebuilding technique: $\begin{cases} \text{if } n < \frac{n_0}{2}: m \leftarrow \frac{m}{2}; \text{REBUILD}(T, h, m); n_0 \leftarrow n \\ \text{if } n > 2n_0: m \leftarrow 2m; \text{REBUILD}(T, h, m); n_0 \leftarrow n \end{cases}$

$$2. \text{ Load factor } \alpha = \frac{m}{n} = 2$$

Cor Dynamical hash-table with simple uniform hashing takes constant time, expected and amortized for each op.

Proof With maintaining $\alpha=2$ time for op = $O(1)$.

$$\text{time}(\text{REBUILD}(\dots)) = \underbrace{\text{time}(\text{n items to delete})}_{O(n)} + \underbrace{\text{time}(\text{n items to insert})}_{O(n)} = O(n)$$

at worst case (all insertions or all deletions):

$$\underbrace{op(k) + \dots + op(k)}_{\stackrel{\uparrow}{O(1)}} + \underbrace{\text{REBUILD}(\dots)}_{O(n)} \Rightarrow \text{spread } O(n) \text{ among } n \text{ operations having cost } = O(1)$$

$\Rightarrow O(1)$ amortized time per op.

Universal hashing

venerdì 10 marzo 2023 19:20

(uniform and independent dist. among keys)

Problem: to generate a good hash function, it should be one of all possible hash functions

$h: U \rightarrow \{0, \dots, m-1\} \Rightarrow m^u$ possible functions $\Rightarrow O(\lg_2 m^u) = O(u \lg_2 m)$ bits
Space to represent it (\rightarrow)

Solution: choose h random in a good set H

Def H is universal iff:

$$\forall x, y \in U \quad |\{h \in H \mid h(x) = h(y)\}| \leq \frac{|H|}{m} \quad \left| \begin{array}{l} \text{Def c-universal if } \\ \leq \frac{|H|}{m} \end{array} \right.$$

$$\text{Prop } h \in H \text{ universal}, \forall x, y \in U \quad p(h(x) = h(y)) \leq \frac{|H|}{m} = \frac{1}{m} \quad \left| \leq \frac{c}{m} \text{ if c-universal} \right.$$

Then hash table $T[0, m-1]$, with $h \in H$ universal. For a generic $k \in S$ $E[\text{length}(h, k)] \leq 1 + \alpha$, with $\alpha = \frac{c}{m}$

Proof $\forall_{k, k' \in S}$

$$I_{k, k'} := \begin{cases} 0 & \text{if } h(k) \neq h(k') \\ 1 & \text{otherwise} \end{cases} \quad \left(p(I_{k, k'} = 1) \leq \frac{1}{m} \right) \text{ (*)}$$

$$\text{Let } N_k = |\{k' \in S \mid k' \neq k, h(k') = h(k)\}|, \quad N_k = \sum_{\substack{k' \neq k \\ k' \in S}} I_{k, k'}$$

$$\begin{aligned} E[\text{length}(h, k)] &= E[N_k + 1] = E[N_k] + 1 = E\left[\sum_{k' \neq k} I_{k, k'}\right] + 1 = \\ &\stackrel{\substack{\text{we must} \\ \text{consider } k}}{=} \left(\sum_{k' \neq k} E[I_{k, k'}]\right) + 1 = \\ &= \left(\sum_{k' \neq k} p(I_{k, k'} = 1)\right) + 1 \stackrel{(*)}{\leq} \frac{m-1}{m} + 1 \leq 1 + \alpha \end{aligned}$$

Thus $T_{m, u}$ is $O(1 + \alpha \cdot m)$

Thm $T_{\lceil \alpha m \rceil}$ hashtable with $\log_2 H$ universal. If we insert $n = O(m)$ keys
 $\Rightarrow E[\text{length}[\text{longest chain}]] = O\left(\frac{\log n}{\log \log n}\right)$
 (NO PROOF)

THE POWER OF TWO CAVES

d sub-tables T_1, \dots, T_d of size m_j with independently hash functions

$$h_1, \dots, h_d : U \rightarrow \{0, \dots, m_j - 1\}$$

| $\text{INSERT}(x, T_1, \dots, T_d, h_1, \dots, h_d)$: | time: $O(d)$

for $i = 1, \dots, d$:

$$\alpha_i \leftarrow \text{LENGTH}(T_i[h_i(x)]),$$

$j \leftarrow \text{ARG-MIN}(\alpha_1, \dots, \alpha_d)$; (j leftmost in case of tie)

$\text{INSERT}(x, T_j, h_j);$

$\text{SEARCH}(x, T_1, \dots, T_d, h_1, \dots, h_d)$

$i \leftarrow 1$; stop \leftarrow false.

while (!stop AND $i < d$) {

if ($\text{SEARCH}(x, T_i, h_i)$)

stop = true;

else

$i++$;

return stop;

$$\text{time} = \begin{cases} O(\log \log n) & \text{if } d = 2 \\ O(d \left(\frac{\log \log n}{\log d} + 1 \right)) & \text{if } d > 2 \end{cases}$$

$\downarrow \cdot \text{time}(\text{SEARCH}(x, T_i)) \leq d \cdot E[\text{length chain}]$

$$= d \cdot \left(\frac{\log \log n}{\log d} + O(1) \right)$$

Examples of universal H

1) Assume $m = 101$ prime.

$$h \in U, \quad u = [u_0, \dots, u_{r-1}] \text{ by } r \text{ bits.} \quad r = \frac{\log m}{\log 2} \quad (|U| = m^r)$$

$$a \in \{1, \dots, m-1\}, \quad a = [a_0, \dots, a_{r-1}] \text{ (as } r \text{)}$$

$$\mathcal{H} = \left\{ h_e(k) = \sum_{i=0}^{r-1} a_i k_i \pmod{m} \mid e \in \{1, \dots, m-1\} \right\}$$

Then \mathcal{H} is an universal class

Proof Let $x \neq y \in U$. Without loss of generality we suppose $x_0 \neq y_0$

$$\begin{aligned} h_e(x) = h_e(y) &\Leftrightarrow \sum_{i=0}^{r-1} a_i x_i \equiv \sum_{i=0}^{r-1} a_i y_i \pmod{m} \\ &\Leftrightarrow a_0(x_0 - y_0) \equiv - \sum_{i=1}^{r-1} a_i(x_i - y_i) \pmod{m} \\ &\Leftrightarrow \xrightarrow[x \neq y \text{ and } m \text{ prime}]{} a_0 \equiv \left[- \sum_{i=1}^{r-1} a_i(x_i - y_i) \right] (x_0 - y_0) \pmod{m} \end{aligned}$$

\Rightarrow 1 choice for a_0 and a_1, \dots, a_{r-1} free

$$\Rightarrow |\{h_e(x) = h_e(y) \mid h_e \in \mathcal{H}\}| = m^{r-1} = \frac{|U|}{m}$$

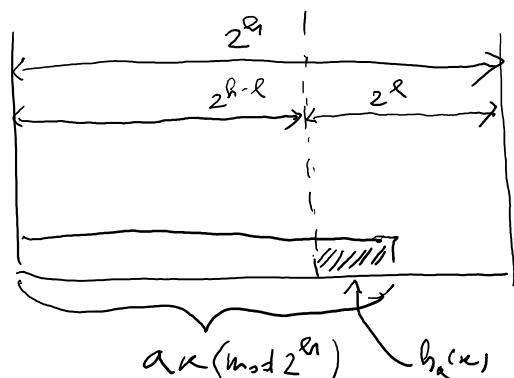
□

2) max $|U|$ size of $T(m, n, t \text{ prime})$. Set $p > |U|$ prime

$$\mathcal{H}_{p,m} = \left\{ h_{e,b}(k) = ((ak + b) \pmod{p}) \pmod{m} \mid \begin{array}{l} a > 0 \\ b > 0 \end{array} \right\}$$

3) Assume $|U| = 2^h$, $m = 2^l < |U|$, $e \equiv 1 \pmod{2}$ and $e < |U|$

$$\mathcal{H}_{e,l} = \left\{ h_{e,l}(k) = (ek \pmod{2^h}) \pmod{2^{h-l}} \right\}$$



Then $\mathcal{H}_{e,0}$ is 2-universal ($\forall x \neq y \in U, h_e \in \mathcal{H}_{e,0} \quad p(h_e(x) = h_e(y)) \leq \frac{1}{2}$)

Then $\mathcal{H}_{h,d}$ is 2-universal ($\forall x \neq y \in U, h_0 \in \mathcal{H}_{h,d} \quad p(h_0(x) = h_0(y)) \leq \frac{2}{m}$)

(No proof)

Perfect hashing, minimal, ordered

lunedì 13 marzo 2023 20:30

Def $h: \mathcal{U} \rightarrow \{0, \dots, m-1\}$ $\xrightarrow{\text{OPPHF}}$ perfect w.r.t. a dictionary $S \subseteq \mathcal{U}$ iff $\forall x \neq x' \in S$
 $h(x) \neq h(x')$

OSS h perfect $\Rightarrow |\mathcal{D}| = |S| = n \leq m$

Def If $m=n$, h perfect $\Leftrightarrow h$ minimal perfect

OSS h MPHF \Rightarrow

Search(x)	• Space: $O(n)$
Insert(x) (with hash table)	• Time: $O(1)$
Delete(x)	

Def h is order preserving w.r.t. S iff $\forall x < x' \in S \Rightarrow h(x) < h(x')$

OSS h order preserving $\Rightarrow h$ perfect

Def If $|S|=n=m$ and h OPHF $\Leftrightarrow h$ OPMPHF

Construction of e OPMPHF (Supposing fixed S)

- Pick randomly h_1, h_2 from \mathcal{H} universal,

$h_1, h_2: \mathcal{U} \longrightarrow \{0, \dots, m'-1\}$ with $m' = cn$, $c \geq 1$ (typically)

$$(m' \geq |S| = |\mathcal{D}|)$$

- Design a function $g: \{0, \dots, m'-1\} \longrightarrow \{0, \dots, n-1\}$ ($m' \geq n$)

s.t. $h(t) := g(h_1(t)) + g(h_2(t)) \pmod{n}$ is OPMPHF $\left(h: \mathcal{U} \rightarrow \{0, \dots, n-1\} \right)$
 $\text{OPMPHF w.r.t. } S$

Start: graph $G = (V, E)$, $V = \{0, \dots, m'-1\}$, $E = \{(h_1(t), h_2(t))\}_{t \in \{0, \dots, n-1\}}$ $\begin{cases} |V| = m' \\ |E| = n \end{cases}$

Acqis adjacency list of $i \in \{0, \dots, m'-1\}$

$\forall (u, v) \in E$ ($u, v \in \{0, \dots, m'-1\}$) has weight $h(u, v) = t$ where $u = h_1(t)$ and $v = h_2(t)$
 "an edge is a key"

LABELACYCLIC-GRAPH(G)

for $v \in V$

$g(v) = \text{undef}$ // initialize labels for nodes

• Extra space: $\Theta(n)$	}
• Time: $\Theta(n)$ on avg.	

(*)

```

for  $v \in V$ 
     $g[v] = \text{undef}$  // initializes labels for nodes
for  $v \in V\{$ 
    if  $g[v] = \text{undef}$ 
        LABEL_FROM( $v, 0$ );
    }
LABEL_FROM( $v, c$ )
if  $g[v] \neq \text{undef}\{$ 
    if  $g[v] \neq c$ 
        return; // the graph is cyclic
    }
 $g[v] = c$ ; // new label of node  $v$ 
for  $u \in A(v)$ 
    LABEL_FROM( $u, h(u, v) - g[v] \bmod n$ );

```

- extra space: $\Theta(n)$
 - Time: $\Theta(n)$ on avg. } (*)
- Note: If it returns a graph with unlabeled nodes (cyclic graph case) pick two new hash functions, rebuild G and restart the alg.

Prop (*) To construct the OMPHF we need to build $\sqrt{\frac{m'}{m'-2n}} = \Theta(1)$ graphs on avg.
So extra space and time are $\Theta(n)$ on avg.

Prof Given $|V| = m'$ and $|E| = n$

$$P(\text{build acyclic graph}) \approx \begin{cases} 0 & \text{if } m' \leq 2n \text{ (few nodes)} \\ \underbrace{\frac{\sqrt{m'-2n}}{m'}}_p & \text{if } m' > 2n \end{cases}$$

$$\Rightarrow E[\# \text{graph built}] = \sum_{l=1}^{\infty} l \cdot (1-p)^{l-1} p = \frac{1}{p} \quad (p(\# \text{graph built} \geq n) = (1-p)^{k-1} p)$$

geometric distribution

$$= \sqrt{\frac{m'}{m'-2n}} = \Theta(1)$$

$$m' = \Theta(n)$$

$$\Rightarrow \text{Time on avg.} = \Theta(1) \cdot \text{time(alg)} = \Theta(n)$$

$$\Rightarrow \text{Space on avg.} = \Theta(1) \cdot \text{space one graph} = \Theta(1) \cdot (\underbrace{m'}_{\text{nodes}} + \underbrace{n}_{\text{edges}}) = \Theta(n)$$

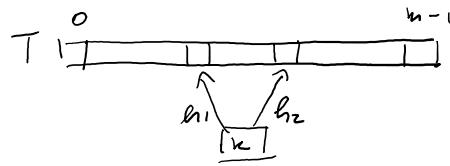
$$m' = \Theta(n)$$

Cuckoo hashing

martedì 14 marzo 2023 01:38

(Efficient for dynamic dictionary)

$$h_1, h_2 : U \rightarrow \{0, \dots, m-1\}$$



SEARCH(K, T, h_1, h_2):

```

 $i_1 \leftarrow h_1(k); i_2 \leftarrow h_2(k);$ 
if (!SEARCH( $K, T[i_1]$ ))
    return SEARCH( $K, T[i_2]$ )
else
    return true;
    
```

DELETE(K, T, h_1, h_2):

```

 $i_1 \leftarrow h_1(k); i_2 \leftarrow h_2(k);$ 
if (SEARCH( $K, T[i_1]$ ))
    return DELETE( $K, T[i_1]$ );
else if (SEARCH( $K, T[i_2]$ ))
    return DELETE( $K, T[i_2]$ );
else
    return  $\overline{T}$ ;
    
```

INSERT(K, T, h_1, h_2, max)

```

if ( $T[h_1(k)] = \emptyset$ ) {
     $T[h_1(k)] = K;$ 
    return  $\overline{T}$ ;
}
if ( $T[h_2(k)] = \emptyset$ ) {
     $T[h_2(k)] = K;$ 
    return  $\overline{T}$ ;
}
 $i \leftarrow T[h_1(k)];$  // we could choose  $T[h_2(k)]$ 
steps  $\leftarrow 1;$ 
while ( $T[i] \neq \emptyset$  AND steps < max)
    { $T, k, i \leftarrow KICKOUT(T[i], k);$  steps++}
if (steps  $\geq max$ )
    return false; // we must replace two new  $h_1, h_2$ 
else
    { $T[i] \leftarrow K; return \overline{T};$ }
    
```

KICKOUT($T[i], k$):

```

 $k' \leftarrow T[i];$ 
 $j \leftarrow h_x(k');$  // alternative pos. of  $k'$  ( $i \rightarrow j \in G$ )
 $T[i] \leftarrow k';$ 
return  $T, k', j$ ;
    
```

	Time	Space
Search(n)	$O(1)$	$(*)$
Delete(n)	$O(1)$	$\Omega(\frac{2n}{\log n})$ bits
Insert(n)	$O(1)$ amortized $O(1)$ expected	$\Theta(n)$ keys

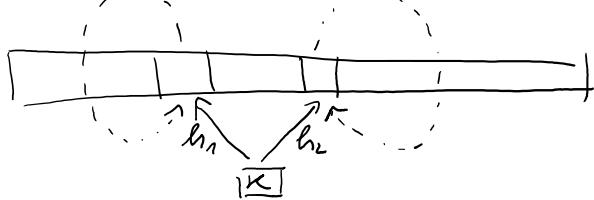
$\text{TrigLc};$
 $\text{return } T, x^!, j;$

Def Cacheo graph $G = (V, E)$ direct graph

$$V = \{i \mid T[i] \neq \phi\}, E = \{i \rightarrow j \mid i = h_x(k), j = h_y(k) \text{ and } T[i] = k\}$$

$(x \neq y \in \{1, 2, 3\})$

OSS → infinite loop iff



Def $\kappa' \neq \kappa$, $\kappa' \in \text{bucket}(\kappa)$ iff \exists undirected path from $b_1(\kappa)$ or $b_2(\kappa)$ to $b_1(\kappa')$ or $b_2(\kappa')$ ($\exists b_1(\kappa)/b_2(\kappa) \xrightarrow{\text{u.path}} b_1(\kappa')/b_2(\kappa')$)

Thm $\forall i, j, \forall c > 1 \text{ if } m \geq 2cn \Rightarrow p(i \stackrel{\text{L}}{\sim} j, \text{shortest}) \leq \frac{1}{c^L m}$

Proof By induction:

$$L=1 \quad p(i \rightarrow j) = p(i \rightarrow j \vee j \rightarrow i) \leq \frac{2}{m^2} \quad \left(p(i \rightarrow j) = \frac{1}{m^2} \right)$$

Unlikely bound

$$\Rightarrow \sum_{n \in S} \frac{2}{m^2} = \frac{2n}{m^2} \leq \frac{1}{cm}$$

$n \leq \frac{m^2}{2c}$

$$\begin{aligned}
 L-1 \Rightarrow L) \quad p(i \overset{L}{\min} j) &= P(\exists z \mid i \overset{L-1}{\min} z \overset{1}{\min} j) \stackrel{\text{def}}{=} \Pr_{\text{choice of } z} p(i \overset{L-1}{\min} z \overset{1}{\min} j) = \\
 &= m \cdot p(i \overset{L-1}{\min} z) \cdot p(z \overset{1}{\min} j) \stackrel{\text{inductive step}}{=} m \cdot \frac{1}{c^{L-1}_m} \cdot \frac{1}{c^1_m} = \frac{1}{c^L_m}
 \end{aligned}$$

Thm $p(x' \in \text{bucket}(x)) = O(\frac{1}{t_m})$, given $x, x' \in S$ (with $m, z \ll n$)

$$\Pr_{\text{rand}}[p(x \in \text{bucket}(w)) = p(\exists L \mid h_i(w) \leq \max_{\substack{i,j \in \mathcal{E}, l \in L \\ i,j \neq x}} h_j(w))] \leq q \cdot \sum_{L=1}^{\infty} \frac{1}{C_L m} = \frac{q}{m} \cdot \frac{1}{C-1} = O\left(\frac{1}{m}\right)$$

A pos. b/c
 choice of
 path

Thm $p(\text{unsuccessful insertion}) \leq \frac{1}{c-1}$ (with m, z, c, n)

Proof

$$\begin{aligned}
 p(\text{infinite loop}) &= p(\exists i, j \text{ such that } i \text{ and } j \text{ are in the same bucket}) \leq p(\exists i \text{ such that } i \text{ is in a bucket}) \\
 &= p(\exists i, \exists L : i \in L) = \sum_{i=0}^{m-1} \sum_{L \ni i} \frac{1}{m c^L} = m \left(\frac{1}{m} \cdot \sum_{L \geq 1} \frac{1}{c^L} \right) \\
 &\rightarrow \frac{1}{c-1}
 \end{aligned}$$

Corollary 1 Setting $c=3$ $\underbrace{p(\text{infinite loop})}_{1-p} \leq \frac{1}{2}$ (with $m \geq 2c\epsilon$)

Corollary 2 $E[\#\text{rehashes}] = \sum_{i=0}^{m-1} i \cdot (1-p)^{i-1} p = \frac{1}{p} \leq 2$
 $p \geq \frac{1}{2}$ (prob of successful insertion)

Corollary 3 Suppose to perform ϵn insertion in a table hash of size m with initially n keys and suppose $m \geq 2c(n+\epsilon n)$. The cost of the ϵn insertion has constant expected amortized time.

Proof Using Cor 1 and Cor 2 (with $m \geq 2c(n+\epsilon n)$) we have constant number of re-hashes expected

$$\begin{aligned}
 E[\text{time}] &= E[\#\text{rehashes}] \cdot \text{time(rehash)} \\
 &= O(1) \cdot O(n) = O(n) \quad (\text{for all } \epsilon n \text{ insertions}) \\
 &\quad \text{compute 2 hashes per key} \\
 \Rightarrow O\left(\frac{n}{\epsilon n}\right) &= O\left(\frac{1}{\epsilon}\right) = O(1) \text{ expected amortized time for one insertion}
 \end{aligned}$$

DS For general cost (generic n and ϵ) use rebuilding technique
 (half or double the table)

Bloom filters

martedì 14 marzo 2023 16:59

When storage of the keys is expensive (if $n \log m$ bits is too large)

- (+) Keys not explicitly stored
- (-) One-side error (only FP errors)

$h_i: U \rightarrow \{0, 1, \dots, m-1\} \quad i=1, \dots, r \quad$ universal hash functions

// Start: $B = \text{ZEROS}(1, m)$

$\text{INSERT}(k, B, h_1, \dots, h_r)$

for $i=1, \dots, r$

$B[h_i(k)] \leftarrow 1$

	Time	Space
Insert(k)	$O(r)$	m bits
Search(k)	$\tilde{\Theta}\left(\frac{m}{n} \log n\right)$	
Delete(k)		

$\text{SEARCH}(H(k, B, h_1, \dots, h_r))$ // possible FP error

$\text{res} \leftarrow \text{true}; i \leftarrow 1;$

while (res AND $i < r$) {

 if ($B[h_i(k)] = 0$)

$\text{res} \leftarrow \text{false};$

}

return $\text{res};$

Lemme $p(B[j] = 1) = 1 - e^{-\frac{n}{m}}$

Prof $p(B[j] = 0) = \left(\frac{m-1}{m}\right)^{n/r} = \left(1 - \frac{1}{m}\right)^{n/r} = \left[\left(1 - \frac{1}{m}\right)^m\right]^{\frac{n}{m}} \underset{n \rightarrow \infty}{\approx} [e^{-1}]^{\frac{n}{m}} = e^{-\frac{n}{m}}$

$\frac{n \text{ keys}}{r \text{ hashes}}$

$\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n = e^x$

$\Rightarrow p(B[j] = 1) = 1 - p(B[j] = 0) = 1 - e^{-\frac{n}{m}}$

□

Prop $p(\text{FP error}) = (1 - e^{-\frac{n}{m}})^r \approx p_{err}$

$$\text{If } p(\text{FP error}) = (1 - e^{-\frac{m}{n}}) \approx p_{\text{err}}$$

$$\begin{aligned} \text{Pf } p(\text{FP error}) &= P(\kappa \notin S, B[h_1(\kappa)] = \dots = B[h_r(\kappa)] = 1) \\ &= P(B[j_1] = \dots = B[j_r] = 1) \\ &= (P(B[j] = 1))^r = (1 - e^{-\frac{m}{n}})^r \end{aligned}$$

□

Corollary Given n and m fixed the best number of hashes is $r^* = \frac{m}{n} \ln 2$ and the respective p_{err} optimal is $(0.6185)^{\frac{m}{n}}$.

$$\begin{aligned} \text{Pf } f(r) &:= (1 - e^{-\frac{m}{n}})^r \quad f'(r) = 0 \iff r = \frac{m}{n} \ln 2 \\ \Rightarrow p_{\text{err}} &= (1 - e^{-\frac{m}{n}r^*})^{r^*} = (1 - e^{-\frac{m}{n} \ln 2})^{\frac{m}{n} \ln 2} = \left[\left(\frac{1}{2}\right)^{\ln 2}\right]^{\frac{m}{n}} = (0.6185)^{\frac{m}{n}} \end{aligned}$$

Prop To ensure $p_{\text{err}} \leq \epsilon$ with ϵ given we must have B of size $m \geq 1,44 \cdot \lg_2\left(\frac{1}{\epsilon}\right) n$ bits

$$\begin{aligned} \text{Pf We want } p_{\text{err}} &\leq \epsilon \iff \left[\left(\frac{1}{2}\right)^{\ln 2}\right]^{\frac{m}{n}} \leq \epsilon \\ &\stackrel{u_{\text{size}}}{\uparrow} \\ &r^* \geq \frac{1}{\epsilon} \end{aligned}$$

$$\ln 2 \cdot \frac{m}{n} \geq \lg_2\left(\frac{1}{\epsilon}\right)$$

$$m \geq \frac{n}{\ln 2} \lg_2\left(\frac{1}{\epsilon}\right) \approx 1,44 \cdot \lg_2\left(\frac{1}{\epsilon}\right) \cdot n$$

□

Spectral Bloom filters

martedì 14 marzo 2023 18:38

Extension of bloom filters to multisets

Def $S = \{(x, f_x)\}$ multiset $x \in U, f_x = \text{multiplicity of } x \text{ in } S$

h is an $h_1, \dots, h_r : U \rightarrow \{0, \dots, m-1\}$

$C[1, m]$ counter vector $C[i] = \sum_{\substack{x \text{ st. } h_j(x)=i \\ j=1, \dots, r}} f_x$

Ops: $\begin{cases} \text{Counting}(x) \\ \text{Insert}(n)/\text{Increment}(x) \\ \text{Delete}(n)/\text{Decrement}(n) \end{cases}$

// Start: $C = \text{zeros}(1, m)$

$\text{INCREMENT}(x, f_x, C, h_1, \dots, h_r)$

for $i = 1, \dots, r$

$$C[h_i(x)] += f_x;$$

$\text{DECREMENT}(x, f_x, C, h_1, \dots, h_r)$

for $i = 1, \dots, r$

$$C[h_i(x)] -= f_x;$$

$\text{Counting}(x, C, h_1, \dots, h_r)$ // overestimates f_x in S

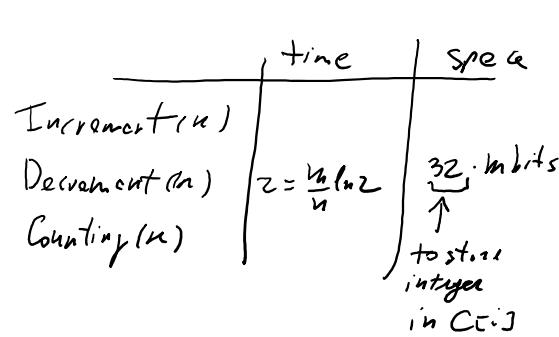
min = 0;

for $i = 1, \dots, r$ // $C[h_i(x)]$ are overestimations

if ($C[h_i(x)] < \text{min}$)

$$\text{min} \leftarrow C[h_i(x)];$$

return min;



prop $p(\text{error counting}) = (0.6185)^m$

$$\begin{aligned}
 p_{\text{coll}} &= p(\text{larger estimation}) = p(\tilde{c}_{\text{high}}) > f_k \quad k=1, \dots, r \\
 &= p(\text{Collision } k, k=1, \dots, r) = (1 - e^{-\frac{n}{m}})^r = \left(\frac{1}{2}\right)^r = \left[\left(\frac{1}{2}\right)^{\rho_m}\right]^{\frac{n}{m}} = (0.6185)^{\frac{n}{m}}
 \end{aligned}$$

↑
see BT
using RT

General

lunedì 13 febbraio 2023 20:06

$$\mathcal{D} = \{s_1, \dots, s_n\} \quad \text{dictionary of strings, } \Sigma \text{ alphabet}$$
$$n = |\mathcal{D}|$$
$$N = \sum_{s_i \in \mathcal{D}} \text{len}(s_i)$$
$$|\Sigma| = \alpha$$

We use 2 special characters \$\\$ and \$\#

such that

\$ \forall x \forall e \in \Sigma,

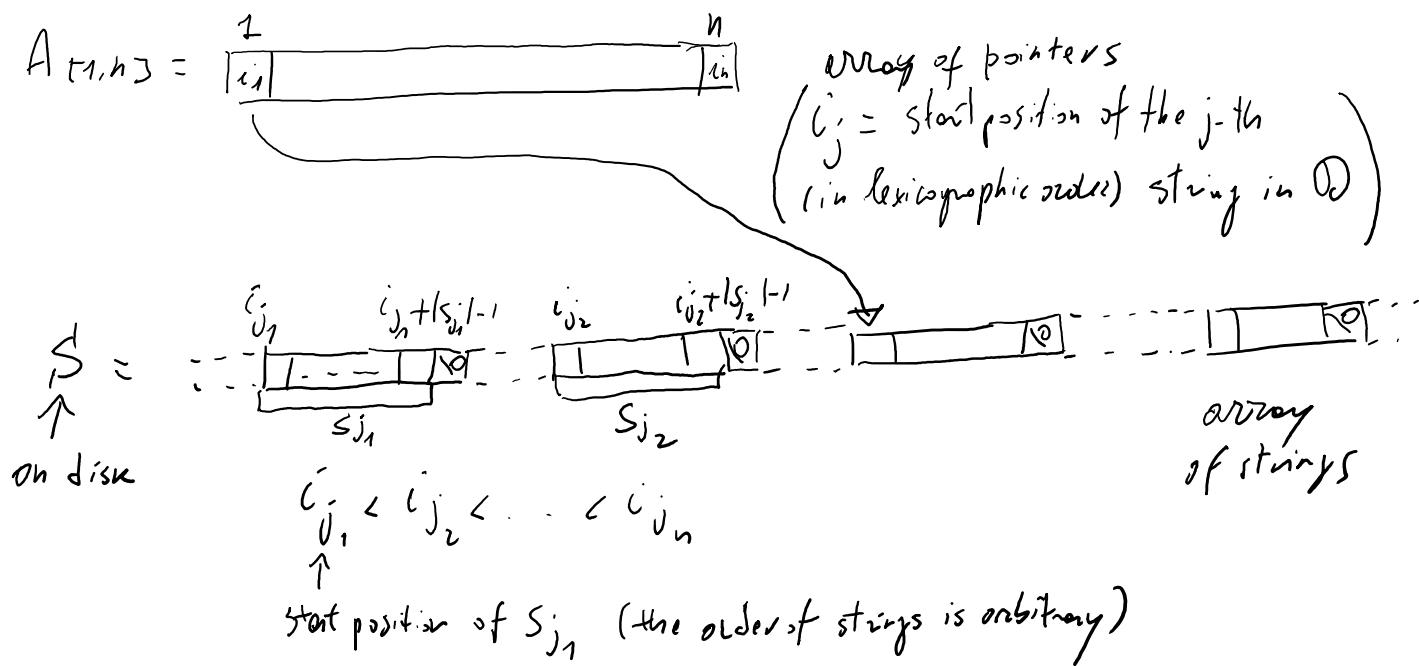
\$\forall a \forall e \in \Sigma\$

2 Operations given string $P \in \mathcal{D}$

- $\text{Count}(P) = \left| \{s_i \in \mathcal{D} \mid s_i = [P, q]\} \right|$
- $\text{Relative}(P) = \{s_i \in \mathcal{D} \mid s_i = [P, q]\}$

Array of string pointers

martedì 14 marzo 2023 19:39



RETRIEVE(P, A, S):

```

l ← BINARYSEARCH(P$, A, S);
r ← BINARYSEARCH(P#, A, S);
return A[l:r];

```

- Space: $N \lg r + (\lg r + w)$ in bits
32/64 bits
 - Time: $O(p \lg n)$
 - I/Os: $O\left(\frac{p}{B} \lg n\right)$

(a)

(4)

Thm time is $O(p \log n)$, I/Os are $O(p_B \log n)$ and the total space is $N + (1+w)n$ bits. Retrieve n_{acc} prefixed by P needs $\Omega(n_{\text{acc}})$ I/Os.

Prof time: Binary search on A: $O(\log n)$

↑
lexicographic
comparison (P_{C_1, P_2})

Some reasons for $O(gly_2)I/O_3$

Space: $N_{\text{strings}} + N_{\text{delimiters}} + w_n$ bits
 strings ↑ delimiters ↑ pointers
 $\underbrace{\hspace{1cm}}$ $\underbrace{\hspace{1cm}}$ $\underbrace{\hspace{1cm}}$

$$= N + (1+w)n \text{ bits}$$

After defining the range $A[l:r]$ each string (prefixed by P) can be visualized using at least 1 I/O because the contiguity in A doesn't imply the contiguity in S.

using at least 1 I/O because the contiguity in A doesn't imply the contiguity in S .

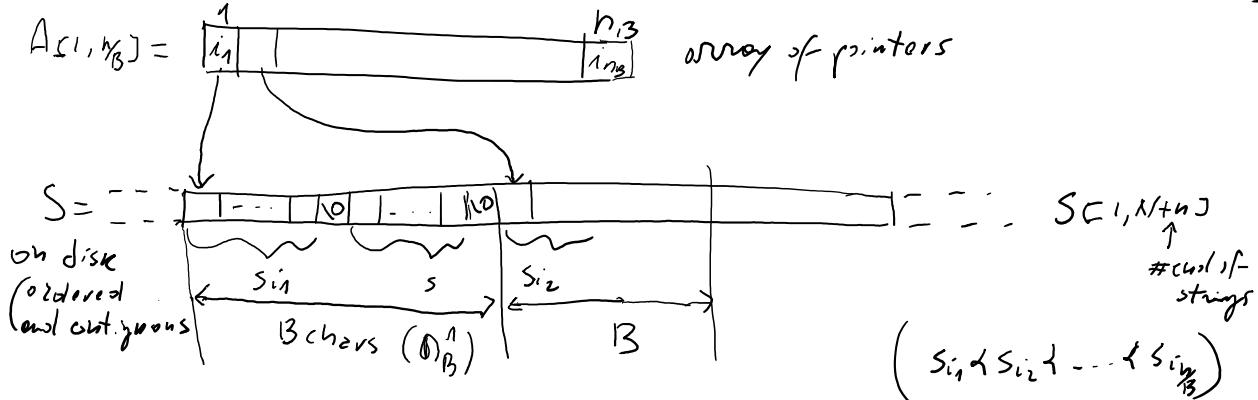


Contiguous allocation of strings

martedì 14 marzo 2023 20:53

$$\text{① } \text{memb} \text{ } \text{①}_{\mathcal{B}} = \{s_{i_1}, \dots, s_{i_{n_{\mathcal{B}}}}\} \subseteq \{s_1, \dots, s_n\} \quad |\text{①}_{\mathcal{B}}| = n_{\mathcal{B}} \leq \frac{N}{B}$$

$s_{ij} \in \text{①}_{\mathcal{B}}$ if i is the first string
of the j -th block



RETRIEVE(p, A, S) // $p \in P$
 $l \leftarrow \text{SEARCH}(p_L, A, S);$
 $r \leftarrow \text{SEARCH}(p_R, A, S);$
 return $\langle l, r \rangle.$

- Space: $(Nn) \lg \omega + w n_B$ bits
- Time: $\mathcal{O}(\varphi \lg \frac{N}{B})$
- I/Os: $\mathcal{O}(\frac{\varphi}{B} \lg \frac{N}{B})$

(4)

SEARCH(Q,A,S); // Q(A,P)

```

j ← BINARY SEARCH(QA); // give the j-th block
pos ← SCAN & COMPARE(Q, SCA[j], AC[j+1]) // Scan & compare
return pos;

```

(A) Thm Time is $O(P \log \frac{N}{B})$, I/Os are $O(\frac{P}{B} \log \frac{N}{B})$ and space is $(N_{\text{rec}}) \lg r + w b$ bits.
 Retrieve the strings prefixed by P needs $O(\frac{N_{\text{rec}}}{B})$ I/Os ($N_{\text{rec}} = \frac{\text{total length of strings prefixed by } P}{B}$)

Proof time: $O(p \lg n_B + B)$ = $O(p \lg n_B) \leq O(p \lg \frac{N}{B})$

↓
steps for BS in A

↓
at most
1 step for
each block

strings comparison

For the same reason $O\left(\frac{f}{B} \cdot g \frac{N}{B}\right) I/Os$.

$\text{Spec} \alpha : (N_{\text{fun}}) \text{ by } \alpha + w \text{ in } n \text{ bits}$

$S \leftarrow$ $(N+n)$ by $\alpha + w \lceil \frac{N}{B} \rceil$ bits
 # bits to store a pointer \rightarrow # pointers in A

After we retrieve $S[l, z]$ we scan the contiguous part to visualize
 all the strings with prefix $P \Rightarrow \mathcal{O}\left(\frac{N_{occ}}{B}\right)$ ($N_{occ} = z - l$)

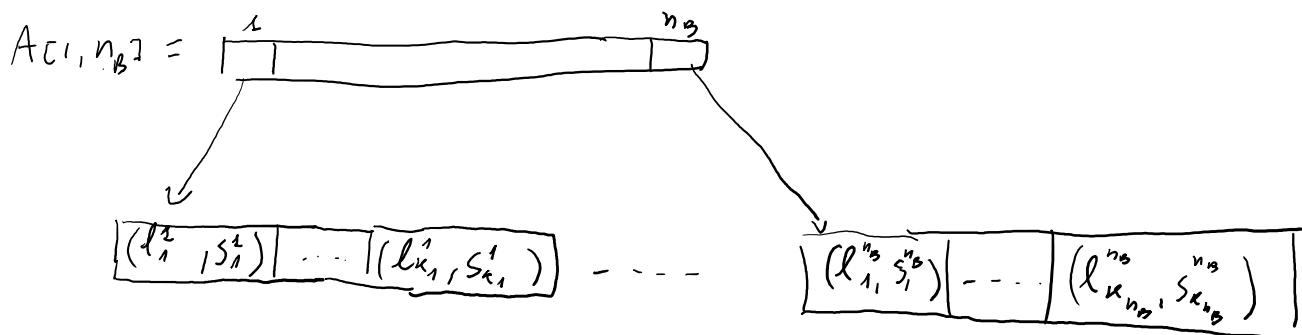
Front coding with bucketing

mercoledì 15 marzo 2023 00:41

$$\mathcal{D} = \{s_1, \dots, s_n\} \xrightarrow{\text{A-compression}} \{(l_1, \hat{s}_1), \dots, (l_n, \hat{s}_n)\}$$

$$\begin{cases} l_1 = 0 \\ l_i = \text{LCP}(s_i, s_{i-1}) \end{cases}, \quad \begin{cases} \hat{s}_1 = s_1 \\ \hat{s}_i = s_i[l_i+1, |s_i|] \end{cases} \Rightarrow s_i = [s_{i-1}[1, l_i], \hat{s}_i]$$

Two-level indexing: Front coding per bucket (FC_B)



Def $FC_B(\mathcal{D})$: space to store all the strings compressed by FC_B

OSS Now $n_B \approx \frac{FC_B(\mathcal{D})}{B} < \frac{N}{B}$ (+) (potentially increase the # strings)
stuffed in a disk page

$\text{RETRIEVE}(P, A, S) // PC1, P)$

// like contiguous allocation of strings

• Space: $\underbrace{FC_B(\mathcal{D})}_{\text{strings}} + \underbrace{(1+w)n_B}_{\text{pointers}}$
 • Time: $O(p \lg_2 \frac{FC_B(\mathcal{D})}{B} + B)$
 • I/Os: $O(\frac{p}{B} \lg_2 \frac{FC_B(\mathcal{D})}{B})$

(*) Then I/Os of prefix search are $O(\frac{p}{B} \lg_2 \frac{FC_B(\mathcal{D})}{B})$. Retrieving the strings prefixed by P needs $O(\frac{FC_B(\mathcal{D}_{\text{occ}})}{B})$ I/Os where \mathcal{D}_{occ} are the strings wanted

Proof time $O(p \lg_2 n_B + \underbrace{\text{size(decompress}(\mathcal{D}))}_{\text{string compression}} + \underbrace{\# \text{steps of BS in } A}_{\text{scan}}) \approx O(p \lg_2 \frac{FC_B(\mathcal{D})}{B} + B)$
 OSS usually $O(B^2)$ very rare

For I/Os: $O(p \lg_2 n_B)$

(15-)
Very rare)

For I/O_s : $O\left(\frac{P}{B} \lg \frac{FC_B(O)}{B}\right)$.

To retrieve all the strings prefixed by P ($\text{tse} \otimes \text{arc}$) we need

$O\left(\frac{FC_B(Doc)}{B}\right) I/Os$ (space to store compressed strings
size of disk page)

Interpolation search

mercoledì 15 marzo 2023 02:01

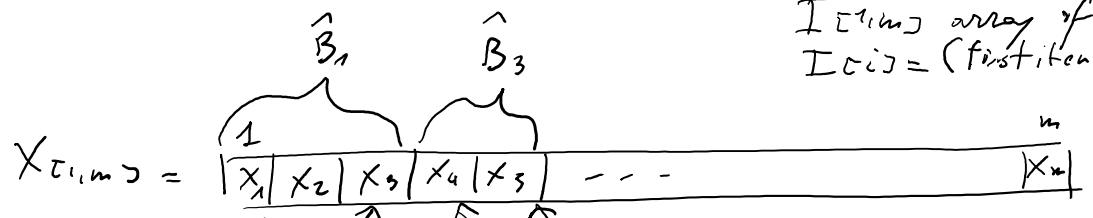
\mathcal{O}_B or integers (if not see string as number in base $\sigma = |\Sigma|$)

$\mathcal{O}_B \rightarrow X_{[1,m]} = x_1, \dots, x_n \xrightarrow{\text{distinct}} \text{integers} \quad x_i < x_{i+1} \quad (n = n_B)$

Transform range $[x_1, x_m]$ in bins B_1, \dots, B_m of length $b = \frac{x_m - x_1 + 1}{m}$

$$B_i := [x_1 + (i-1)b, x_1 + ib], \quad \hat{B}_i = X \cap B_i$$

$I_{[1,m]}$ array of pointers
 $I_{[c,d]} = (\text{first item of } \hat{B}_c, \text{last item of } \hat{B}_d)$



$$I_{[1,m]} = \underbrace{[c_{1,3} | (0,0) | (4,5)]}_{\dots}$$

SEARCH(y, X, I): // search integer y in X

$$j \leftarrow \lfloor \frac{y - x_1}{b} \rfloor + 1; \quad // \text{compute the } B_j \text{ s.t. } y \in B_j$$

$l, r \leftarrow I_{[j]}; \quad // \text{determine the pos. of first and last item}$

return BINARY SEARCH($y, X[l:r]$);

• Extra Space: $O(m)$
 • Time: $O(\lg \Delta)$

(*) Then searching a string in a dict of size m , take $\mathcal{O}(\lg \Delta)$ time,

where $\Delta = \frac{\max_{i=2, \dots, m} (x_i - x_{i-1})}{\min_{i=2, \dots, m} (x_i - x_{i-1})}$. The extra space is $O(m)$

Time: $\mathcal{O}(1 + \underbrace{\lg_2 |\hat{B}_j|}_{\substack{\text{select} \\ B_j}}) = \mathcal{O}(\lg_2 |\hat{B}_j|) \quad (< \mathcal{O}(\lg_2 b))$
 A generic \hat{B}_j contains at most $\frac{b}{s}$ integers where $s = \min_{i=2, \dots, m} (x_i - x_{i-1})$.

And

$$\max_{i=2, \dots, m} (x_i - x_{i-1}) \geq \frac{\sum_{i=2}^m (x_i - x_{i-1})}{m-1} = \frac{x_m - x_1}{m-1} \geq \frac{x_m - x_1 + 1}{m} = b$$

$$\Rightarrow |\vec{B}_j| \leq \frac{b}{S} \leq \Delta \Rightarrow \text{time} = O(\lg_2 |\vec{B}_j|) = O(\lg_2 \Delta). \quad \frac{a}{b} \geq \frac{a+1}{b+1} \quad (\text{if } a > b)$$

$$S_{\text{space}} = \underbrace{S_{\text{space}}(x)}_{O(m)} + \underbrace{S_{\text{space}}(\vec{x})}_{O(n)} = O(m+n)$$

oss $\Delta = \begin{cases} 1 & \text{if uniform dist.} \\ \geq 1 & \text{otherwise} \end{cases}$ and time is $O(\lg(\min\{\Delta, m\}))$
 $\Rightarrow O(\lg(m)) \leq O(\lg \Delta) \leq O(\lg m)$
 \uparrow
time BS on X

^(*) Thm If X is uniformly distributed $\Delta = \text{poly}(\lg(m))$ with high probability,
 and time became $O(\lg(\lg(m)))$ with high probability.

Proof (Hint) $\Delta = \text{poly}(\lg(m)) = P(\lg(m)) = a_k \lg^k(m) + \dots + a_1 \lg(m) + a_0$
 $\Rightarrow O(\lg_2 \Delta) = O(\lg_2(a_k \lg^k(m) + \dots)) = O(\lg_2(\lg^k(m)))$
 $= O(\underbrace{\lg_2 a_k + k \lg_2(\lg(m))}_{\in O(1)}) = O(\lg_2(\lg(m)))$

Compacted trie

mercoledì 15 marzo 2023 15:25

$$\mathbb{D} = \{s_1, \dots, s_n\}, \Sigma \text{ alphabet}$$

Def A trie is a binary search tree ($\Sigma = |\Sigma|$). Labels of edges are characters $\in \Sigma$ and labels of internal nodes are integers representing partial length of the substring. Leaves are strings of \mathbb{D} . Given an int. node s_{cur} is the string obtained concatenating the labels of both roots s_i .

Def A compacted trie CT is a trie without any internal node.

Oss For a CT: no internal unary node \Rightarrow #internal node \leq leaves

Optimization: • the substring (label of compacted trie edges) can be represented

in $O(n)$ space: $s' \mapsto \langle i, \text{start}, \text{end} \rangle$
if $s' = s_i[\text{start}, \text{end}]$

• the node can be represented by a perfect hash table which stores
only the branching characters \Rightarrow $\begin{cases} O(n) \text{ branching time} \\ \text{space optimal} \end{cases}$

Trie (T)	Space
	$O(N+n)$ = $O(N)$ <small>int. nodes \nwarrow leaves</small>
Compacted trie (CT)	$O(n+n)$ = $O(n)$ <small>int. nodes \nwarrow leaves</small>

2-level indexing

1) Strings in Disk:



$\text{RETRIEVE1}(Q, CT, D) // Q[1, p], D \text{ disk}$

$M \leftarrow \text{DOWNWARD_TRAVERSAL}(Q, CT); // \text{if int. node } s_{\text{cur}} = Q[1, p]$

return LEAVES(SUBTREE($M, CT), D);$

- Extra space: $O(\frac{n}{B})$
- Time: $O(p + n_{\text{sc}})$
- I/Os: $O(p + \frac{n_{\text{sc}}}{B})$

$\text{SEARCH}(Q, CT, D) // Q[1, q]$

$M, p \leftarrow \text{DOWNWARD_TRAVERSAL}(Q, CT); // s_{\text{cur}} = Q[1, p], p \leq q$
(no more steps into an edge)

$j \leftarrow \text{COMPARE}(\text{NEXT}(s_{\text{cur}}), Q[p+1]); // S_j \neq Q[p+1]$

return SCAN(R[:]). // result in ...

- Extra space: $O(\frac{n}{B})$ leaves
- Time: $O(p + FC_B(O_B))$

$\cup \dots \cup \text{NEXT}(\text{SCAN}(Q), Q_{[P+1]})$; // $S_i \neq Q \neq S_{i+1}$
 return $\text{SCAN}(B_j)$; // return position of Q

$O(p + FC_B(O_B))$
 \uparrow
 CT traversal
 \uparrow
 scan the block

$\text{RETRIEVE}(Q, (T, D)) \approx \omega_{CT, p}$

$l \leftarrow \text{SEARCH}(Q\$, CT, D)$,
 $r \leftarrow \text{SEARCH}(Q\#, CT, D)$,
 return $D[l:r]$; // all strings on disk prefixed by Q

$\cdot I/Os: O(p + \frac{FC_B(O_B)}{B})$

$\cdot Extra\ space: O(\frac{n}{B})$

$\cdot Time: O(2p + \frac{2FC_B(O_B)}{B})$
 \uparrow
 2 SEARCH

$\cdot I/Os: O(2p + \frac{2FC_B(O_B)}{B})$

(*)

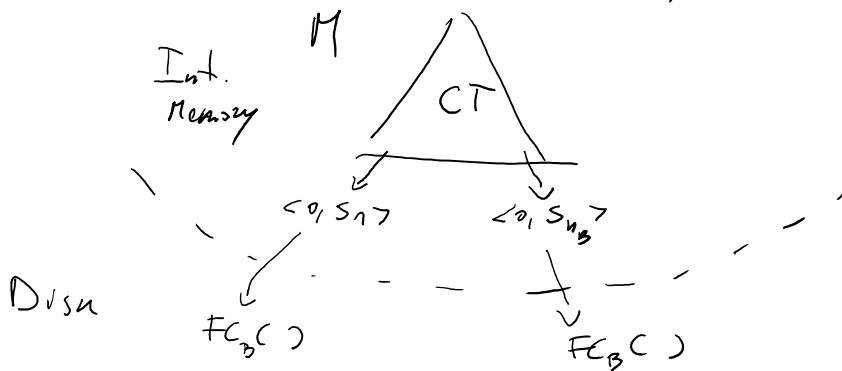
Thm (RETRIEVE) time is $O(p + n_{occ})$, #I/Os is $O(p + \frac{n_{occ}}{B})$ for prefix-search and retrieving strings prefixed by Q takes $O(FC_B(O_{occ}))$ time and $O(\frac{FC_B(O_{occ})}{B})$ I/Os. The space is $O(\frac{n}{B})$ for the CT with 2-level indexing.

Proof time: $O(p + n_{occ})$
 \uparrow
 downward traversal
 \uparrow
 # leaves
 \uparrow
 descending by α ($\text{SCAN} = Q$)

I/Os: $O(p + \frac{n_{occ}}{B})$
 \uparrow
 no jump
 \uparrow
 on disk
 are contiguous and ordered

To retrieve all the occurrence we must display them, therefore we need $FC_B(O_{occ})$ time and $\frac{FC_B(O_{occ})}{B}$ I/Os
 (using front loading, contiguous and ordered storage on Disk)

2) strings s_1, \dots, s_{n_B} on internal memory



$\text{SELECT}(Q, CT, D)$

//like search of 1

$\text{RETRIEVE}(Q, CT, D)$

//like RETRIEVE of 1

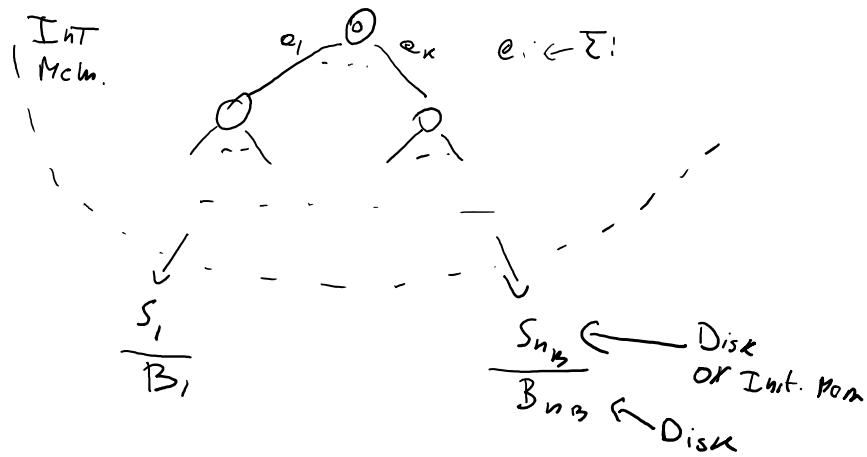
- time: $O(p + fC(B))$
- I/Os: $O(1)$ ← download traversal
confuses with string
in memory
- time: $O(2p + 2fC(B))$
- I/Os: $O(1)$ ($\frac{2 \text{ I/Os}}{K_2 \text{ searches}}$)

Patricia trie

mercoledì 15 marzo 2023 17:04

Does not depends if strings are in IM or disk.

Def patricia trie(PT) is a CT with an edge label only the first char. ($\in \Sigma$)



RETRIEVE(Q, PT, D) // $Q \in \Sigma^*$
 $l \leftarrow \text{SEARCH}(Q\$, PT, D);$
 $r \leftarrow \text{SEARCH}(Q\#, PT, D);$
 Return $D[l, r];$

Space: $O(\frac{n}{B})$
 Time: $O(p)$
 I/Os: $O(1 + F_B)$

SEARCH(Q, PT, D) // $Q \in \Sigma^*, q \in \Sigma$

$S, p \leftarrow \text{DOWNWARD_TRAVERSAL}(Q, PT)$ // S leaves of the subtree of u
 (a.s.t. $s_{[u]} = Q\epsilon, p \leq q$)
 $mismatch \leftarrow \text{LCP}(S, Q) + 1;$

$j \leftarrow \text{UPWARD_TRAVERSAL}(Q[mismatch], PT)$ // upper $Q[mismatch]$

$pos \leftarrow \text{SCAN}(Q, B_j);$
 return $pos;$

with $S^{[u+1]}$
 j is s.t. $S_j < Q < S_{j+1}$

The prefix-search need $O(p)$ time and $O(1 + F_B)$ I/Os, and $O(\frac{n}{B})$ space.

Proof time: $O(\underbrace{p}_{\text{downward traversal}} + \underbrace{p + 1}_{\text{LCP}(S, Q)}) = O(p)$

$$I/O_s: \underbrace{\frac{O}{B} I/O_s}_{\text{. . .}} + \underbrace{\frac{1}{B} I/O_s}_{\text{. . .}} + \frac{p}{B} I/O_s = O(F_B + 1) I/O_s$$

$$\overline{I/O_s} : \underbrace{O(I/O_s)}_{\text{disk seek}} + \underbrace{I/O_s}_{\substack{\text{pick the} \\ \text{leaf (jump} \\ \text{to disk)}}} + \frac{P}{B} I/O_s = O(\frac{P}{B} + 1) I/O_s$$

space: $O(\frac{h}{B})$ leaves, boundary interval table $\Rightarrow O(\frac{h}{B})$ int. table
 $\Rightarrow O(\frac{h}{B})$ total space.

General

lunedì 13 febbraio 2023 20:06

Given text string $T[1..n]$, alphabet Σ^t s.t. $|\Sigma^t| = \sigma$ and $T[n+1] = \$$ (the tree)

We want to resolve 2 ops:

• Counting ($P_{1..p}$) = #occ of $P_{1..p}$ as substring of $T[1..n]$

• Retrieve ($P_{1..p}$) = $\{i \in \{1, \dots, n\} \mid T[i, i+p-1] = P_{1..p}\}$

Def $Suff_i := T[i..n]$, $SUF(T) = \{Suff_i(T)\}_{i=1, \dots, n}$

Oss $P_{1..p}$ substring of $T[1..n] \Leftrightarrow P_{1..p}$ prefix of $s \in SUFF(T)$

Proof $P_{1..p}$ substring of $T[1..n] \Leftrightarrow T[i, i+p-1] = P_{1..p}$
 $\Leftrightarrow \exists i \text{ s.t. } P_{1..p}$ prefix of $Suff_i$
 $\Leftrightarrow P_{1..p}$ prefix of $s \in SUFF(T)$. \square

Suffix array

giovedì 16 marzo 2023 01:17

$T[1, n]$, $SA[1, n] := [i_1, \dots, i_n]$ s.t. $\text{Suff}_{i_1}(T) \leq \dots \leq \text{Suff}_{i_n}(T)$

```

// BINARY SEARCH OF P as prefix of suff(T)
SUBSTRING-SEARCH(P, SA(T)) // P[1, p]
l ← 1; z ← n;
while (l ≠ z) {
    m ← ⌊(l+z)/2⌋;
    if (P > SuffixSA[m][1, p])
        l ← m+1;
    else
        z ← m;
}
if (P = SuffixSA[l][1, p])
    return l;
else
    return -1;

```

- Extra space: $O(n)$ and $O(n \lg n)$ bits
- Time: $O(p \lg n)$
- I/Os: $O(\lg n)$

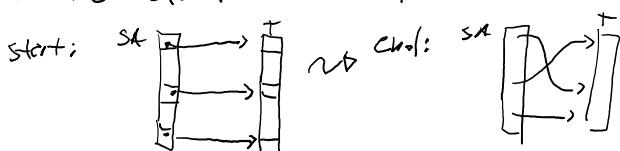
Lemma: given $T[1, n]$ e $SA[1, n]$, substring search of P in T needs
 $O(p \lg n)$ time, $O(\lg n)$ I/Os and $O(n \lg n + \lg n)$ bits for space ($O(n \lg n)$ extra space).
 Retrieving all the occurrences need $O(\text{occ})$ time.

Print time: $O(p \lg n)$ | I/Os: $O(\lg n)$ jumping at most for BS
 ↗ string comparison ↗ binary search

Total space: $O(n \lg n + n \lg n)$ | Retrieving: scan occ time ($\Rightarrow O(\text{occ})$)
 ↗ text ↗ SA (extra space) ↗ starting from l obtained
 by SUBSTRING-SEARCH

Suffix array construction

C-like method: T as array of char and SA as array of pointers to char



```

SA-BUILD(char* T, int n, char** SA) // C-like
for (i = 0; i < n; i++)
    SA[i] = T+i; // pointer to T[i, n]

```

QSORT(SA, n, sizeof(char*), SUFFIX-CMP);

SUFFIX-CMP(char** p, char** q)

- Working space: $O(n \lg n)$ bits
- Time: $O(n^2 \lg n)$
- I/Os: $O(\frac{n}{B} n \lg n)$

(+) using string pointers we don't need T 's position in memory
 (-) no locality in suffix comparison
 \Rightarrow not efficient in I/Os

SUFFIX-CMP($\text{char} \times p, \text{char} \times q$)

return STR-CMP($\&p, \&q$); // number > 0 if string $p \neq$ string q
= 0 if =
< 0 <

| \Rightarrow not efficient in I/Os

(*)

Then Using a comparison-based sort to build an array need $O(n^2 \log n)$ time
 $O(\frac{n}{B} n \log n)$ I/Os and $O(n \log n)$ working space.

Proof time: $O(n) + O(n) \cdot O(n \log n) = O(n^2 \log n)$

$\uparrow \quad \uparrow \quad \underbrace{\quad \quad \quad}_{\text{Qsort}}$

For length of string $T[i, j]$

I/O: $O(\underbrace{\frac{n}{B} n \log n}_{\text{T Qsort}})$

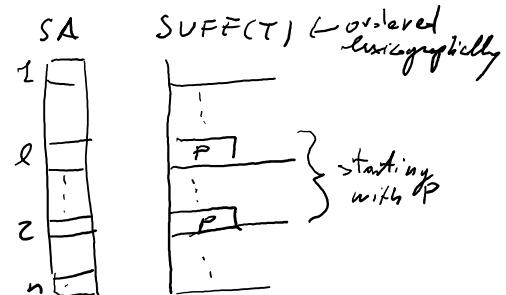
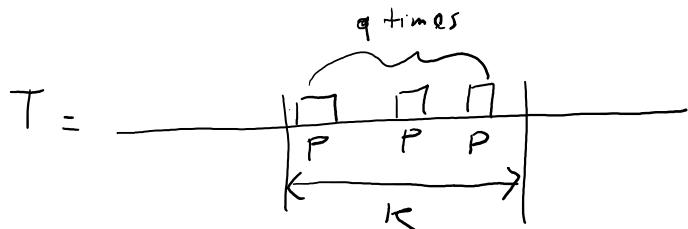
String are on disk

Working space:
 $O(n \log n)$ to store the pointers

Text mining

giovedì 16 marzo 2023 03:06

Example given text $T[1..n]$, his SA array and a pattern $P[i..j]$, establish whether it exists a substring of T of length k which includes at least q occurrences of P



CHECK(T, SA, P, k, q)

```

l, r ← SEARCH(P, SA); // binary search of P in SA
occ ← r - l + 1; // all occurrences of P in T
A[1..occ] ← ZEROS(1..occ);
for i = 1, ..., occ
    A[i] ← SA[l+i-1]; // copy the starting pos. of P in T

```

```

A ← SORT(A);
c ← 1; res ← false;
while (c < occ - q + 1 AND NOT(res)) { // search q occurrences of P that over k pos apart
    if (A[c+q-1] - A[c] ≤ k)
        res ← true;
    c++;
}
return res;

```

Prop time complexity is $O(p \lg n + occ \cdot ly occ)$

Proof

$$\text{time: } O(p \lg n) + O(occ) + O(occ \cdot ly occ) + O(occ - q) = O(p \lg n + occ \cdot ly occ)$$

\uparrow \uparrow \uparrow \uparrow
 bin. search of P in SA copy SA[l..r] in A[1..occ] sort A[1..occ] last for

□

General

lunedì 13 febbraio 2023 20:06

Source $S \rightarrow s_1, s_2, s_3, \dots$ ($s_i \in \Sigma^+$), Alphabet Σ

$$\text{Def Entropy of } S : H(S) := \underset{\substack{RV \\ \text{self information} \\ \text{of } s}}{\sum_{s \in \Sigma} p(s)} \cdot \lg_2 \frac{1}{p(s)} = \underset{RV}{\mathbb{E}[-\lg_2 p(s)]} \quad (\mathbb{E}[g(X)] = \sum_x g(x) \cdot p(x))$$

OSS $0 \leq H \leq \lg_2 |\Sigma|$, and H is max when p is uniform

Proof Using Lagrange multipliers with constraint $\sum_{s \in \Sigma} p(s) = 1$, we obtain that

H is maximum when $p(s) = \frac{1}{|\Sigma|}$ (uniformly distributed).

if p is uniform $H = \mathbb{E}[-\lg_2 \frac{1}{|\Sigma|}] = \mathbb{E}[\lg_2 |\Sigma|] = \lg_2 |\Sigma|$

$$p(s) = \frac{1}{|\Sigma|}$$
constant \square

Def given Coding C . Average Length of codeword is

$$L_C := \sum_{s \in \Sigma} p(s) |c_{ws}(s)|, \quad c_{ws}(s) = "codeword of s"$$

$p(s)$ can be $\begin{cases} \text{static (prob. distr. given)} \\ \text{semi-static (we know the frequency } f_{xs}, \forall x \in S) \\ \text{dynamic (we must compute the frequency dinamically)} \end{cases}$

Shannon thm for all prefix-free code C : $L_C \geq H$

OSS1 C is optimal $\Leftrightarrow L_C = H \Leftrightarrow |c_{ws}(s)| = \lg_2 \left(\frac{1}{p(s)}\right) \forall s \in \Sigma$

OSS2 C is optimal for prob. dist. $p(x) = 2^{-k_{wx}}$

Def A codeword C is universal iff $|c_{wx}| = O(\lg_2(x)) \quad \forall x \in \Sigma$

Def A word \mathcal{C} is universal iff $|\mathcal{C}_{x_1}| = O(g(x_1)) \forall x \in \mathbb{Z}$

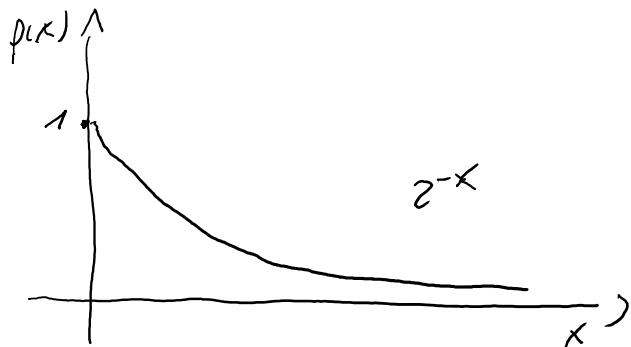
Unary code

giovedì 16 marzo 2023

17:24

Def for $x \geq 0$ the unary code of x is

$$U(x) := 0^{x-1} 1$$



Prop $|U(x)| = x$

Prop U is optimal for $p(x) = 2^{-x}$

Prop $|c_w(x)| = |U(x)| = x \Rightarrow$ optimal for $p(x) = 2^{-x}$

Decode(T):

```
next ← NEXT_BIT(T).
while (next ≠ EOF) {
    y ← 1;
    while (next ≠ 1) {
        y++;
        next ← NEXT_BIT(T).
    }
    OUTPUT(y);
}
return;
```

Fixed-length binary code

giovedì 16 marzo 2023 17:34

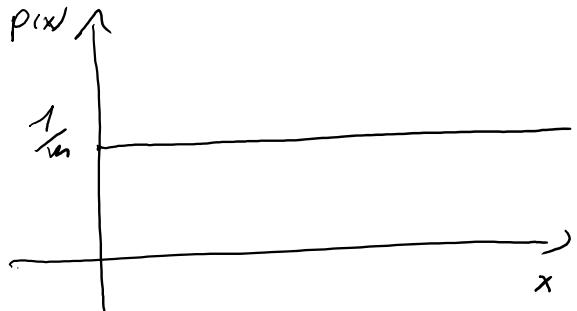
Def Given $m = \max_{s_j \in S} \{s_j\}$ and $M = 1 + \lceil \lg_2 m \rceil$, fixed-length binary code

of x is $B_M(x)$.

Prop $|B_M(x)| = M$

Prop It's optimal $\Rightarrow p(x) \propto \frac{1}{m}$

Proof $p(x) = 2^{-M} \propto 2^{-f_m} = \frac{1}{m}$ \square



DECODE(T)

$i \leftarrow 1;$

while ($i < \text{LEN}(T)$) {

$x \leftarrow T[i:i+M];$

$y \leftarrow \text{DECIMAL}(x, M);$ // from binary to decimal representation

$i \leftarrow i+M;$

OUTPUT(y);

}

return;

Gamma-code (Elias)

giovedì 16 marzo 2023 17:43

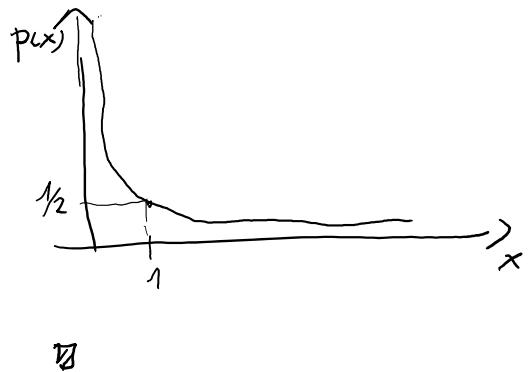
Def Given $x > 0$, $\gamma(x) := 0^{\ell-1}B(x)$, $\ell = |B(x)|$ ($\gamma(x) = 0^{|B(x)|-1}B(x)$)

Prop $|\gamma(x)| = 2\lfloor \log_2 x \rfloor + 1$

Proof $|\gamma(x)| = |B(x)| - 1 + |B(x)| = \lfloor \log_2 x \rfloor + \lfloor \log_2 x \rfloor + 1 = 2\lfloor \log_2 x \rfloor + 1$ \square

Prop Optimal $\Rightarrow p(x) \approx \frac{1}{2x^2}$

Proof $p(x) = 2^{-|\gamma(x)|} = 2^{-2\lfloor \log_2 x \rfloor - 1} \approx \frac{1}{2x^2}$



oss $\gamma(x)$ is universal prefix-free code

Proof $|\gamma(x)| = 2\lfloor \log_2 x \rfloor + 1 = O(\log x)$

DECODE(T)

next \leftarrow NEXT BIT(T); start $\leftarrow 1$;

while (next \neq EOF) {

$l \leftarrow 0$;

while (next $= 0$)

$\{ l \leftarrow l + 1, \text{next} \leftarrow \text{NEXT BIT}(T) \}$

start \leftarrow start + l ;

$X \leftarrow T[\text{start}, \text{start} + l + 1]$;

$Y \leftarrow \text{DECIMAL}(X, l + 1)$;

OUTPUT(Y);

start, \leftarrow start + $l + 1$.

```
OUTPUT(y);  
start ← start + l + 1];  
}  
}
```

Delta-code (Elias)

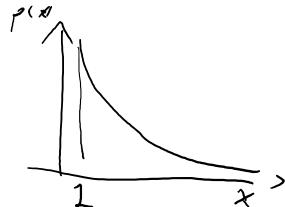
giovedì 16 marzo 2023 18:24

$$\text{Def } F_{\geq 1} \quad \delta(x) := \gamma(|B(x)|) B(x) \quad (= 0^{|B(|B(x)|)} B(|B(x)|) B(x))$$

$$\underline{\text{Pf}} \quad |\delta(x)| \approx 1 + \log_2 x + 2 \log_2 (\log_2 x) \text{ bits}$$

$$\underline{\text{Pf}} \quad |\delta(x)| = \underbrace{\lfloor \log_2(x) \rfloor + 1}_{B(x)} + \underbrace{2 \lfloor \log_2(\lceil \log_2(x) \rceil + 1) \rfloor + 1}_{\gamma(|B(x)|)} \approx 1 + \log_2 x + 2 \log_2 (\log_2 x) \quad \square$$

$$\underline{\text{Pf}} \quad \delta(x) \text{ optimal} \Leftrightarrow p(x) \approx \frac{1}{2x(\log_2 x)^2}$$



$$\underline{\text{Pf}} \quad 2^{-|\delta(x)|} \approx 2^{-(1 + \log_2 x + 2 \log_2(\log_2 x))} = \frac{1}{2^{2 \log_2 x} \cdot 2^{\log_2(\log_2 x)^2}} = \frac{1}{2x(\log_2 x)^2} \quad \square$$

Pf $\delta(x)$ is universal

$$\underline{\text{Pf}} \quad |\delta(x)| = 1 + \log_2(x) + 2 \log_2 \log_2(x) = O(\log_2(x)) + o(\log_2(x)) = O(\log_2(x))$$

DELTA-DECODE(T)

```

next ← NEXTBIT( $T$ );
start ← 1;
while (next ≠ EOF) {
    l ← 0;
    while (next = 0) {
        l++;
        next ← NEXTBIT( $T$ );
    }
    start ← start + l;
}

```

3. $\text{DECIMAL}(T[1:n], l)$
start ← start + l;
 $L \leftarrow \text{DECIMAL}(T[1:n], start, start + l + 1);$
start ← start + l + 1;
 $y \leftarrow \text{DECIMAL}(T[1:n], start, start + L);$
OUTPUT(y);

Rice code

giovedì 16 marzo 2023 18:54

Def given $x > 0$ and a parameter $\kappa \in \mathbb{Z}, \kappa > 0$

$$R_\kappa(x) := V_{(q+1)} B_\kappa(r)$$

$$(x-1 = q \cdot 2^\kappa + r)$$

$$\begin{cases} \text{quotient} & \text{remainder} \\ q = \left\lfloor \frac{x-1}{2^\kappa} \right\rfloor, \quad r = x-1 - q \cdot 2^\kappa \end{cases}$$

$\exists f \ x > 0$
 $c \in S$

$$R_\kappa(x) = V_{(q+1)} B_\kappa(r)$$

$$(x = q \cdot 2^\kappa + r)$$

Prop $|R_\kappa(x)| = \left\lfloor \frac{x-1}{2^\kappa} \right\rfloor + 1 + \kappa$ bits

Proof $|R_\kappa(x)| = \underbrace{q+1}_{V_{(q+1)}} + \underbrace{\kappa}_{B_\kappa(r)} = \left\lfloor \frac{x-1}{2^\kappa} \right\rfloor + 1 + \kappa$

Prop Optimal iff $p(x) = (1-p)^{x-1} p$ (geometric distribution $p(G=x)$)
 with parameter p of success

(N) PROOF)

Variabile-byte code

giovedì 16 marzo 2023 19:16

Def $x \geq 0$ $B(x) = B_k(x)B_7(x)\dots B_1(x)$, $k \in \mathbb{Z}$ and $7 + \dots + 7 + k = \lfloor \lg_2(x) \rfloor + 1$

$$VB(x) := 1 \underbrace{B_7(x)}_1 1 B_7(x) \dots 1 B_7(x) 0 B_7(x)$$

\uparrow
 $B_k(x)$ padded
with 0 to left

Prop $|VB(x)| = 8 \cdot \lceil \frac{|B(x)|}{7} \rceil$ bits.

Proof $B(x)$ is divided in blocks of 7 bits (the leftmost is padded) ($\lceil \frac{|B(x)|}{7} \rceil$ blocks)
and for each block is added 1 or 0 \Rightarrow 8 bit per each block \square

Prop $VB(x)$ optimal iff $p(x) \propto \sqrt[7]{\frac{1}{x^8}}$

$$\text{Proof } 2^{-|VB(x)|} = 2^{-8 \lceil \frac{|B(x)|}{7} \rceil} = 2^{-\frac{8}{7}(\lg_2(x) + 1)} \approx \sqrt[7]{\frac{1}{x^8}} \quad \square$$

Analysis es (S, C) -dense code

(S, C) -code with $S+C=2^b$, $b=8$

$$S = \left\{ \underbrace{b_1 \dots b_7}_{\text{strings}} \mid b_i \in \{0, 1\} \right\} = \{0, 1, \dots, 127\} = 128$$

$$C = \left\{ 1b_1 \dots b_7 \mid b_i \in \{0, 1\} \right\} = \{128, \dots, 255\} = 128$$

$128+128=2^8=256$

	# representat.	set
S	128	$\{0, \dots, 127\}$
CS	$128 \cdot 128$	$\{128, \dots, 128 + 128^2 - 1\}$
CCS	$128 \cdot 128 \cdot 128$	$\{128 + 128^2, \dots, 128 + 128^2 + 128^3 - 1\}$

etc...

(s,c)-dense code

venerdì 17 marzo 2023 01:06

$$s+c = 2^b, \begin{cases} s = |\text{Stoppers}| \\ c = |\text{Continuers}| \end{cases}$$

$$\begin{aligned} \text{Stoppers} &= \left\{ 0 \leq B_b(x) = x_1 \dots x_b \leq s \right\}_{x_i \in \{0,1\}} \\ \text{Continuers} &= \left\{ s \leq B_b(x) = x_1 \dots x_b \leq 2^b \right\}_{x_i \in \{0,1\}} \end{aligned}$$

Def $D(s, c)(x) := e_1 \dots e_s s, \quad e_i \in \text{Continuers}, \quad s \in \text{Stoppers}$

$D(s, c)$	# representati	Set
s	s	$\{0, \dots, s-1\}$
$e_1 s$	$c \cdot s$	$\{s, \dots, s+c \cdot s - 1\}$
$e_1 e_2 s$	$c^2 \cdot s$	$\{s+c \cdot s, \dots, s+c \cdot s + c^2 \cdot s - 1\}$

PForDelta code

venerdì 17 marzo 2023 01:18

$$S = S_1, \dots, S_n, \dots \text{ end} / \text{base} := \min_{S_i \in S} \{S_i\}$$

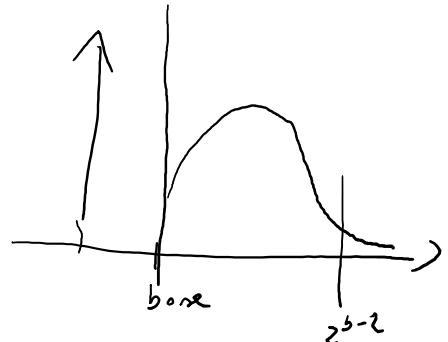
Choose b s.t. most of $S_i \in \{\text{base}, \dots, \text{base} + 2^b - 2\}$
 $(\approx 90\%)$

Def $\text{PFD}(x - \text{base}) = \begin{cases} B_b(x - \text{base}) & \text{if } x - \text{base} \in \{0, \dots, 2^b - 2\} \\ \text{ESC} & \text{if } x - \text{base} \geq 2^b - 1 \end{cases}$

$B_{b+w}(x - \text{base}) \leftarrow \text{fixed-size representation}$

Prop $|\text{PFD}(x - \text{base})| = \begin{cases} b \text{ bits} & \text{if } x - \text{base} \in \{0, \dots, 2^b - 2\} \\ b+w \text{ bits} & \text{if } x - \text{base} \geq 2^b - 1 \end{cases}$

\approx Good for Gaussian-like distribution



Interpolative code

venerdì 17 marzo 2023 01:51

- Adaptive (depends on distribution of the seq.)
- NO prefix-free

$$S = s_1, \dots, s_n \quad \forall i: s_i < s_{i+1} \quad \left(\text{otherwise } S' \text{ s.t. } s_i = \sum_{j=1}^i s_j \right)$$

// Start: $l=1, r=n, low=s_1, hi=s_n$
 $IC(S, l, r, low, hi)$ // interpolative-coding

```

if  $r < l$ 
    return {};
if  $l = r$ 
    return BC(S[l], low, hi);

```

```

m ← ⌊ $\frac{l+r}{2}$ ⌋;
A1 ← BC(S[m], low+m-l, hi-m+m); // better low and hi
A2 ← IC(S, l, m-1, low, S[m-1]);
A3 ← IC(S, m+1, r, S[m]+1, hi);
return [A1, A2, A3];

```

$\forall i: s_{i+1} > s_i \quad \forall i: l+1, \dots, r$	<ul style="list-style-type: none"> (+) Very compact (-) very slow (-) recursive (no cache friendly) (-) $s_{i+1} > s_i \quad \forall i: l+1, \dots, r$
---	--

BC(x, a, b) // "succinct" binary code

```

return B(x-a),
    ↗yi(b-a+i)

```

Elias-Fano code

venerdì 17 marzo 2023 16:48

$$S = s_1, \dots, s_n \quad s_i < s_{i+1}$$

Let $\begin{cases} m := (\max S) + 1 = s_n + 1 \\ b := \lceil \lg_2(u) \rceil \\ l := \lceil \lg_2(\frac{u}{n}) \rceil \\ h = b - l \end{cases}$

$$\forall x \in S \quad x \rightsquigarrow B_b(x) = \underbrace{H(x)}_h \underbrace{L(x)}_l$$

Def Given $S = s_1, \dots, s_n$ and constants m, b, h and l $EF(S) = (L, H)$

where $\begin{cases} L := L(s_1) \dots L(s_n) \\ H := B_0 \dots B_{2^h-1}, B_j := \sum_{i=0}^{s_j} L(s_i) \quad (x_j := |\{s_i : H(s_i) = j\}|) \\ \text{with } B_0 \dots B_{2^h-1} \text{ binary strings} \end{cases}$

Prop Given S with $|S| = n$, $|EF(S)| \leq 2n + h \lceil \lg_2 \frac{u}{n} \rceil$ bits, so for $x \in S$ $|EF(x)| \leq 2 + \lceil \lg_2 \frac{u}{n} \rceil$ bits on avg (regardless the distribution of S).

Proof $|S| = n \quad |EF(S)| = |L| + |H| = \underbrace{n}_L \underbrace{l}_H + \#L + \#H = n \lceil \lg_2 \frac{u}{n} \rceil + \underbrace{n + 2^h}_{\#1 = \#S} + \underbrace{\#0 = \#binarys}_{\#0 = \#1}$

$$= n \lceil \lg_2 \left(\frac{u}{n} \right) \rceil + n + \underbrace{2^{\lceil \lg_2 \frac{u}{n} \rceil} - \lceil \lg_2 \frac{u}{n} \rceil}_{\#0 = \#1}$$

$$\leq n \lceil \lg_2 \left(\frac{u}{n} \right) \rceil + n + n$$

For $x \in S \quad |EF(x)| = \underbrace{|EF(S)|}_n \text{ on avg} \leq 2 + \lceil \lg_2 \left(\frac{u}{n} \right) \rceil$

OSS If S is uniformly distributed in $[0, u]$, EF requires less than 2 bits in addition by optimal per integer.

Pr If S unif-distr. in $[0, u] \Rightarrow$ On avg $x \in S \quad |\mathcal{C}(x)| = \lceil \lg_2 \left(\frac{u}{n} \right) \rceil$ where \mathcal{C} is optimal

Augmenting the Elias-Fano coding (using RANK/SELECT data structure)
we can perform two operations:

- $\text{Acces}(i) = s_i \quad H_i = s_{i+1} \dots s_n$
- $\text{NextGEQ}(x) = \min \{ s_i | s_i \geq x \} \quad \forall x = 0, \dots, n$

The auxiliary RANK/SELECT data structure support efficiently the RANK or SELECT ops

$\left(\begin{array}{l} x \in [0, h], B \text{ binary vector} \\ \cdot \text{Rank}_x(i, B) = \# x \text{ in } B[1, i] \quad \forall i \leq |B| \\ \cdot \text{Select}_x(i, B) = \text{position of } i^{\text{th}} x \text{ in } B \end{array} \right)$

$\text{ACCESS}(i, L, H, b)$:

return $B_{\text{en}}(\text{SELECT}_1(i, H) - i) \sqcup [(i-1)l+1, il]$,
 to find the number of 0 before
 the i^{th} 1 \Rightarrow to find the bucket

Extra space: $O(|H|)$ bits
 $= O(n)$ bits
 Time: $O(1)$

$\text{NEXTGEQ}(x, H, L)$:

$H(x), L(x) \leftarrow B_b(x); // H(x) first h bits, L(x) last l bits$
 $p \leftarrow \text{SELECT}_0(H(x), H) + 1;$

$r \leftarrow p - H(x); // \text{SELECT}_0(H(x), H) - H(x) \text{ gives } \# \text{ of elements}$
 repeat {
 of the first $H(x)$ buckets, so adding 1
 we have the position i of the first element of
 $y \leftarrow \text{ACCESS}(i); \quad (H(x))\text{-th bucket}$
 $i++.$
 $(j+1)\text{-th bucket} = B_{\text{en}}(j))$

3 until ($y \geq x$) // return first value $\geq x$ in S

return y .

Extra space: $O(|H|) = O(n)$ bits

Time: $O(1)$

Worst case $O(\frac{n}{h})$

Worst case

+ binary search
 instead of
 scan (for), $O(\lg(\frac{n}{h}))$

Thm (*)_2 Using auxiliary select₂ data structure on H, ACCESS need $O(n)$ extra space and $O(1)$ time
Proof for this space we use an auxiliary select₂ data structure on H
 that occupy $O(|H|)$ extra space.

$$|H| = \# 1 + \# 0 = n + 2^{l_H} = h + 2^{\lfloor \lg_2 n \rfloor - \lfloor \lg_2 \frac{n}{h} \rfloor} \approx n + 2^h = 2n = O(n)$$

$\Rightarrow O(|H|) = O(n)$ extra space.

the time is $O(1)$ exploiting aux. select₂ data structure

(*)² Then using auxiliary Select, J.S. structure on H, NEXTFET needs $O(n)$ extra space and $O(n)$ time on avg. and $O(\frac{m}{n})$ in worst case ($O(\lg \frac{m}{n})$ in worst case using binary search in bucket)

Proof Space: using auxiliary select J.S. that $\text{array } O(|H|) = O(n)$ extra space

Time: perform at most $\text{size(bucket)}(n)$ times access operations (cycle FOR).

$$\text{So on avg. } \text{size(bucket)}(n) = \frac{\# \text{buckets}}{n} = \frac{\# O}{n} = \frac{2^k}{n} \approx \frac{2^k}{2^n} = O_{(1)}$$

In the worst case the $\text{size(bucket)} = 2^k = 2^{\lceil \lg \frac{m}{n} \rceil} = O(\frac{m}{n})$

\uparrow
after 2^k we have a new bucket

Using binary search in the bucket we have $O(\lg \frac{m}{n})$

Rank and Select

lunedì 13 febbraio 2023 20:09

$B \in \{0, 1\}^m$ binary array

Def $\text{Rank}_1(B, k) := \sum_{i=1}^k B[i]$ ($\text{Rank}_0(B, k) := k - \text{Rank}_1(B, k)$)

Def $\text{Select}_1(B, k) :=$ position of the k -th occurrence of 1 ($\text{Select}_0(B, k) =$ pos. of k -th 0 in B)

Thm \exists Rank data structure and Select data structure that takes $O(|B|) = O(m)$ bits and support respectively Rank and Select operation in $O(1)$ time

Proof See for rank succinct ranking implementation (Rank data structure)
For Select similar to rank (Select data structure) .



Succint solution for Rank (via Rank data structure)

giovedì 23 marzo 2023 02:35

$B[1..n]$ binary array

- Split $B[1..n]$ in $\frac{n}{z}$ block of size z :

$$B_i = B[z(i-1)+1, z \cdot i] \quad i = 1, \dots, \frac{n}{z}$$

$$|B_i| = z$$

For each B_i we store the absolute rank r_i ($\neq 1$ from start to end of B_i)

$$\begin{cases} r_0 = 0 \\ r_i = \text{Rank}_1(B_i, z) + r_0 & \text{if } i = 1, \dots, \frac{n}{z} \end{cases}$$

- Split each B_i in $\frac{z}{w}$ blocks of size w :

$$B_{i,j} = B[z(i-1) + w \cdot (j-1) + 1, z \cdot (i-1) + w \cdot j] \quad j = 1, \dots, \frac{z}{w}$$

$$|B_{i,j}| = w$$

For each $B_{i,j}$ we store the relative rank $r_{i,j}$ ($\neq 1$ from the start of $B_{i,1}$ to end of $B_{i,j}$)

$$r_{i,j} = \text{Rank}_1(B_{i,1}, w) + \dots + \text{Rank}_1(B_{i,w}, w)$$

$$i = 1, \dots, \frac{n}{z}$$

$$j = 1, \dots, w$$

Four Russians trick (Rank data structure)

- Construct Look-up table R of size 2^w rows $\times w$ columns

block	1	2	...	k	...	w
00...0						
00...01						
1						
$B[1..w]$						
1						
11...1						
$\text{Rank}(B[1..w], k)$						

$R[1..2^w][1..w] =$

all possible
block of size w

Def the look-up table R for all possible blocks of size w is the Rank data structure.

Ques Given a particular $B_{(z, w)}$ block of size w , we have

$$R[B_{(z, w)}][k] = \text{Rank}(B_{(z, w)}, k) \quad \text{for } k=1, \dots, w$$

Succinct implementation of Rank

// R = rank data structure, $B_{(z, h)}$, x index $1 \leq x \leq h$
 $\text{RANK}(B, x)$

$B_{i,j} \leftarrow \text{SEARCH}(x, B);$ // find the block where the index k is in
 $(z(i-1)+w \cdot j-1)+1 \leq k \leq z(i-1)+w \cdot j)$

$k \leftarrow \text{RELATIVE_POSITION}(x, B_{i,j});$ // position of index k in $B_{i,j}$
 $(\text{Start counting from } 1)$

return $r_{i-1} + r_{i,j-1} + R[B_{i,j}, k];$

• Extra space: $O(n)$ bits
 • Time: $O(1)$

(*) Thm Extra space for succinct rank is $O(n)$ bits using $z = \lg^2 n$ and $w = \frac{1}{2} \lg n$.
 Time needed updating the Rank d.s. is $O(1)$.

Proof Extra space = Space($\{r_i\}$) + Space($\{r_{i,j}\}$) + Space(R) (in bits)

$$= O\left(\underbrace{\frac{h}{z} \lg n}_{\# r_i \text{ each } r_i} + O\left(\underbrace{\left(\frac{n}{z} \cdot \frac{z}{w}\right)}_{\# r_{i,j} \text{ each } r_{i,j}} \cdot \lg z\right) + \text{Space}(R)\right) \quad \text{Rank data structure}$$

$$= O\left(\frac{h}{z} \cdot \left(\lg n + \frac{z}{w} \lg z\right)\right) = O\left(\frac{h}{\lg^2 n} \cdot \left(\lg n + \frac{\lg^2 n}{\frac{1}{2} \lg n} \lg(\lg^2 n)\right)\right)$$

$$= O\left(\frac{h}{\lg^2 n} \cdot \left(\lg n + 2 \lg n \cdot \lg(\lg^2 n)\right)\right) = O\left(\frac{h}{\lg^2 n} \cdot \left(\lg n + 4 \lg n \cdot \lg(\lg n)\right)\right)$$

$$= O\left(\frac{h}{\lg^2 n} \cdot \left(1 + 4 \lg(\lg n)\right)\right) = O\left(\frac{n \lg(\lg n)}{\lg n}\right) = O(n) + \text{Space}(R)$$

Space(R) = # rows \cdot # columns \cdot # bits per row (bits per cell)

$$= 2^w \cdot w \cdot \lg w = \sqrt{n} \cdot \frac{\lg n}{2} \cdot \lg\left(\frac{\lg n}{2}\right) \approx$$

$$= O(\sqrt{n} \cdot \lg h \cdot \lg(\lg n)) = o(h)$$

$$\Rightarrow \text{Extra space} = o(h) + o(h) = o(h)$$

Time : Look-up table + sum of values ($r_i, r_{i,j}$ + cell in R)
 $\Rightarrow O(1)$

Corollary Total space is $n + o(h)$ bits

Proof total space $|B| + o(|B|)$

Compressed Rank/Select based on Elias-Fano coding

giovedì 23 marzo 2023 02:35

B binary array $B = \{1, m\}$ $|B| = n$, $n = \#\{1\} \in B$ ($n = \text{Rank}(B, m)$)
 $\delta = \frac{n}{m}$ density

Fact Using Rank data structure and Select data structure respectively
 for Rank and Select operations we need

$\begin{cases} \text{Total Space: } |B| + o(|B|) = m + o(m) \text{ bits} \\ \text{Time: } O(1) \text{ for Rank or Select data structure} \end{cases}$

(*) If B is sparse ($\delta = \frac{n}{m} \approx 0$) we can do better.

Sparse solution via EF-code

$B_{\{1, m\}} \xrightarrow{\text{transform}} A_{\{1, n\}}$, $n = \#\{1\} \in B$ $A_{[i]} = \text{SELECT}_1(B, i)$

CONSTRUCT_A(B) // $B_{\{1, m\}}$

for ($i = 1, \dots, n$)

$A_{[i]} = \text{SELECT}_1(B, i)$.

$\text{SELECT}_1(i, B)$

return $\text{ACCESS}(i, A)$; // $A_{[i]}$

$\begin{cases} \text{Total space: } 2n + n \lceil \lg_2 \frac{m}{n} \rceil + o(n) \text{ bits} \\ \text{Time: } O(1) \end{cases}$

$\text{RANK}_1(i, B)$

return $\text{INDEX}(\text{NEXTGEQ}(i+1, A), A) - 1$;

$\begin{cases} \text{Total space: } 2n + n \lceil \lg_2 \frac{m}{n} \rceil + o(n) \text{ bits} \\ \text{Time: } O(\lg \frac{m}{n}) \end{cases}$

(*) Then transforming B in A and applying EF-coding to B we perform $\text{SELECT}_1(i, B)$ in $O(1)$ time taking $2n + n \lceil \lg_2 \frac{m}{n} \rceil + o(n)$ bits

Proof Time: $O(1)$ time for $\text{ACCESS}(i, A)$ exploiting Select data structure

for H ($\text{EF}(A) = (L, H)$) and select value on L (see ACCESS)

Total space: $|\text{EF}(A)| + (\text{Select}_1 \& c)$

$$\begin{aligned}
 &= |L| + |H| + (\text{Select}_1 \text{ d.c.}) = \\
 &\leq \underbrace{h \lceil \lg_2 \frac{m}{n} \rceil}_{|L|} + \underbrace{2h}_{|H|} + \underbrace{o(n)}_{O(|H|)} \quad \text{bits}
 \end{aligned}$$

20

Thm (*) As before $\text{RANK}_1(i, B)$ take $O(\lg \frac{m}{n})$ time explicitly selects and Select_1 data structure build on H . The total space used is $2h + h \lceil \lg_2 \frac{m}{n} \rceil + o(n)$ bits

Prf If time is given by $\text{NEXTBET}(i+1, A)$, so in worst case $O(\lg \frac{m}{n})$, explicitly Select_0 and Select_1 data structure build on H ($u = \max A + 1$ see EF-code).
 $u = \max A + 1 \leq m + 1 \Rightarrow$ in worst case $O(\lg \frac{m}{n})$.

- Total space: $|EF(A)| + \text{Select}_0 \text{ for } H + \text{Select}_1 \text{ for } H$

$$\leq \underbrace{h \lceil \lg_2 \frac{m}{n} \rceil}_{\text{L}} + \underbrace{2h}_{\text{H}} + o(n)$$

Analysis

R/S Data Structure			Via Elias-Fano Code		
	Time	Space		Time	Space
Rank	$O(1)$	$O(h) + o(h)$	$= \text{NextBETQ}$	$O(\lg \frac{m}{n})$	$2h + h \lceil \lg_2 \frac{m}{n} \rceil + o(n)$
Select	$O(1)$	$O(h) + o(h)$	$= \text{Access}$	$O(1)$	$2h + h \lceil \lg_2 \frac{m}{n} \rceil + o(n)$

- R/S data structure on B vs. EF-code on A w.r.t. density of B :

	B dense ($\frac{m}{n} \approx 1$)				B sparse ($\frac{m}{n} \ll 1$)			
	R/S take		EF-coste		R/S take		EF-coste	
	Time	Space	Time	Space	Time	Space	Time	Space
Rank	+	+	~	-	+	-	-	+
Select	=	+	=	~	=	-	=	+

Succint representation of binary tree (LOUDS-like)

giovedì 23 marzo 2023 02:37

Binary tree T , node $x \in T$ with label ℓ

Def x is a binary node
 \Leftrightarrow has two children

ENCODE (T)

$T' \leftarrow \text{COMPLETE_WITH_DUMMY_NODES}(T)$; // complete all non-binary nodes with dummy leaves

$T' \leftarrow \text{ENUMERATE_NODES}(T')$; // left to right and downward (all nodes)

for $i = 1, \dots, 2n+1$ {

if $T'[i] \neq \text{dummy}$ // i -th node

$B[i] = 1$,

else

$B[i] = 0$;

}

$T' \leftarrow \text{ENUMERATE_NO_DUMMY}(T')$;

for $i = 1, \dots, n$

$L[i] = \text{LABEL}(T'[2i], T'[2i+1])$; // $T'[2i]$ \uparrow $T'[2i+1]$ \uparrow (first node child even.

Def A binary tree is full \Leftrightarrow each node has 0 or 2 children

OSS T' is a full binary node

Prof each original node has 2 children, leaves (dummy node) have no children \square

OSS B binary array $|B| = 2n+1$ ($B[1, 2n+1]$)

L label array $|L| = n$ ($L[1, n]$)

Prof $|B| = m = \# \text{nodes of } T' = \# \text{internal nodes} + \# \text{leaves} = \# \text{original nodes} + \# \text{dummy nodes}$
 $= n + \underbrace{(h+1)}_{= h+1} = 2n+1$

T' is a full binary tree \square

OSS1 $T[\text{Rank}_1(B, i)] = T'[i]$

OSS2 Given a node i in T ($T[i]$) $T'[2i]$ is the left child in T'
 and $T'[2i+1]$ is the right child in T' . (T' complete)

OSS3 $T'[\text{Select}_1(B, i)] = T[i]$

OSS4 Given a node i in T' ($T'[i]$) $T[1 \dots i]$ is the parent in T .

QSS4 Given a node i in $T^1(T[i:i])$, $T[\lfloor \frac{i}{2} \rfloor]$ is the parent in T .

```

// Input: i = i-th node of T
LEFTCHILD(i) // Using oss1 and oss2
if BC[i] = 1
    return RANK1(B, 2i);
else
    j ← ⌊SELECT1(i, B)⌋;
    if BC[j] = 1
        return RANK1(B, 2j);
    else
        return RANK1(B, 2j + 1);

// Input: i = i-th node of T
RIGHTCHILD(i) // Using oss1 and oss2
if BC[2i + 1] = 1
    return RANK1(B, 2i + 1);
else
    j ← ⌊SELECT1(i, B)⌋;
    if BC[j] = 1
        return RANK1(B, 2j + 1);
    else
        return RANK1(B, 2j + 2);

// Input: i = i-th node of T
PARENT(i) // Using oss3 and oss4
j ← ⌊SELECT1(i, B)⌋;
if BC[j] = 1
    return j;
else
    return 2j;

```

Space: $2n+1+o(n)$ bits
Time: $O(1)$

(4)

Thm A binary tree of n nodes can be represented in $2n+1+o(n)$ bits and supports PARENT, LEFTCHILD, RIGHTCHILD in $O(1)$ time. (without an extra array of labels L)

Proof Space: $2n+1+o(n)$ bits
 $|B|$ R/S data structure for B ($o(|B|)$)

Time: PARENT, LEFTCHILD, RIGHTCHILD $O(1)$ using rank or select and exploiting R/S data structure.

10

General

lunedì 13 febbraio 2023 20:07

(Text compression)

Alphabets $\bar{Z}_i = \{\text{characters}\}$, $\sigma = |\bar{Z}_i|$
 ↑
 atomic items

Huffman coding

sabato 18 marzo 2023 00:19

- prefix-free code

$$\forall s \in \Sigma^*, p(s) \text{ given } \sum_{s \in \Sigma^*} p(s) = 1$$

HUFFMANN-TREE(Σ^*):

candidates $\leftarrow \{\Sigma^*\}$; // queue with priority $p(s)$

while ($|\text{candidates}| > 1$) {

$u_1, u_2 \leftarrow \text{POP}(\text{candidates}, 2)$; // extract 2 items with minimum prob.

$u \leftarrow \text{PARENT}(u_1, u_2)$; // create parent with child u_1 and u_2

$\text{ADD}((u, p(u_1 + u_2)), \text{candidates})$, // edges with 0 end/1
priority

}

space: $\Theta(|\Sigma^*|)$ space

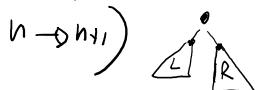
\hookrightarrow space high if $|\Sigma^*|$ high

(*) Proof the Huffman tree takes $\Theta(|\Sigma^*|)$ space

Proof Huffman tree is a full binary tree (each node has 2 or 0 children)

We prove that T full binary tree with $n = |\Sigma^*|$ leaves \Rightarrow # internal nodes = $n-1$.

$n=1$) $T = \bullet$ 0 internal node ok.



L, R two full subtrees of the root. ($\gamma(L) = \# \text{internal nodes of } L$)

L have respectively k_1 and k_2 leaves. ($\gamma(R) = \# \text{internal nodes of } R$)

leaves. $\Rightarrow \gamma(L) = k_1 - 1$, $\gamma(R) = k_2 - 1$ and $k_1 + k_2 = n+1$

$$\Rightarrow \gamma(T) = \gamma(L) + \gamma(R) + 1 = k_1 - 1 + k_2 - 1 + 1 = k_1 + k_2 - 1 = (n+1) - 1 = n$$

Huffman tree has $|\Sigma^*|$ leaves $\Rightarrow |\Sigma^*|-1$ internal nodes $\Rightarrow |\Sigma^*| + |\Sigma^*|-1 = |\Sigma^*| = \Theta(|\Sigma^*|)$ space

HUFFMANN-CODE($s \in \Sigma^*$):

$T \leftarrow \text{HUFFMANN-TREE}(\Sigma^*)$,

$\text{LABELING}(\Sigma^*)$,

$\text{DOWNWARD_TRAVERSAL}(T, s)$.

// print the list of edges path (labels of s)

(-) downward traversal
traversed
=> (choice
misses)

HUFFMANN-DECODE(B, T) // $B \in \mathbb{N}$, T Huffman tree

$m \leftarrow T$; // root

$i \leftarrow 1$;

while ($i < n$) {

while ($m \notin \Sigma^*$) {

$m \leftarrow \text{DOWNWARD}(T, B[i])$; // go through
edge $B[i]$

$i++$;

} $\text{OUTPUT}(m)$;

Def T_e is a binary tree that identify a code C by downward traversing. L_e is the average length of codeword and the average depth of the leaves of T_e

$$L_e = \sum_{s \in \Sigma^*} p(s) |C(s)| = \sum_{s \in \text{leaves}(T)} p(s) \text{depth}(s)$$

Fact We can build different Huffman tree $T_1 \neq T_2$ obtaining different coding $H_1 \neq H_2$.

We'll have $L_{H_1} = L_{H_2}$ but $\max_{s \in \Sigma} |H_1(s)| \neq \max_{s \in \Sigma} |H_2(s)|$

Aim Minimize $\left(\max_{s \in \Sigma} |H(s)| \right)$

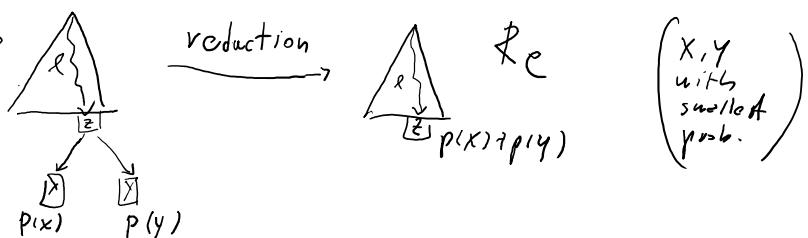
Strategy Among equal probability select the oldest node (the deepest one)

Lemma 1 $\mathcal{F} = \left\{ T_e \mid L_{T_e} = \min_{T} L_T, T \text{ with } |\Sigma| \text{ leaves} \right\} \Rightarrow \exists T'_e \in \mathcal{F} \text{ s.t.}$

$T'_e = \begin{array}{c} \triangle \\ \diagdown \quad \diagup \\ z \\ \square \quad \square \\ x \quad y \end{array}$ and $p(x), p(y)$ minimum prob. among $\{p(s) \mid s \in \Sigma\}$
 $x = \max \text{depth}, z = \text{parent}(x, y)$ (we can have other leaves at next depth)

Proof By contradiction Supposing that a leaf with minimum prob. is not at max depth.
By swapping it with a leaf at max depth but with higher prob, we obtain
a tree T with L_T lower \Downarrow

Lemma 2 Performing reduction on a tree T_e



$$\Rightarrow L_T = L_R + (p(x) + p(y))$$

Proof $L_T = \sum_{s \in \Sigma} \text{depth}(s) \cdot p(s) = \sum_{\substack{s \neq x, y \\ s \in \Sigma}} \text{depth}(s) \cdot p(s) + \text{depth}(x) \cdot p(x) + \text{depth}(y) \cdot p(y)$

$$= \sum_{\substack{s \neq x, y \\ s \in \Sigma}} \text{depth}(s) \cdot p(s) + (l+1)(p(x) + p(y))$$

$$L_T = \sum_{s \in \Sigma} \text{depth}(s) \cdot p(s) = \sum_{s \neq z} \text{depth}(s) \cdot p(s) + \text{depth}(z) \cdot p(z)$$

$$= \sum_{s \neq x} \text{depth}(s) \cdot p(s) + \ell \cdot (p(x) + p(y))$$

$$\Rightarrow L_T - L_R = p(x) + p(y)$$

□

Then H Huffman code $\Rightarrow L_H \leq L_e$ & prefix-free code C
 (equivalently: $L_{T_H} \leq L_{T_e}$)

Proof By induction on $|\Sigma'|$ ($\#$ leaves)

$$|\Sigma'|=2 \quad T_H = \begin{array}{c} \bullet \\ \diagdown \quad \diagup \\ \boxed{x} \quad \boxed{y} \end{array} \quad \text{true (Huffman tree is full binary tree)}$$

$|\Sigma'|-1 \rightarrow |\Sigma'|$ for Lemma 1 $\exists T_e$ optimal with $|\Sigma'|$ leaves s.t.

$$T_e = \begin{array}{c} \triangle \\ \downarrow \\ z \\ \diagdown \quad \diagup \\ \boxed{x} \quad \boxed{y} \end{array} \quad \begin{matrix} px, py \geq \min. prob. \text{ and } x, y \text{ at max depth} \\ \text{with } z = \text{parent}(x, y) \end{matrix}$$

$\Rightarrow L_{T_e} \leq L_{T_H}$ given T_H Huffman tree of $|\Sigma'|$ leaves and for the characteristic of H , T_H is s.t.

$$\begin{array}{c} \triangle \\ \downarrow \\ z \\ \diagdown \quad \diagup \\ \boxed{x} \quad \boxed{y} \end{array} = T_H$$

\leftarrow min prob.

Reducing T_e en T_H respectively in R_e and R_H , we have

$$L_{T_e} \stackrel{\text{Lemma 2}}{\geq} L_{R_e} + (p(x) + p(y)) \stackrel{\text{ind. hyp.}}{\geq} L_{R_H} + (p(x) + p(y)) \stackrel{\text{Lemma 2}}{=} L_{T_H}$$

\uparrow Huffman tree
with $\#$ leaves $|\Sigma'|$

$$\Rightarrow L_{T_H} \leq L_{T_e} \Rightarrow L_{T_H} = L_{T_e} \Rightarrow T_H \text{ is optimal}$$

□

oss C is Huffman $\Rightarrow C$ is optimal (L_e minimum)

OSS C is Huffman $\Rightarrow C$ is optimal (L_C minimum)

Thm $S = s_1, \dots, s_n \Rightarrow H \leq L_H < H+1$, where H is the entropy of S (At most 1 bit lost for each $s \in S$)

Proof

$$L_H = \sum_{s \in \Sigma} p(s) \cdot (H + c_s) = \sum_{s \in \Sigma} p(s) \cdot \underbrace{\log_2 \frac{1}{p(s)}}_{\in \mathbb{N}} \quad \left. \begin{array}{l} \text{and } H \text{ is optimal} \\ \text{if } H = \sum_{s \in \Sigma} p(s) \log_2 \left(\frac{1}{p(s)} \right) \end{array} \right\} \in \mathbb{R}$$

OSS The compress ratio is $\geq \frac{1}{\log_2 |\Sigma|}$

Proof $|H(x)|$ is at least 1 bit for $x \in \bar{\Sigma}$, and by $|\bar{\Sigma}|$ is the longest number of bit for one item \Rightarrow compression is $\geq \frac{1}{\log_2 |\Sigma|}$ \square

Optimizations Consider $\bar{\Sigma}^k = \{s_1, \dots, s_m | s_i \in \bar{\Sigma}\}$ alphabet of blocks of symbols
 $|\bar{\Sigma}^k| = |\bar{\Sigma}|^k$, if $x \in \bar{\Sigma}^k$ $H(x)$ lost only one extra bit
 $\Rightarrow H \leq L_H < H + \frac{1}{k}$ (at most $\frac{1}{n}$ bit lost for each symbol)

Canonical Huffman coding

sabato 18 marzo 2023 16:50

- Implicit building of the Huffman tree to save space

```
// Start: num array, symb array of lists, fc array
IMPLEMENT_CANONICAL_HUFFMANN_TREE( $\Sigma_i$ , H) // H is an Huffman code
    max ← MAX(|H(s)|, s ∈  $\Sigma_i$ ); // max length of codeword
    for l = 1, ..., max do
        symbol[l, :] ← {s | |H(s)| = l}; // starting from index 0
        num[l] ← LEN(symbol[l, :]); // lexicographic order
    }
    F[0] ← 0
    for l = max - 1, ..., 1
        F[l] ←  $\frac{f(l+1) + num[l+1]}{2}$ ; // Specie:  $O(|\Sigma_i| \log |\Sigma_i|)$ 
```

CODE (s) // $s ∈ \Sigma_i$:

```
(i, j) ← SEARCH( $s$ , symbol); // symbol[i, j] = s, i = level, j = j-th symbol
return  $B_i(f_{i+1} + j)$ ;
```

DECODE(C) // binary text

```
V ← NEXT_BIT(C);
while (V ≠ EOF) {
```

$l \leftarrow 1$;

while ($V < f[l]$) { // when $V \geq f[l]$, l is the level of V

$V \leftarrow 2V + \text{NEXT_BIT}(C)$;

$l++$;

}

OUTPUT(symbol[l, V - f[l]]);

V ← NEXT_BIT(C);

return;

(+) NO tree traversal \Rightarrow NO cache misses

(*) This Canonical Huffman coding, not building the tree, needs $O(\max^2 + |\Sigma_i| \cdot (\lfloor \lg |\Sigma_i| \rfloor + 1))$ bits, saving $O(|\Sigma_i|(\lfloor \lg_2 |\Sigma_i| \rfloor + 1))$ bits.

... , saving the country to the people .

Proof Don't need to store num array but only fc array and symb table

$$\Rightarrow \text{Space: } O(\max \cdot \max + |\Sigma| \cdot (\lfloor \log_2(|\Sigma|) \rfloor + 1)) = O(\max^2 \cdot |\Sigma| \cdot (\lfloor \log_2(|\Sigma|) \rfloor + 1))$$

$\# \text{level}$
 $\max \# \text{bits per item}$
 $\# \text{item}$
 $\max \# \text{bits per item}$
 bits

fC
 Σ
 symbol

$$Space(HuffmanTree) = O\left(\left(|\Sigma| - 1\right) \cdot 2w + |\Sigma| \cdot \lceil \lg(|\Sigma|) \rceil\right) \text{ bits}$$

int. nodes # bits per pointer leaves

12

Arithmetic coding

sabato 18 marzo 2023 20:18

- Prefix code
- Compressed output is not a concatenation of codewords
- a bit of the output can represent more than one input symbol

bit stream $B = b_1 \dots b_n$

$$\phi : \{0,1\}^* \longrightarrow [0,1] \subseteq \mathbb{R}$$

$$b_1 \dots b_n \longmapsto 0.b_1 \dots b_n = \sum_{i=1}^n b_i \cdot 2^{-i}$$

Def Dyadic fraction $\frac{x}{2^n}, x \in \mathbb{N}^+$

fact 1 $\forall b_1 \dots b_n, \phi(b_1 \dots b_n) = \underline{\text{val}}(b_1 \dots b_n) \in [0,1]$

Proof $\underline{\text{val}}(b_1 \dots b_n) = \frac{1}{2^n} \left(\sum_{i=1}^n b_i \cdot 2^{n-i} \right) = \sum_{i=1}^n b_i \cdot 2^{-i} = \phi(x)$

Fact 2 $\forall x \in [0,1], x = \frac{x}{2^n} \Rightarrow \phi^{-1}(x) = B_n(x)$

Proof As for fact 1

Fact 3 $\forall x \in [0,1], \phi^{-1}(x) = \text{CONVERTER}(x, k) = b_1 \dots b_k \quad (*)$

(*) CONVERTER (x, k) // $x \in [0,1], k \in \mathbb{N}^+$

(print first k bits
possibly the stream
has infinite length)

```

    i < 1;
    while (i <= k) {
        x <- 2 * x;
        if (x < 1)
            OUTPUT(0);
        else {
            OUTPUT(1);
            x <- x - 1;
        }
        i++;
    }
  
```

Compressing

$$S = \sigma_1 \dots \sigma_n, \sigma_i \in \Sigma, \quad \phi(\sigma_i) \text{ given } \forall i = 1, \dots, |\Sigma|$$

$$f_{\sigma_i} := \begin{cases} 0 & \text{if } i=1 \\ f_{\sigma_{i-1}} + p(\sigma_{i-1}) & \text{if } 2 \leq i \leq n \end{cases}$$

CODE (S)

```
(ln, sn) ← SEARCH(H_INTERVAL(S));
(x, rk) ← DYADIC([ln, ln+sn]); // select  $x \in [l_n, l_n+s_n]$ ,  $x = \frac{v}{2^k}$ 
B ← Bk(r); // w(r)
return OUTPUT((B, n));
```

CODE EFFICIENT(S)

```
(ln, sn) ← SEARCH(H_INTERVAL(S));
X ← ln +  $\frac{s_n}{2^k}$ ;
B ← CONVERTER(X,  $\lceil \log_2 \frac{v}{s_n} \rceil$ ); // w(s)
return OUTPUT((B, n));
```

// p(σ_i), for given and $n = |S|$

SEARCH_INTERVAL(S)

$s_0 \leftarrow 1$; $l_0 \leftarrow 0$; // starting size and lower

$i \leftarrow 1$;

while ($i \leq n$) {

$s_i = s_{i-1} \cdot p(\sigma_i)$; // new size

$l_i = l_{i-1} + s_{i-1} \cdot f_{\sigma_i}$; // new lower

$i++$;

return (l_n, s_n); // interval is [l_n, l_n+s_n]

Delegir enish

(B, n) B = b₁ ... b_n b_i ∈ {0, 1}, n = # iter desired

// Start: p(σ_i), f $_{\sigma_i}$ given

DECODE(B, n)

$x \leftarrow \phi(B)$; // $x = \frac{\text{val}(y)}{2^n}$

s₀ ← 1; l₀ ← 0;

i ← 0;

while ($i < n$) {

{l_i, l_i + s_i · f $_{\sigma_1}$, ..., l_i + s_i · f $_{\sigma_n}$ } ← PARTITION([l_i, l_i + s_i]);

```

 $\{l_i, l_i + s_i \cdot f_{\sigma_1}, \dots, l_i + s_i \cdot f_{\sigma_n}\} \leftarrow \text{PARTITION}([l_i, l_i + s_i]);$ 
 $j \leftarrow \text{SEARCH}(x, \{l_i, l_i + s_i \cdot f_{\sigma_1}, \dots, l_i + s_i \cdot f_{\sigma_n}\}); // x \in [l_i + f_{\sigma_j} \cdot s_i, l_i + f_{\sigma_{j+1}} \cdot s_i)$ 
 $b_i \leftarrow p(\sigma_j);$ 
 $s_{i+1} \leftarrow s_i \cdot p(\sigma_j);$ 
 $l_{i+1} \leftarrow l_i + s_i \cdot f_{\sigma_j};$ 
}
return;

```

Analysis

OSS (Compress) the size of the final interval does not depend on the particular ordering of input $S = \sigma_0, \dots, \sigma_n$

$$\text{Pr of } S_n = s_{n+1} \cdot p(\sigma_n) = \dots = \underset{\substack{\uparrow \\ 1}}{s_0 \cdot p(\sigma_1) \cdot \dots \cdot p(\sigma_n)} = \prod_{i=1}^n p(\sigma_i)$$

Def $x = 0.b_1 \dots b_n \in [0, 1]$, $\text{trunc}_d(x) := 0.b_1 b_2 \dots b_d$ (for $1 \leq d \leq n$)

Lemma $\forall x \in [0, 1], \text{trunc}_d(x) \in [x - 2^{-d}, x]$

$$\text{Pr of } x \in [0, 1] \Rightarrow x = \sum_{i=1}^{+\infty} b_i \cdot 2^{-i}, \text{trunc}_d(x) = \sum_{i=1}^d b_i \cdot 2^{-i}$$

$$\Rightarrow x - \text{trunc}_d(x) = \sum_{i=d+1}^{+\infty} b_i \cdot 2^{-i} = \sum_{i=1}^{+\infty} b_{d+i} \cdot 2^{d-i} \leq \sum_{i=1}^{+\infty} 2^{d-i} = 2^{-d} \cdot \left(\sum_{i=1}^{+\infty} 2^{-i} \right)$$

$$b_{d+i} = 1 \quad \forall i$$

$$= 2^{-d} \cdot \underbrace{\left(\frac{1}{1-\frac{1}{2}} - 1 \right)}_1 = 2^{-d}$$

$$\Rightarrow x - 2^{-d} \leq \text{trunc}_d(x) \leq x$$

\uparrow
trivial

Corollary $\text{trunc}_{\lceil \lg \frac{l+s}{2} \rceil}(l + \frac{s}{2}) \in [l, l+s] \leftarrow \text{final interval}$

$$\text{Pr of } \text{By lemma} \quad \left\{ \begin{array}{l} \text{Trunc}_{\lceil \lg \frac{l+s}{2} \rceil}(l + \frac{s}{2}) \geq l + \frac{s}{2} - 2^{\lceil \lg \frac{l+s}{2} \rceil} \geq l + \frac{s}{2} - 2^{\lg \frac{3}{2}} = l \end{array} \right.$$

Proof By lemma $\begin{cases} \text{Trunc}_{\lceil \lg \frac{n}{2} \rceil}(l + \frac{s}{2}) \geq l + \frac{s}{2} - 2^{\lceil \lg \frac{n}{2} \rceil} \geq l + \frac{s}{2} - 2^{\lceil \lg \frac{n}{2} \rceil} = l \\ // \leq l + \frac{s}{2} \end{cases}$

$$\Rightarrow \text{Trunc}_{\lceil \lg \frac{n}{2} \rceil}(l + \frac{s}{2}) \in [l, l + \frac{s}{2}] \subseteq [l, l + s)$$

□

Theorem $|A(S)| \leq 2^{nH}$, where $n = |S|$ and H is entropy of S

Proof $|A(S)| = \lceil \lg_2 \frac{n}{2} \rceil = \lceil 1 + \lg_2 \frac{1}{2} \rceil_p = \lceil 1 - \lg_2 s_n \rceil < 2 - \lg_2 s_n$
 $(s = s_n)$

$$\begin{aligned} S = \sigma_1, \dots, \sigma_n \quad &= 2 - \lg_2 \left(\prod_{i=1}^n p(\sigma_i) \right) = 2 - \sum_{i=1}^n \lg_2 p(\sigma_i) = 2 - \sum_{\sigma \in S} \text{occ}(\sigma) \lg_2 p(\sigma) \\ (s_n = 1) \quad &= 2 + \sum_{\sigma \in S} \text{occ}(\sigma) \lg \frac{1}{p(\sigma)} = 2 + n \sum_{\sigma \in S} \frac{\text{occ}(\sigma)}{n} \lg \frac{1}{p(\sigma)} \\ &\stackrel{n}{\approx} 2 + n \sum_{\sigma \in S} p(\sigma) \frac{1}{p(\sigma)} = 2 + nH \\ \frac{\text{occ}(\sigma)}{n} = \text{freq}(\sigma) \approx p(\sigma) \end{aligned}$$

□

Corollary $L_{AC} \leq H + \frac{2}{n}$ ($\text{only } \frac{2}{n} \text{ bit wasted on avg. per symbol}$)

avg $|A(\sigma)|$

	Compressor	On the fly	time	flexibility
Canonical Huffman	< (-)	No (-)	< (+)	Yes (+)
Arithmetical Coding	> (+)	Yes (+)	> (-)	No (-)

\nearrow
can use
dynamic modeling

for decompression
 any portion
 provided its
 first codeword

General

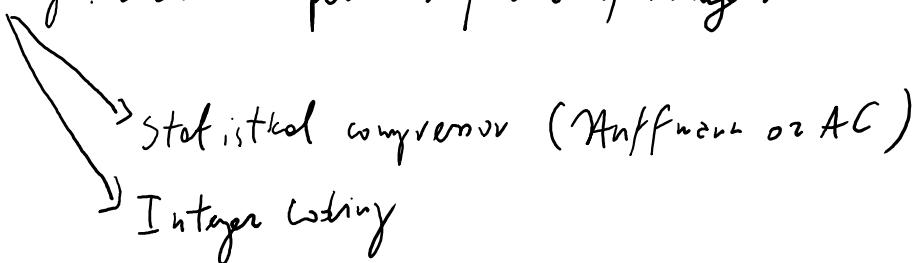
lunedì 13 febbraio 2023 20:08

- No probabilities of the source S
- given a dictionary of strings, we replace them with tokens which identify them in the dictionary

Compression

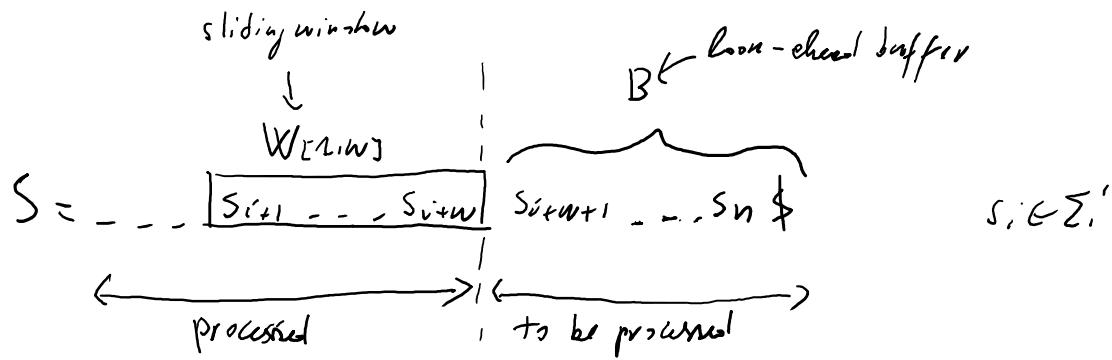
1) Parsing: input stream $S \longrightarrow S'$ sequence of integers
(atomic items)

2) Encoding: encode the parsed sequence of integers S'



Compressione

• Parsing:



$$S \longrightarrow S' = \underbrace{\langle d_1, |d_1|, c_1 \rangle, \dots, \langle d_k, |d_k|, c_k \rangle}_{\text{triplet } \in N \times N \times \Sigma'}$$

$\left\{ \begin{array}{l} d_i = \text{distance from start in } W \\ |d_i| = \text{length of prefix} \\ c_i = \text{next char in } B \end{array} \right.$

PARSING(S): // $S = s_1 \dots s_n$

$W \leftarrow \{ \}; B \leftarrow S; i \leftarrow 1;$

while ($i \leq n$) {

$\langle d, |d|, c \rangle \leftarrow \text{PARSE}(B, WB);$

$i \leftarrow i + |d| + 1;$

$W \leftarrow W[B[i, n)];$

$B \leftarrow B[i, n];$

 OUTPUT($\langle d, |d|, c \rangle$);

PARSE(WB, B) // the common prefix d in WB must be stored in W

$d \leftarrow 0; |d| \leftarrow 0; c \leftarrow B[1]; i \leftarrow 0;$

 while ($i < |B|$) {

$i++;$

 if ($LCP(W[i..i+|d|], WB, B) > |d|$) {

$|d| \leftarrow LCP(W[i..i+|d|], WB, B);$

$c \leftarrow B[|d| + 1];$

$d \leftarrow i;$

 }

}

3

$\} \quad \text{return } \langle d, |d|, c \rangle;$

• Encoding $S' = \langle d_1, |d_1|, c_1 \rangle, \dots, \langle d_n, |d_n|, c_n \rangle$

Encode $\{d_1, \dots, d_n\}$ (i.e. using Huffman)

Σ_i^1	occ	$\text{freq} \approx p$
d_1	$\#d_1$	$\#d_1/k$
\vdots	\vdots	\vdots
d_n	$\#d_n$	$\#d_n/k$

Output: $\langle H(d_1), |d_1|, c_1 \rangle, \dots, \langle H(d_n), |d_n|, c_n \rangle$

OSS $w=|W| \uparrow \Rightarrow$ Computation \uparrow , time complexity \uparrow
 $w=|W| \downarrow \Rightarrow$ \uparrow , \uparrow \downarrow

Decompression

$\text{DECOMPRESS}(S') : // S' = \langle d_1, |d_1|, c_1 \rangle, \dots, \langle d_n, |d_n|, c_n \rangle$

$S \leftarrow \{\}; i \leftarrow 1;$
 $\text{for } (i < |S'|) \{$

$S \leftarrow [S, \text{DECODE}(S, \langle d_i, |d_i|, c_i \rangle)];$
 $i \leftarrow i + |d_i| + 1;$

$\text{return } S;$

$\text{DECODE}(S, \langle d, |d|, c \rangle)$

$\text{for } i = 0, \dots, |d|-1 // \text{out} = \text{zeros}[1, |d|]$

$\text{out}[i+1] \leftarrow S[i|S|-d+1];$

$\text{return } \text{out}$

Compresssion

$$\text{PARSE}(CB, WB) = \begin{cases} \langle 0, BC \rangle; i \leftarrow i+1 & \text{if } d=0 \\ \langle d, id \rangle; i \leftarrow i+d+1 & \text{if } d>0 \end{cases}$$

* Parsing:

$\text{PARSING}(S) : // S = s_1, \dots, s_n$

$WB \leftarrow \{\}; B \leftarrow S; i \leftarrow 1;$

while ($i \leq n$) {

$\langle d, x \rangle \leftarrow \text{PARSE}(B, WB);$

if $d=0$ {

$BC[i] \leftarrow x; // \langle 0, BC \rangle$

$i \leftarrow i+1; \beta$

else if $d > 0$ {

$id \leftarrow x; // \langle d, id \rangle$

$i \leftarrow i+id+1; \beta$

$WB \leftarrow WB[i:n]; B \leftarrow BC[i:n]$

OUTPUT($\langle d, id \rangle, c$);

} return;

$\text{PARSE}(WB, B)$ {

$d \leftarrow 0; id \leftarrow 0, c \leftarrow BC[0];$

while ($d < |B|$)

$d++;$

if $LCP(W[i:i+d], B, B) > id$

$|d| \leftarrow LCP(W[i:i+d], B, B);$

if $d=0$

return $\langle d, BC[d+1] \rangle; // \text{return } \langle 0, BC \rangle$

else if $d > 0$

return $\langle d, id \rangle; // r$

* Encoding: A_c for LZ77 encode $\{d_1, \dots, d_n\}$

Decompression

$\text{DECOMPRESS}(S') // S' = \langle d_1, x_1 \rangle, \dots, \langle d_n, x_n \rangle \quad (x_n \text{ char or int})$

$i \leftarrow 1.$

DECOMPRESS (S^1) // $S^1 = \langle d_1, x_1 \rangle, \dots, \langle d_n, x_n \rangle$ (x_n char or int)

```
i ← 1;  
S ← {};  
while i < |S| {  
    out ← DECODE (S,  $\langle d_i, x_i \rangle$ );  
    S ← [S, out];  
}  
return out;
```

DECODE ($S, \langle d, x \rangle$)

```
if d = 0  
    out[i], j ← x; // x = char  
else if d > 0 {  
    for i = 0, ..., |d| - 1 // out = zeros [1, d]  
        out[i+1] ← S[|S|-d+i];  
}  
return out.
```

Gzip (fast LZ77)

lunedì 20 marzo 2023 18:31

- Hash-table (T) for speed up search on triplets

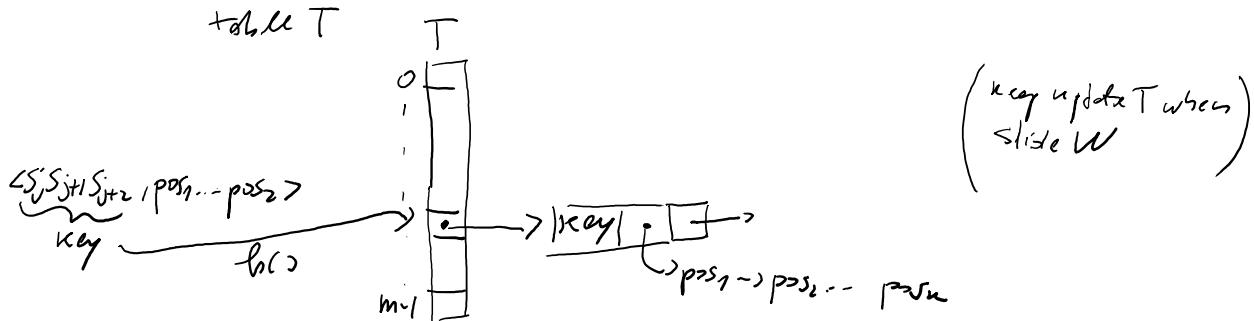
Compress

- Parsing:

$$S = \dots | \underbrace{s_{i+1} \dots s_{i+w}}_{W[i:w]} | \underbrace{B}_{s_{i+w+1} \dots s_n \$}$$

$$W[i:w] \rightarrow \text{dictionary } \mathcal{D} = \left\{ \langle \underbrace{\text{pos}_1, \dots, \text{pos}_n}_{\text{satellite data}}, \underbrace{\overbrace{s_i s_{i+1} s_{i+2}}^{\text{triplet}}}_{\text{key}} \rangle \right\}$$

\mathcal{D} stored in a hash table T



PARSING(S): // $S = s_1 \dots s_n$

$W \leftarrow \{ \}; B \leftarrow S; i \leftarrow 1;$

$T \leftarrow \text{EMPTY}(m)$; // empty hash-table of size m

while ($i \leq n$) {

$d \leftarrow B[i, \max(i+2, |B|)]$;

$\langle d, x \rangle \leftarrow \text{PARSE}(d, T)$;

 if $d = 0$ {

$B[i] \leftarrow x$; // $\langle 0, B[i] \rangle$

$l \leftarrow i+1$;

 } else if $d > 0$ {

$l \leftarrow x$; // $\langle d, l \rangle$

$l \leftarrow i+l$;

 } $\mathcal{D}_{del} \leftarrow \text{TRITRANS}(W[i:l])$;

$\mathcal{D}_{add} \leftarrow \text{TRIGRANS}(B[i:l])$;

$T \leftarrow \text{DELETE}(\mathcal{D}_{del}, T)$;

$T \leftarrow \text{INSERT}(\mathcal{D}_{add}, T)$;

$i \leftarrow l$;

} return T .

$\text{PARSE}(\alpha, T, WB)$

if $(\text{SEARCH}(\alpha, T))$

$\text{Pos} \leftarrow \text{SEARCH}(\alpha, T)$, // $\text{Pos} = pos_1, \dots, pos_n$

else

return $\langle \emptyset, \mathcal{B}_{C, j} \rangle$

for $i = 1, \dots, \text{LEN}(\text{Pos})$

$L_{C, i} \leftarrow \text{LCP}(WB[pos_i], WB[B, B])$ // $L_{C, 1, n}$ array of lengths of prefixes

$(i, |d|) \leftarrow \max(L)$, // $|d| = \max \{l\}$, $i = \arg \max \{l\}$

$d \leftarrow |WB| - pos_i + 1$,

return $\langle d, |d| \rangle$

• Encoding:

Building 2 Huffman

1) H_1 for alphabet $\Sigma_{i, 1} = \{l_{a, i}, c_j\}$

$\begin{cases} l_{a, i}, c_j \\ \uparrow \quad \nearrow \\ \text{from } \langle d_i, |d_i| \rangle \quad \text{for } \langle \emptyset, c_j \rangle \end{cases}$

same b, c
+ distinguish
 $\langle \emptyset, c \rangle$ with $\langle d, |d| \rangle$

2) H_2 for alphabet $\Sigma_{i, 2} = \{d_n\}$

$\uparrow \quad \nearrow \\ \text{from } \langle d_n, |d_n| \rangle$

Decompression

DECOMPRESS(B)

$S \leftarrow \emptyset;$

while ($blobn \neq EOF$) {

$x \leftarrow \text{Decode}(H_1, blobn)$

if x is literals // $x = c_j$ ($\langle \emptyset, c_j \rangle$)

$S \leftarrow [S, x]$.

else // x is a length of prefix ($\langle d_i, |d_i| \rangle$)

$d \leftarrow \text{Decode}(H_2, blobn)$

$S \leftarrow [S, \text{DECOMPRESS}(S, \langle d, x \rangle)]$.

}

OUTPUT(S)

- Build incrementally a dictionary

$$S = \underbrace{s_1 \dots s_j}_{W} \quad \underbrace{s_{j+1} \dots s_n}_{B}$$

permitted

not permitted

We use a dictionary of W $\mathcal{D} = \text{Dictionary}(W)$
 $\mathcal{D} = \{ \langle \text{id} : f \rangle \}$

Possing

+ empty string
 // Start: $\emptyset = \{ < \circ : \epsilon > \}$, $B = B[1, h] = S$
 $\text{PARSE}(S, \emptyset)$
 while ($B \neq \emptyset$) {
 $d, <\text{id}(d), c> \leftarrow \text{PARSE}(B, \emptyset)$;
 $\text{OUTPUT}(<\text{id}(d), c>);$ // $\text{id}(d)$ id of LCP in D
 $\emptyset \leftarrow \text{INSERT}(<\max(\text{id})+1 : \alpha(c), \emptyset);$
 $B \leftarrow B[d|+1, h];$
 }

PARSE(B, Q)

```

<id;f> ← LCP(B, ⊕) // f lcp with B in D (f = ⊕)
c ← B[|f| + 1];
return f, <id, c>;

```

To manage dictionary keep $\text{TRIE}(\emptyset) = (V, E)$
 $u \in V \quad u = \text{id}(\alpha)$

$$\exists u \xrightarrow{\alpha} u' \text{ s.t. } \alpha' = \alpha \subset \\ \begin{array}{c} u \\ \parallel \\ \text{id}(\alpha) \end{array} \quad \begin{array}{c} u' \\ \parallel \\ \text{id}(\alpha') \end{array}$$

Decoding

Start: $T = \langle id_1, c_1 \rangle, \dots, \langle id_n, c_n \rangle, \emptyset = \{ \langle \sigma : \epsilon \rangle \}$

DECODE(T, \emptyset):

```

for i = 1, ..., n {
    <id, x> ← T[i];
    f ← SEARCH(id, O); // <id:f> in O
    O ← INSERT(<maxid+1:f>, O);
    OUTPUT(fx);
}

```

General

martedì 21 marzo 2023 16:53

- Not a compression algorithm "per se"
- It's a permutation of the input string s.t. the obtained string is more suitable to be compressed via simple algorithms.
- The permutation forces locality homogeneous property to string
- Eventually Huffman or AC is performed to squeeze the bit-stream obtained from the simpler algorithms.

Forward transform

martedì 21 marzo 2023 17:12

Construct $s_i \longrightarrow bw(s)$

$S = s_1, \dots, s_n$ $s_i \in \Sigma^*$ ordered alphabet, & s.t. $\$ \not\in S \wedge s \in \Sigma^*$

// start: $M \in \mathbb{N}^{|\Sigma| \times |\Sigma|}$ matrix

FORWARDTRANSFORM(S)

$S \leftarrow S\$$; // append $\$$

for $i=0, \dots, n$ { // $M[i, :] = \text{string } S \text{ cycled}$

$M[i, :] = s_{1+i} \dots s_n \$ s_1 \dots s_{1+i-1}$

}

$M' \leftarrow \text{SORTROW}(M, \Sigma \cup \{\$\})$; // sort lexicographically rows of M w.r.t. $\Sigma \cup \{\$\}$

$L \leftarrow M'[:, n]$; // last column of M'

$(\hat{L}, r) \leftarrow \text{REMOVE}(\$, L)$; // remove $\$$ from L , $r = \#\$$ if $\$$ in \hat{L}

return (\hat{L}, r) ; // $bw(s)$

Backward transform

martedì 21 marzo 2023 17:26

$$S_b = s_1 \dots s_n \$$$

construct $bw(S) \xrightarrow{\quad} S$

$$M = \begin{bmatrix} s_1 & \dots & s_n & \$ \\ s_2 & \dots & s_n & s_1 \\ \$ & s_1 & \dots & s_n \end{bmatrix} \xrightarrow{\text{SORT ROW}} M' = \begin{bmatrix} F \\ \vdots \\ L \end{bmatrix}$$

Def $F := M'[c:][0:n]$, $L := M'[c:][n:](F = F[c, n], L = L[c, n])$

OSS1 $\forall i = 1, \dots, n \quad M'[c:i][r:] = \underset{\substack{\uparrow \\ S[k,n]}}{\text{Suff}_n} [S[c, k-1]]$, $k = SA[c:i]$

OSS2 $L[c:i] = S[k,n] \Rightarrow F[c:i] = S[c:c+1] \quad \forall i = 0, \dots, n, \quad i \neq n (L[n] = \$)$
 $(L[c:i] \text{ is the predecessor of } F[c:i] \text{ in } S)$

OSS3 $\exists \pi_F, \pi_L \text{ permutations s.t. } F = \pi_F(S\$), \quad L = \pi_L(S\$) \quad (\pi_c(i) = SA[c:i])$

OSS4 $\sigma \in S (\sigma \in \Sigma^*) \quad \text{occ}(\sigma, S) = \sigma_1, \dots, \sigma_n \quad \sigma_i = i\text{-th occurrence of } \sigma \text{ in } L$

$$\underset{(i \leq j)}{\text{pos}(\sigma_i, L)} < \text{pos}(\sigma_j, L) \iff \text{pos}(\sigma_i, F) < \text{pos}(\sigma_j, F)$$

Def LF mapping $LF: \{0, \dots, n\} \xrightarrow{\quad} \{0, \dots, n\} \text{ permutation}$
 $i \xrightarrow{\quad} j$

if $L[c:i] = \sigma_k \Rightarrow F[LF[c:i]] = F[c:j] = L[c:i] = \sigma_k$
 $\kappa\text{-th occurrence}$
 $\text{of } \sigma \text{ in } L$

//Start $C = C[0, |\Sigma|] = (0, \dots, 0)$ indexed by $\Sigma \cup \{\$\}$ in lexicographic order
CONSTRUCT_LF (Σ, L)

For $i = 0, 1, \dots, n-1$

$([C[i]]++) // (C[i]) = \# \text{occ}(\sigma) \quad (L[c:i] = \sigma)$

$\text{temp} \leftarrow 0, \text{sum} \leftarrow 0$

for $j = 0, 1, \dots, |\Sigma| \{ \quad // (C[i]) = \sum_{j < i} (C[j])$

$\text{temp} \leftarrow C[i];$

$C[i] \leftarrow \text{sum};$

$\text{sum} += \text{temp};$

}

for $i = 0, 1, \dots, n-1 \{$

• Extra space: $O(|\Sigma| + 1)$
 • Time: $O(n)$

```

    }  

    for i := 0, 1, ..., n-1 {  

        LF[i] = C[L[i]];  

    }  

        C[L[i]]++; // increment for the new occurrence
}

```

Backward Transform $bw(S) \xrightarrow{\quad} S$

BACKWARD (bw, Σ)

$(L, r) \leftarrow bw; // derive L from bw(S)$

$L \leftarrow \hat{L}[1, r-1] \# \hat{L}[n+1, n];$

$LF[0, n] \leftarrow \text{CONSTRUCT_LF}(\Sigma, L);$

$k \leftarrow 0;$

$i \leftarrow n-1; // from n-1 to not output &$

while $i \geq 0 \{ // L[LF[k]] precursor of F(LF[k]) = L[i]\}$

$S[i] \leftarrow L[i], // (at start: $S[n-1] = L[0]$)$

$k \leftarrow LF[k];$

$i--;$

}

• Space: $O(n)$

• Time: $O(n)$

Theorem S can be reconstructed from its BWIT in $O(n)$ time and space, eliciting possibly 1 cache-misses per symbol.

Move-To-Front transform (MTF)

martedì 21 marzo 2023 23:09

$$S = s_1 \dots s_n, \quad \Sigma^i = \{s_1, \dots, s_i\} \quad s_i \leq s_{i+1}$$

// Start $\mathcal{L} = (s_1, \dots, s_n)$ list of symbols, MTF = {}

$\text{MTF}(S) // S = s_1 \dots s_n$

for $i = 1, \dots, n$

$p \leftarrow \text{SEARCH}(\Sigma^i, \mathcal{L})$; // position of s_i in \mathcal{L}

$\text{MTF} \leftarrow [\text{MTF}, p]; S = \text{LCP}(p);$

if $p > 0$

$\mathcal{L}[0:p] \leftarrow S;$

for $j = 0, \dots, p-1$

$\mathcal{L}[i+1] \leftarrow \mathcal{L}[i:j];$

}

• Working space: $O(|\Sigma^i|)$

• Time: $O(|S|) = O(n)$

(*)
Then MTF takes $O(|\Sigma^i|)$ working space and $O(n)$ time

Proof Working space: $O(|\Sigma^i|) \leftarrow$ fast

Time: $O(|S|) = O(n) \leftarrow$ for cycle

Then $L_{MTF+\gamma} \leq 2H+1$ and $L_{H+\gamma} \leq 2L_H$ $\begin{pmatrix} \gamma = \gamma - \text{code} \\ H = \text{Huffman} \\ H = \text{entropy} \end{pmatrix}$
(No prefix)

Then $\exists S$ s.t. $|[MTF+\gamma](S)| < |H(S)|$ (of sfondo $\Omega(4^n)$)

Proof $S = 1^n 2^n \dots n^n \quad |S| = n^2, \quad p(1) = \dots p(n) = \frac{1}{n}$

For an optimal free-prefix code C $|C(x)| = \log_2 \frac{1}{p(x)} = \lg n$

$|C(S)| = |S| \cdot \lg n = n^2 \lg n$ bits $\approx |H(S)|$
optimal

MTF + γ : $S = 1^n 2^n \dots n^n \rightsquigarrow MTF(S) = 0^n 10^{n-1} 20^{n-2} 30^{n-3} \dots (n-1)0^{n-1}$

$\rightsquigarrow MTF(S)+1 \rightsquigarrow \gamma(MTF(S)+1)$

$$|MTF + \gamma(S)| = O(n^2) + O(n \lg_2 n) = O(n^2 + n \lg_2 n)$$

\uparrow
 $\textcircled{4} \quad (n^2) \geq 0$
 $\Rightarrow O(n^2) \text{ times } \gamma_{\substack{(1) \\ \text{1-bit}}}$

$n = \#x, x > 0, x < n$
 $|\gamma(\cdot)| = O(\lg_2 \cdot) \text{ bits} \Rightarrow |\gamma(\cdot)| \leq O(\lg_2 n)$

$$\Rightarrow |[MTF + \gamma](S)| = O(n^2 + n \lg_2 n) \overset{O(n^2)}{\leq} n^2 \lg_2 n = H(S)$$

$$\text{and } \frac{|H(S)|}{|[MTF + \gamma](S)|} = \frac{n^2 \lg_2 n}{O(n^2)} = \Omega(\lg_2 n)$$

Run-Length Encoding transform (RLE)

mercoledì 22 marzo 2023 00:24

$$S(1, n) = \underbrace{s_1 \dots s_1}_{n_1} \underbrace{s_2 \dots s_2}_{n_2} \dots \underbrace{s_k \dots s_k}_{n_k}$$

$$\sum_{i=1}^k n_i = n$$

$$RLE(S) = \langle n_1, s_1 \rangle, \dots, \langle n_k, s_k \rangle$$

$$RLE(S)$$

```

old ← S[1]; count ← 1;
for i = 1, ..., n { j ← i+1; new ← S[i];
    while (new = old) {
        count++;
        j++;
        new ← S[j];
    }
    i ← j;
    old ← new;
}

```

Thm $\exists S$ s.t. $|[RLE+\gamma](S)| \leq |H(S)|$ by a factor $\mathcal{O}(n)$

Proof $S = 1^n 2^n \dots m^n \quad |S| = n^2 \quad |H(S)| \approx |S| \cdot 2^{\frac{n}{2}} = n^2 \frac{1}{\sqrt{2}} \uparrow p(x) = \frac{1}{x}$

$$RLE(S) = \langle n, 1 \rangle, \dots, \langle n, m \rangle$$

$$|\gamma(\langle n, c \rangle)| = O(\lg n) \Rightarrow |[RLE+\gamma](S)| = O(n \lg_2 n)$$

$$\Rightarrow \frac{|H(S)|}{|[RLE+\gamma](S)|} = \frac{\Theta(n^2 \lg_2 n)}{O(n \lg_2 n)} = \mathcal{O}(n)$$

$BZP(S)$ $(\hat{L}, r) \leftarrow \text{FORWARDTRANSFORM}(S); // bw(S)$ $S' \leftarrow \text{MTF}(\hat{L})$

for $i = 1 \dots |S'|$ { // $S = S' + 1$ for $S[i] \neq 0$
 if $S'[i] \neq 0$

 $S'[i] = S'[i] + 1;$

{}

 $S' \leftarrow \text{RLFO}(S'); // RLE only on 0s of S'$ $S' \leftarrow \text{WHEELER}(S'); // Wheeler code only on lengths of RLFO$ return S' ; $\text{WHEELER}(n)$ Return $B_{1|B_{n+1}-1}(n+1); // we take n+1 so that we can drop the first bit ($n+1 \geq 2$)$ Efficiency

- Avoid construction of M ($\Theta(n^2)$ working space) for BWT (Forward transform)

OSS $L[i] = \begin{cases} S[SAC[i]-1] & \text{if } SAC[i] \neq 0 \\ \$ & \text{otherwise} \end{cases}$

Proof $F[i] = S[SAC[i]]$ and $L[i]$ precedes $F[i]$ in S

①

Theorem S s.t. $|S|=n$. Explicitly oss the forward transform time working space, time and I/Os needed for the SA construction (via qsort) (see SA-BUILD)

$$\begin{cases} \text{working space: } O(n) \text{ (O}(n \lg n \text{ bits)}) \\ \text{time: } O(n^2 \lg n) \\ \text{I/Os: } O(n \lg n) \end{cases}$$

$$| \pm \theta_s : O\left(\frac{h}{b} n \log n\right)^{\theta_s}$$