

APPROXIMATION ALGORITHMS

SUMMARIES FOR MULTISETS

RASMUS PAGM

UNIVERSITY OF COPENHAGEN



TODAY

- MOTIVATING EXAMPLE: WEBSITE VISITS
- MISRA-GRIES SUMMARY
- COUNT-MIN SKETCH

DATA SET: BROWSER LOG

2021-05-31 10:06: diku.edu
2021-05-31 10:07: google.com
2021-05-31 10:09: wikipedia.org
⋮

TRIVIAL SOLUTION:
USE A DICTIONARY
MAPPING SITE NAMES
TO COUNTS

1 USER: SMALL DATA

10^9 USERS: BIG DATA

↑
LINEAR SPACE
NOT FEASIBLE

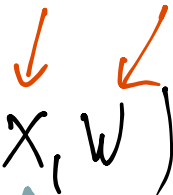
- QUESTIONS:
- WHAT ARE THE MOST POPULAR WEB SITES?
 - IN NUMBER OF PAGE VIEWS
 - IN TIME SPENT ON PAGES

GOAL: ANSWER SUCH QUESTIONS USING SUBLINEAR SPACE

MISRA-GRIES SUMMARY

- RECEIVES STREAM S OF PAIRS (x, w)

ITEM WEIGHT



SAME ITEM CAN APPEAR SEVERAL TIMES

- **TRUE** WEIGHT OF x : $w_x^* = \sum_{(x,w) \in S} w$

- STORES A SUMMARY WITH SPACE FOR k PAIRS (x, w_x)

- NEED TO MAINTAIN SUMMARY UNDER UPDATES

CAN BE MUCH SMALLER THAN #DISTINCT ITEMS

NOT NECESSARILY EQUAL TO TRUE WEIGHT OF x IN S

- ABILITY TO MERGE TWO SUMMARIES S_1, S_2 TO FORM SUMMARY OF THE "MULTISET UNION":

$$w_x^* = \sum_{(x,w) \in S_1 \cup S_2} w = \sum_{(x,w) \in S_1} w + \sum_{(x,w) \in S_2} w$$

MISRA-GRIES UPDATES

CONSIDER $k=3$, ALL WEIGHTS 1

INPUT STREAM: a, b, a, c, d, e, a, d, f, a, d

PAUSE AND THINK:
HOW CAN WE MERGE
TWO MG SUMMARIES
INTO ONE OF SAME SIZE?

T

ITEM	COUNT
a	2
f	0
d	1

NO SPACE
IN SUMMARY
- SMALLEST
COUNT MUST GO!

SUBTRACT
 $m=1$ FROM
ALL COUNTS,
KICK OUT
ZERO COUNT
ITEMS

INSERT (x, w) :

- IF $x \in T$, INCREASE COUNTER BY w
- ELSE, IF A SLOT WITH COUNT ZERO EXISTS, WRITE (x, w)
- ELSE, DECREASE ALL COUNTERS BY $m = \min w(x, w) \in T$

AND THEN INSERT $(x, w - m)$

MISRA-GRIES ANALYSIS

• INVARIANTS: $w_x \leq w_x^*$

ASSUMING POSITIVE UPDATES ($w \geq 0$)
 WE ONLY INCREASE w_x WHEN
 A PAIR (x, w) IS SEEN.

• PROOF:

LET $M = \sum_{(x, w) \in T} w_x$, THEN WE CLAIM THAT

$$w_x \geq w_x^* - \frac{\|w^*\|_1 - M}{k+1}$$

$$w_x \geq w_x^* - \frac{\|w^*\|_1}{k+1}$$

LAPTOP
 10^7 PAIRS

ERROR
 ≤ 1000

OFTEN MUCH
 BETTER IN PRACTICE

$$\|w^*\|_1 = \sum_x |w_x^*|$$

WEBSITE VISITS

TOTAL SUM
 OF ALL TRUE
 WEIGHTS
 (ABS. VALUE)

TOTAL WEIGHT OF DISCARDED ELEMENTS

NOT IN T ← THE CASE $x \notin T$
 EASILY HOLDS BY INDUCT.

WHAT HAPPENS AT UPDATE WITH NEW ELEMENT (x, w) ?

- NEW VALUE $M' = M - (k+1)m + w$

- FOR $x \in T$: $w_x' = w_x - m \geq w_x^* - \frac{\|w^*\|_1 - M}{k+1} - m = w_x^* - \frac{\|w^*\|_1 - M'}{k+1}$

INDUCTION HYPOTHESIS

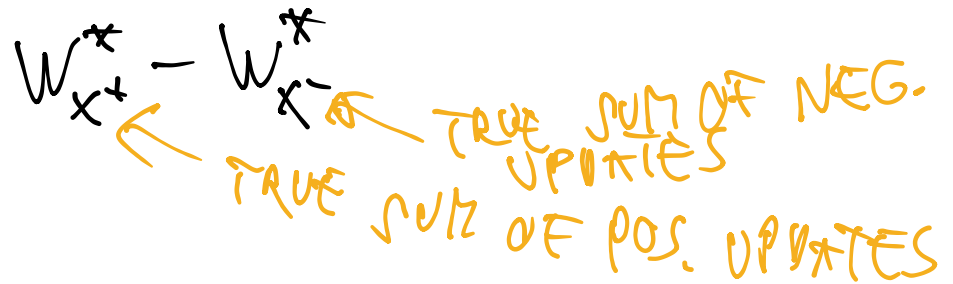
MERGE:
 SIMILAR
 ANALYSIS,
 SEE BOOK

MISRA-GRIES VARIANTS

- KEEP TRACK OF HIGHEST WEIGHT SEEN FOR EACH ITEM x IN T . (HEURISTIC TO DECREASE ERROR IN "TYPICAL" CASE)

- KEEP TRACK OF $\|w^*\|_1$ TO SUPPORT UPPER BOUNDS ON w_x^* , USING $w_x^* \leq w_x + \frac{\|w^*\|_1 - M}{k+1}$.

- SUPPORT NEGATIVE UPDATES BY THINKING OF w_x^* AS A DIFFERENCE

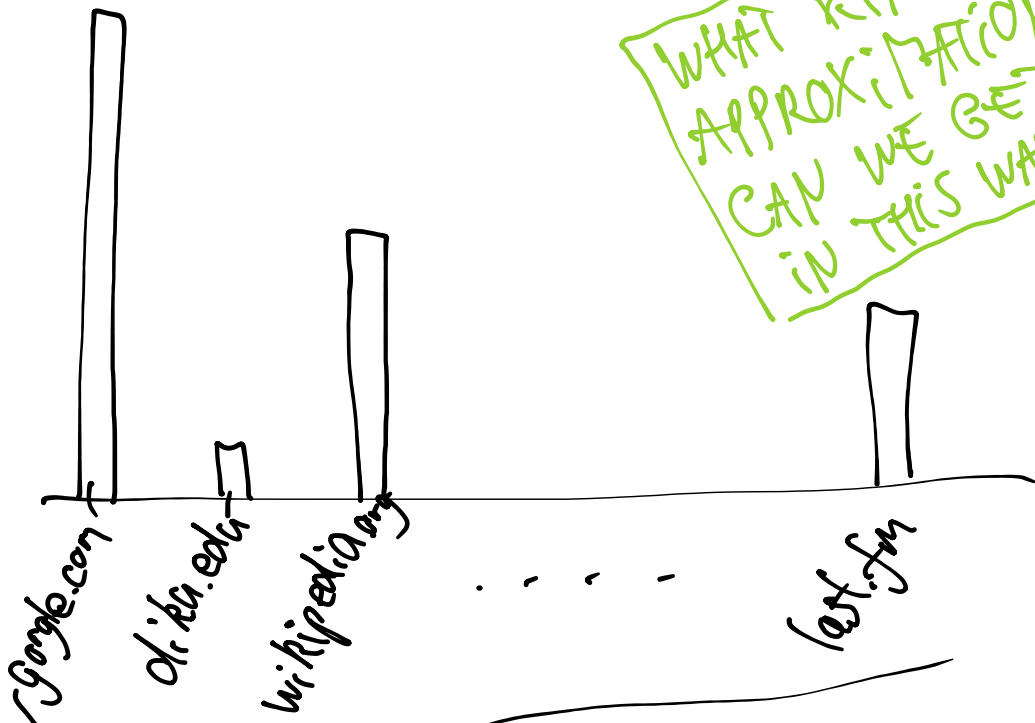


- "SPACESAVING" ALGORITHM:

KEEP TRACK OF UPPER BOUNDS RATHER THAN LOWER BOUNDS (SEC. 3.3 IN BOOK) EQUIVALENT TO MISRA-GRIES

RANDOMIZED HISTOGRAMS

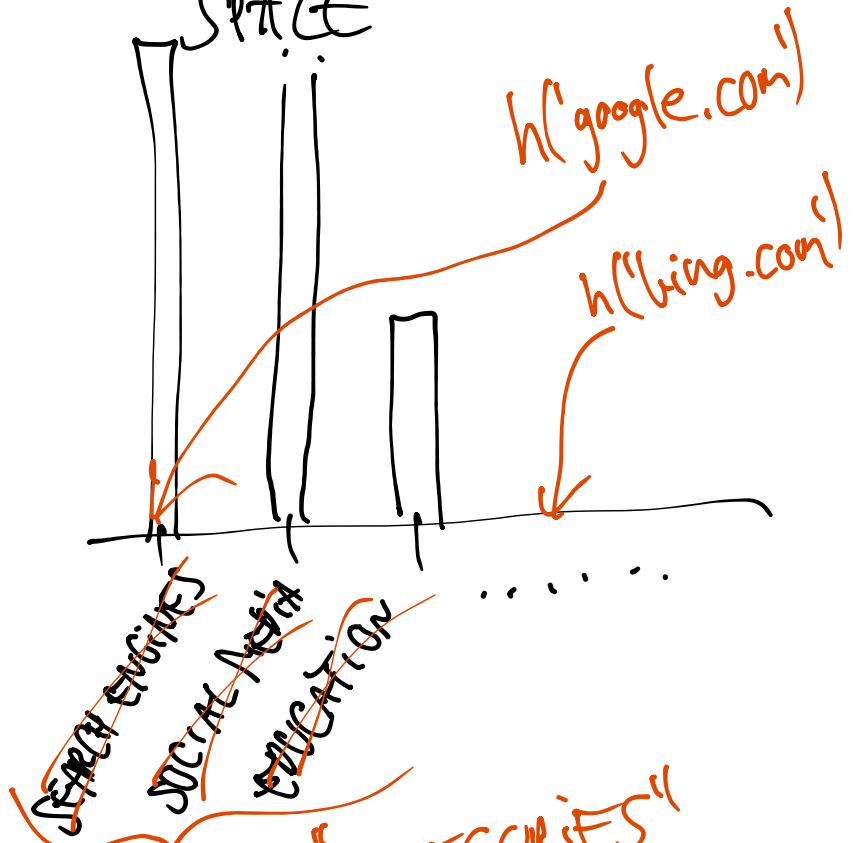
RECORD A NUMBER FOR EACH ITEM IN SOME SET



WHAT KIND OF APPROXIMATION CAN WE GET IN THIS WAY?

SPACE = # DISTINCT ITEMS

CAN WE LUMP THINGS TOGETHER TO REDUCE SPACE

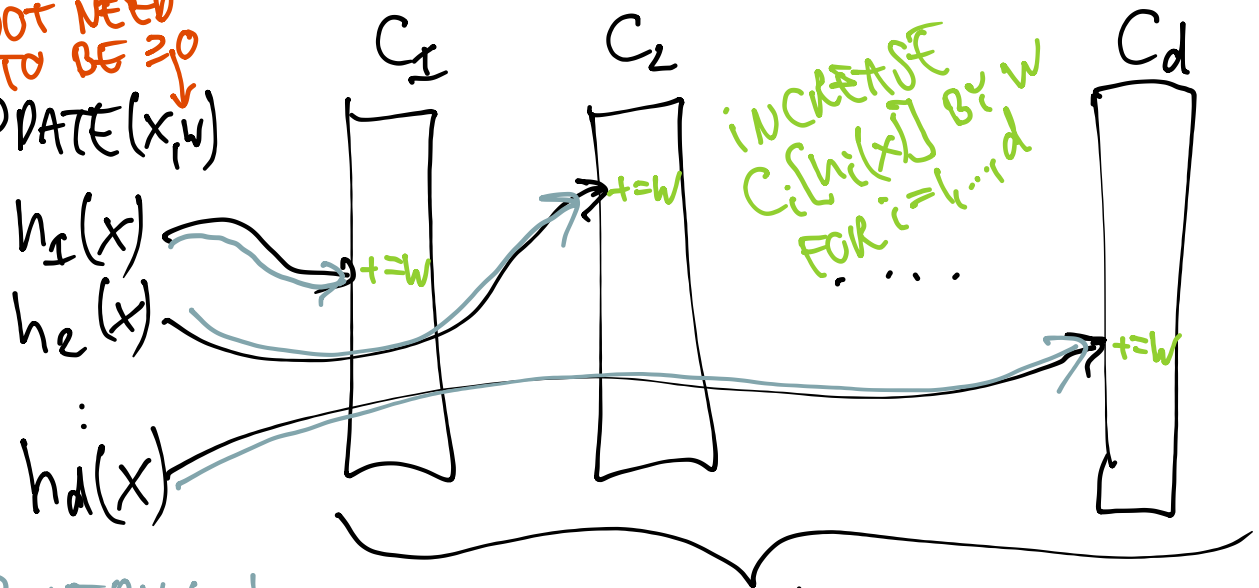


CHOOSING "CATEGORIES" AT RANDOM USING HASHING

COUNT-MIN SKETCH

NB! NOTATION DIFFERENT FROM BOOK

DOES NOT NEED TO BE ≥ 0
UPDATE(x, w)



PAUSE AND THINK:
- HOW CAN TWO COUNT-MIN SKETCHES BE MERGED?
- CAN YOU "REMOVE" ONE SKETCH FROM ANOTHER?

size $t = m/d$ EACH

IF YOU KNOW ALL ELEMENTS OF ONE IS IN THE OTHER.

QUERY(x)

d ARRAYS

RETURN $\min_i C_i[h_i(x)]$

$$= w_x^* + \min_i \left(\sum_{\substack{y \neq x \\ h(y) = h_i(x)}} w_y^* \right)$$

"SMALLEST NOISE"

$$C_i[h_i(x)] = \underbrace{w_x^*}_{\text{TRUE SUM}} + \sum_{\substack{y \neq x \\ h(y) = h_i(x)}} w_y^*$$

"NOISE"

COUNT-MIN ANALYSIS

ASSUME $w_x^* \geq 0$ FOR ALL x , $\|w_x^*\|_1 = \sum_x w_x^*$

NOT NEEDED, BUT SIMPLIFIES

LET $\hat{w}_x = \min_i C_i[h_i(x)]$ BE THE ESTIMATOR OF w_x^* .

WE HAVE $\hat{w}_x = w_x^* + \min_i \sum_{\substack{y \neq x \\ h_i(y)=h_i(x)}} w_y^* \geq w_x^*$ BY ASSUMPTION $w_y^* \geq 0, \forall y$

EXPECTED ERROR IN $C_i[h_i(x)]$:

$$E\left[\sum_{\substack{y \neq x \\ h_i(y)=h_i(x)}} w_y^*\right] = \sum_{y \neq x} \Pr[h(y)=h(x)] \cdot w_y^* = \sum_{y \neq x} \frac{1}{\epsilon} w_y^* \leq \|w^*\|_1 / \epsilon.$$

SIMILAR TO ERROR IN MG WITH $R=\epsilon$

MARKOV'S INEQUALITY: $\Pr\left[\sum_{\substack{y \neq x \\ h_i(y)=h_i(x)}} w_y > \underbrace{2 \cdot \frac{\|w^*\|_1}{\epsilon}}_{\text{TWICE EXPECTATION}}\right] < \frac{1}{2}$

COUNT-MIN ANALYSIS, CONT.

$$\Pr \left[\hat{w}_x > w_x^* + 2 \frac{\|w^*\|_1}{t} \right]$$

$$= \Pr \left[C_i[h_i(x)] > w_x^* + 2 \frac{\|w^*\|_1}{t} \text{ FOR } i = 1, \dots, d \right]$$

indep.

$$\stackrel{\text{indep.}}{\Downarrow} \prod_i \Pr \left[C_i[h_i(x)] > w_x^* + 2 \frac{\|w^*\|_1}{t} \right] < \left(\frac{1}{2}\right)^d$$

← CHOOSE $d = \log_2\left(\frac{1}{\delta}\right)$
TO SUCCEED
WITH PROB. $1 - \delta$

"LINEAR SKETCH"

COMPARISON:	SPACE ϵ	ERROR BOUND	ERROR PROB.	SUBTRACTION ELIMINATES INSERT.
COUNT-MIN	$d \cdot t$	$2\ w^*\ _1/t$	2^{-d}	YES
MISRA-GRIES	k	$\ w^*\ _1/(k+1)$	0	No

↑
DATA STRUCTURE
DETERIORATES WITH
NUMBER OF SUBTRACTIONS

UNBIASED ESTIMATOR FROM COUNT-MIN

- IDEA: EXPECTED "NOISE" IN $C_i[h_i(x)]$ IS THE SAME AS THE EXPECTED VALUE OF $C_i[j]$ FOR $j \neq h_i(x)$

- ESTIMATOR: $E[C_i[h_i(x)] - C_i[h_i(x)+1 \bmod t]] = w_x^*$

- DECREASING ERROR PROBABILITY:

$$\hat{w}_x = \underset{i}{\text{median}} (C_i[h_i(x)] - C_i[h_i(x)+1 \bmod t])$$

PARTICULAR IMPLEMENTATION OF THIS IDEA: COUNT SKETCH (SEE BOOK)