

CarnegieMellon

NOT IN THE EXAM

15-826: Multimedia Databases and Data Mining

Lecture #32: BONUS LECTURE
Approximate Counting
C. Faloutsos

1

CarnegieMellon

Material

- Christopher Palmer, Phillip B. Gibbons and Christos Faloutsos, [*ANF: A Fast and Scalable Tool for Data Mining in Massive Graphs*](#), KDD 2002
- [*Efficient and Tunable Similar Set Retrieval*](#), by Aristides Gionis, Dimitrios Gunopulos and Nikos Koudas, SIGMOD, 2001.
- [*New sampling-based summary statistics for improving approximate query answers*](#), by Phillip B. Gibbons and Yossi Matias, ACM SIGMOD, 1998.

15-826 Copyright (c) 2019 C. Faloutsos 2

2

CarnegieMellon

Outline


- **Flajolet-Martin (and Cohen) – vocabulary size (Problem #1)**
- Application: Approximate Neighborhood function (ANF)
- other, powerful approximate counting tools (Problem #2, #3)

15-826 Copyright (c) 2019 C. Faloutsos 3

3

CarnegieMellon

Problem #1



- Given a multiset (eg., words in a document)
- find the vocabulary size (#, after dup. elimination)

AAABABACAB

Voc. Size = 3 = $|\{A, B, C\}|$


15-826 Copyright (c) 2019 C. Faloutsos 4

4

CarnegieMellon

Thanks to


- Chris Palmer (Vivisimo->IBM)



15-826 Copyright (c) 2019 C. Faloutsos 5

5

CarnegieMellon



Problem #2

- Given a multiset
- compute approximate high-end histogram = hot-list query = (k most common words, and their counts)


AAABABACABDDDDD

(for $k=2$:
A#: 6
D#: 5)

15-826 Copyright (c) 2019 C. Faloutsos 6

6

CarnegieMellon




Problem #3

- Given two documents
- compute quickly their similarity ($\frac{\text{\#common words}}{\text{\#total-words}}$) == Jaccard coefficient

15-826 Copyright (c) 2019 C. Faloutsos 7

7

CarnegieMellon



Problem #1

- Given a multiset (eg., words in a document)
- find the vocabulary size V ($\#$, after dup. elimination)
- using space $O(V)$, or $O(\log(V))$

(Q1: Applications?)
(Q2: How would you solve it?)


15-826 Copyright (c) 2019 C. Faloutsos 8

8

CarnegieMellon

Basic idea (Cohen)

large bit string, initially all zeros



A

A

C

15-826 Copyright (c) 2019 C. Faloutsos 9

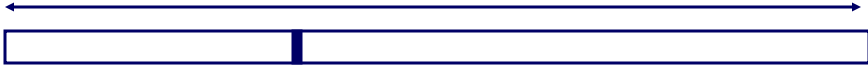
The diagram shows a long horizontal rectangle representing a bit string. Above it is a double-headed arrow with the text "large bit string, initially all zeros". Below the rectangle, the letters "A", "A", and "C" are listed vertically. A small vertical line is positioned at the left end of the rectangle, with a horizontal line extending to the left towards the first "A".

9

CarnegieMellon

Basic idea (Cohen)

large bit string, initially all zeros



A

A

C

15-826 Copyright (c) 2019 C. Faloutsos 10

The diagram is similar to the previous one, but with a thick vertical bar on the bit string rectangle. A horizontal line extends from the first "A" to the left, and a vertical line goes up from the thick bar to meet it, with the word "hash!" written next to the vertical line.

10

CarnegieMellon

Basic idea (Cohen)

large bit string

The diagram shows a long horizontal rectangle representing a bit string. A double-headed arrow above it is labeled "large bit string". A single vertical bar is positioned near the left end of the rectangle. Three horizontal lines with arrows point from the labels "A", "A", and "C" on the left to the vertical bar. The top "A" line is the highest, the middle "A" line is lower, and the "C" line is the lowest.

15-826 Copyright (c) 2019 C. Faloutsos 11

11

CarnegieMellon

Basic idea (Cohen)

large bit string

The diagram shows a long horizontal rectangle representing a bit string. A double-headed arrow above it is labeled "large bit string". Two vertical bars are positioned near the left end of the rectangle. Three horizontal lines with arrows point from the labels "A", "A", and "C" on the left to the two vertical bars. The top "A" line points to the first bar, the middle "A" line points to the second bar, and the "C" line points to the first bar.

15-826 Copyright (c) 2019 C. Faloutsos 12

12

CarnegieMellon

Basic idea (Cohen)

large bit string

A ————— ↑

A ————— ↑

C ————— ↑

the rightmost position depends on the
vocabulary size
(and so does the left-most)

Repeat, with several hashing
functions, and merge the estimates

15-826 Copyright (c) 2019 C. Faloutsos 13

13

CarnegieMellon

Basic idea (Cohen)

large bit string

A ————— ↑

A ————— ↑

C ————— ↑

the rightmost position depends on the
vocabulary size
(and so does the left-most)

Can we do it in less space??

15-826 Copyright (c) 2019 C. Faloutsos 14

14

CarnegieMellon

Basic idea (Cohen)

large bit string

A ———— ↑

A ———— ↑

C ———— |

the rightmost position depends on the
vocabulary size
(and so does the left-most)

Can we do it in less space??
YES

15-826 Copyright (c) 2019 C. Faloutsos 15

15

CarnegieMellon

How?

15-826 Copyright (c) 2019 C. Faloutsos 16

16

CarnegieMellon

Basic idea (Flajolet-Martin)

$O(\log(V))$ bit string (V: voc. size)

first bit: with prob. $\frac{1}{2}$
 second: with prob. $\frac{1}{4}$
 ...
 i-th: with prob. $\frac{1}{2} * i$

15-826 Copyright (c) 2019 C. Faloutsos 17

17

CarnegieMellon

Basic idea (Flajolet-Martin)

$O(\log(V))$ bit string (V: voc. size)

again, the rightmost bit
'reveals' the vocabulary size

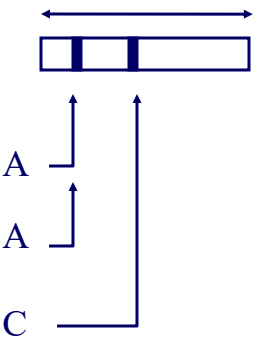
15-826 Copyright (c) 2019 C. Faloutsos 18

18

CarnegieMellon

Basic idea (Flajolet-Martin)

$O(\log(V))$ bit string (V : voc. size)



again, the rightmost bit 'reveals' the vocabulary size

Eg.: $V=4$, will probably set the 2nd bit, etc

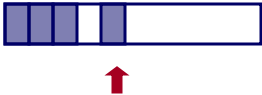
15-826 Copyright (c) 2019 C. Faloutsos 19

19

CarnegieMellon

Flajolet-Martin

- Hash multiple values of X to same signature
 - Hash each x to a bit, using exponential distr.
 - $\frac{1}{2}$ map to bit 0, $\frac{1}{4}$ map to bit 1, ...
- Do several different mappings and average
 - Gives better accuracy
 - Estimate is: $2^b / .77351 / BIAS$
 - $b \sim$ rightmost '1', and actually:



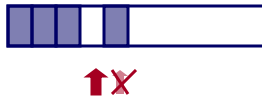
15-826 Copyright (c) 2019 C. Faloutsos 20

20

CarnegieMellon

Flajolet-Martin

- Hash multiple values of X to same signature
 - Hash each x to a bit, using exponential distr.
 - $\frac{1}{2}$ map to bit 0, $\frac{1}{4}$ map to bit 1, ...
- Do several different mappings and average
 - Gives better accuracy
 - Estimate is: $2^b / .77351 / BIAS$
 - b : average least zero bit in the bitmask
 - $bias$: $1 + .31/k$ for k different mappings
- Flajolet & Martin prove this works



15-826 Copyright (c) 2019 C. Faloutsos 21

21

CarnegieMellon

FM Approx. Counting Alg.

```

Assume  $X = \{ 0, 1, \dots, V-1 \}$ 
FOR  $i = 1$  to  $k$  DO  $bitmask[i] = 0000\dots 00$ 
Create  $k$  random hash functions,  $hash_i$ 
FOR each element  $x$  of  $M$  DO
  FOR  $i = 1$  to  $k$  DO
     $h = hash_i(x)$ 
     $bitmask[i] = bitmask[i] \text{ LOR } h$ 
Estimate:  $b = \text{average least zero bit in } bitmask[i]$ 
 $2^b \cdot .77351 / (1 + .31/k)$ 
  
```

- How many bits? $\log V + \text{small constant}$
- What hash functions?

15-826 Copyright (c) 2019 C. Faloutsos 22

22

CarnegieMellon

Random Hash Functions

- Can use linear hash functions. Pick random (a_i, b_i) and then the hash function is:
 - $lhash_i(x) = a_i * x + b_i$
- Gives uniform distribution over the bits
- To make this exponential, define
 - $hash_i(x) =$ least zero bit in $lhash_i(x)$
- Hash functions easy to create and fast to use

15-826 Copyright (c) 2019 C. Faloutsos 23

23

CarnegieMellon

Conclusions

- Want to measure # of distinct elements
- Approach #1: (Flajolet-Martin)
 - Map elements to random bits
 - Keep bitmask of bits
 - Estimate is $O(2^b)$ for least zero-bit b
- Approach #2: (Cohen)
 - Create random permutation of elements
 - Keep least element seen
 - Estimate is: $O(1/le)$ for least rank le

15-826 Copyright (c) 2019 C. Faloutsos 24

24

CarnegieMellon

Approximate counting

- Flajolet-Martin (and Cohen) – vocabulary size
- **Application: Approximate Neighborhood function (ANF)**
- other, powerful approximate counting tools

15-826 Copyright (c) 2019 C. Faloutsos 25

25

CarnegieMellon

Fast Approximation of the “neighborhood” Function for Massive Graphs details

Christopher R. Palmer
Phillip B. Gibbons
Christos Faloutsos

KDD 2001

26

CarnegieMellon

details

Motivation

- What is the diameter of the Web?
- What is the effective diameter of the Web?
- Are the telephone caller-callee graphs for the U.S. similar to the ones in Europe?
- Is the citation graph for physics different from the one for computer science?
- Are users in India further away from the core of the Internet than those in the U.S.?

15-826 Copyright (c) 2019 C. Faloutsos 27

27

CarnegieMellon

details

Proposed Tool: neighborhood

Given graph $G=(V,E)$
 $N(h)$ = # pairs within h hops or less
= **neighborhood function**

15-826 Copyright (c) 2019 C. Faloutsos 28

28

CarnegieMellon

details

Proposed Tool: neighborhood

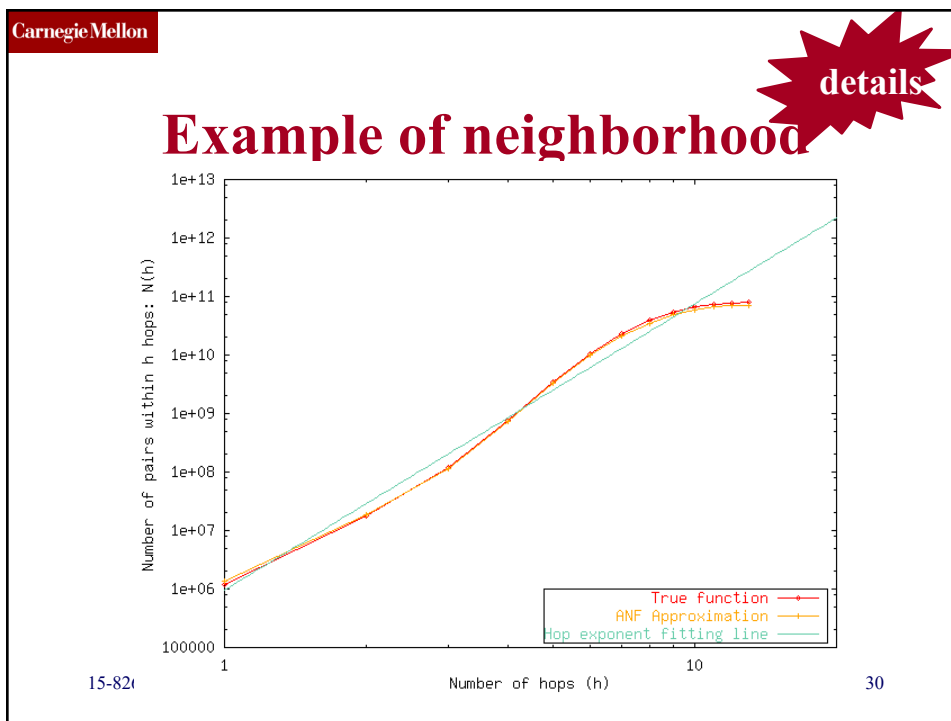
Given graph $G=(V,E)$

$N(h)$ = # pairs within h hops or less
= **neighborhood function**

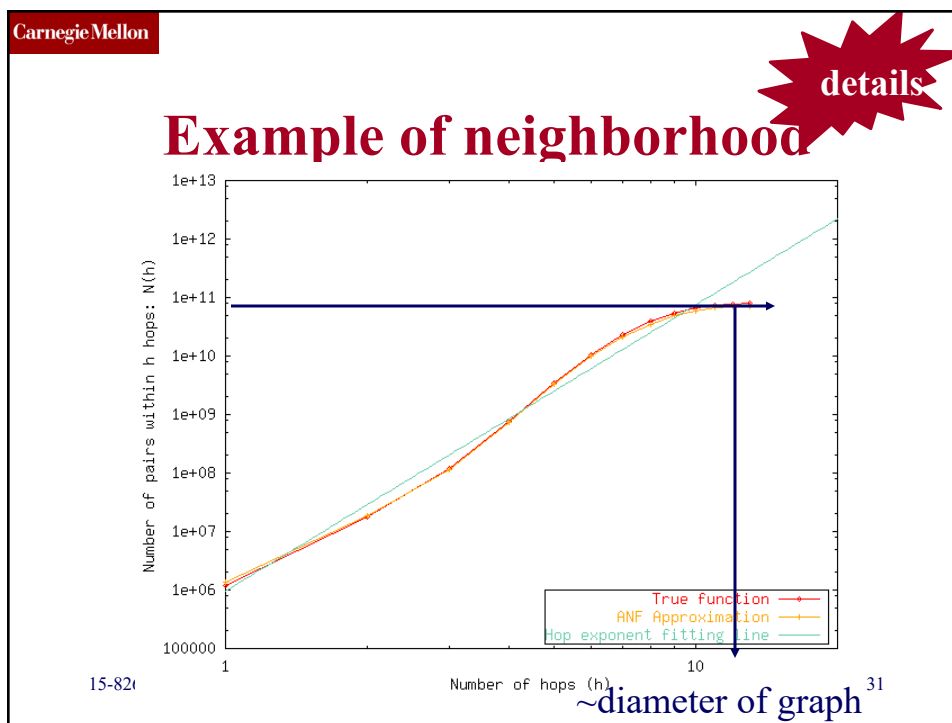
$N(u,h)$ = # neighbors of node u , within h hops or less

15-826 Copyright (c) 2019 C. Faloutsos 29

29



30



31

CarnegieMellon

Requirements (for massive graphs) details

- *Error guarantees*
- *Fast:* (and must scale linearly with graph)
- *Low storage requirements:* massive graphs!
- *Adapts to available memory*
- *Sequential scans of the edges*
- Also estimates *individual neighborhood functions* $|S(u,h)|$
 - These are actually quite useful for mining

15-826

Copyright (c) 2019 C. Faloutsos

32

32

CarnegieMellon

details

How would you compute it!

- Repeated matrix multiply
 - Too slow $O(n^{2.38})$ at the very least
 - Too much memory $O(n^2)$
- Breadth-first search
 - FOR each node u DO
 - bf-search to compute $S(u,h)$ for each h
 - Best known exact solution!
 - We will use this as a reference
- Approximations? Only 1 that we know of which we will discuss when we evaluate it.

15-826 Copyright (c) 2019 C. Faloutsos 33

33

CarnegieMellon

details

Intuition

- Guess what we'll use?
 - Approximate Counting!
- Use very simple algorithm:
 - FOR each node u DO $S(u,0) = \{ (u,u) \}$ initialize to self-only
 - FOR $h = 1$ to *diameter* of G DO
 - FOR each node u DO $S(u,h) = S(u,h-1)$ can reach same things
 - FOR each edge (u,v) in G DO and add one more step
 - $S(u,h) = S(u,h) \cup \{ (u,v') : (v,v') \in S(v,h-1) \}$

15-826 Copyright (c) 2019 C. Faloutsos 34

34

CarnegieMellon
Intuition

details

- Guess what we'll use?
 - Approximate Counting!
- Use very simple algorithm:

(distinct) neighbors of u,
within h hops

FOR each node u DO $S(u,0) = \{ (u,u) \}$ initialize to self-only

FOR $h = 1$ to diameter of G DO

FOR each node u DO $S(u,h) = S(u,h-1)$ can reach same things

FOR each edge (u,v) in G DO and add one more step

$S(u,h) = S(u,h) \cup \{ (u,v') : (v,v') \in S(v,h-1) \}$

(distinct) neighbors of v,
within h-1 hops

15-826
Copyright (c) 2019 C. Faloutsos
35

35

CarnegieMellon
Trace

details

$h=0$

$\{(1,1)\}$

$\{(2,2)\}$

$\{(3,3)\}$

$\{(4,4)\}$

```

graph TD
    1 --- 2
    1 --- 3
    2 --- 3
    3 --- 4
    
```

15-826
Copyright (c) 2019 C. Faloutsos
36

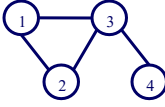
36

CarnegieMellon

details

Trace

h=0	h=1
$\{(1,1)\}$	$\{(1,1)\}$
$\{(2,2)\}$	$\{(2,2)\}$
$\{(3,3)\}$	$\{(3,3)\}$
$\{(4,4)\}$	$\{(4,4)\}$



15-826
Copyright (c) 2019 C. Faloutsos
37

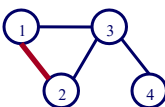
37

CarnegieMellon

details

Trace

h=0	h=1
$\{(1,1)\}$	$\{(1,1)\}$
$\{(2,2)\}$	$\{(2,2)\}$
$\{(3,3)\}$	$\{(3,3)\}$
$\{(4,4)\}$	$\{(4,4)\}$



15-826
Copyright (c) 2019 C. Faloutsos
38

38

CarnegieMellon

details

Trace

h=0	h=1
$\{(1,1)\}$	$\{(1,1), (1,2)\}$
$\{(2,2)\}$	$\{(2,2)\}$
$\{(3,3)\}$	$\{(3,3)\}$
$\{(4,4)\}$	$\{(4,4)\}$

```

graph TD
    1((1)) --- 2((2))
    1 --- 3((3))
    3 --- 4((4))
    style 1 stroke:#f00,stroke-width:2px
    style 2 stroke:#f00,stroke-width:2px
  
```

15-826
Copyright (c) 2019 C. Faloutsos
39

39

CarnegieMellon

details

Trace

h=0	h=1
$\{(1,1)\}$	$\{(1,1), (1,2), (1,3)\}$
$\{(2,2)\}$	$\{(2,2)\}$
$\{(3,3)\}$	$\{(3,3)\}$
$\{(4,4)\}$	$\{(4,4)\}$


```

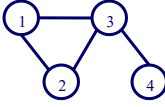
graph TD
    1((1)) --- 2((2))
    1 --- 3((3))
    3 --- 4((4))
    style 1 stroke:#f00,stroke-width:2px
    style 3 stroke:#f00,stroke-width:2px
  
```

15-826
Copyright (c) 2019 C. Faloutsos
40

40

CarnegieMellon
Trace




h=0	h=1	
$\{(1,1)\}$	$\{(1,1), (1,2), (1,3)\}$	
$\{(2,2)\}$	$\{(2,2), (2,1), (2,3)\}$	
$\{(3,3)\}$	$\{(3,3), (3,1), (3,2), (3,4)\}$	
$\{(4,4)\}$	$\{(4,4), (4,3)\}$	

15-826
Copyright (c) 2019 C. Faloutsos
41

41

CarnegieMellon
Intuition



- Guess what we'll use?
 - Approximate Counting!
- Use very simple algorithm:
 - FOR each node u DO $S(u,0) = \{ (u,u) \}$ initialize to self-only
 - FOR $h = 1$ to *diameter* of G DO can reach same things
 - FOR each node u DO $S(u,h) = S(u,h-1)$ and add one more step
 - FOR each edge (u,v) in G DO
 - $S(u,h) = S(u,h) \cup \{ (u,v') : (v,v') \in S(v,h-1) \}$

(distinct) neighbors of u , within h hops

15-826
Copyright (c) 2019 C. Faloutsos
42

42

CarnegieMellon

details

Intuition

- Guess what we'll use?
 - Approximate Counting!
- Use very simple algorithm:
 - FOR each node u DO $S(u,0) = \{u,u\}$ initialize to self-only
 - FOR $h = 1$ to *diameter* of G DO can reach same things
 - FOR each node u DO $S(u,h) = S(u,h-1)$ and add one more step
 - FOR each edge (u,v) in G DO
 - $S(u,h) = S(u,h) \cup \{(u,v') : (v,v') \in S(v,h-1)\}$
- Too slow and requires too much memory
- Replace expensive set ops with bit ops

(distinct) neighbors of u , within h hops

15-826 Copyright (c) 2019 C. Faloutsos 43

43

CarnegieMellon

details

ANF Algorithm #1

FOR each node, u , DO

$M(u,0) =$ concatenation of k bitmasks of length $\log n + r$
each bitmask has 1 bit set (exp. distribution)

DONE

FOR $h = 1$ to *diameter* of G DO

FOR each node, u , DO $M(u,h) = M(u,h-1)$

FOR each edge (u,v) in G DO

$M(u,h) = (M(u,h) \text{ OR } M(v,h-1))$

Estimate $N(h) = \text{Sum}(N(u,h)) = \text{Sum } 2^{b(u)} / .77351 / (1+.31/k)$
where $b(u) =$ average least zero bit in $M(u,ih)$

DONE

15-826 Copyright (c) 2019 C. Faloutsos 44

44

CarnegieMellon

details

ANF Algorithm #1

FOR each node, u , DO
 $M(u,0)$ = concatenation of k bitmasks of length $\log n + r$
 each bitmask has 1 bit set (exp. distribution)
 DONE

FOR $h = 1$ to *diameter* of G DO
 FOR each node, u , DO $M(u,h) = M(u,h-1)$
 FOR each edge (u,v) in G DO
 $M(u,h) = (M(u,h) \text{ OR } M(v,h-1))$

Estimate $N(h) = \text{Sum}(N(u,h)) = \text{Sum } 2^{b(u)} / .77351 / (1+.31/k)$
 where $b(u)$ = average least zero bit in $M(u, it)$

DONE
 15-826

Copyright (c) 2019 C. Faloutsos

45

45

CarnegieMellon

details

ANF Algorithm #1

whatever u can reach with h hops
 plus whatever v can reach with $h-1$ hops
 Duplicates: automatically eliminated!

$M(u,h) = (M(u,h) \text{ OR } M(v,h-1))$


15-826

Copyright (c) 2019 C. Faloutsos

46

46

CarnegieMellon



Properties

- **Has error guarantees:** (from F&M)
- **Is fast:** $O((n+m)d)$ for n nodes, m edges, diameter d (which is typically small)
- **Has low storage requirements:** $O(n)$
- **Easily parallelizable:** Partition nodes among processors, communicate after full iteration
- **Does sequential scans of edges.**
- **Estimates individual neighborhood functions**
- **DOES NOT work with limited memory**

15-826 Copyright (c) 2019 C. Faloutsos 47

47

CarnegieMellon

Conclusions

- Approximate counting (ANF / Martin-Flajolet) take minutes, instead of hours
- and discover interesting facts quickly

15-826 Copyright (c) 2019 C. Faloutsos 48

48

CarnegieMellon

Outline

- Flajolet-Martin (and Cohen) – vocabulary size (Problem #1)
- Application: Approximate Neighborhood function (ANF)
- other, powerful approximate counting tools


➔ **(Problem #2, #3)**

15-826 Copyright (c) 2019 C. Faloutsos 49

49

CarnegieMellon

Problem #2



- Given a multiset
- compute approximate high-end histogram = hot-list query = (k most common words, and their counts)

AAABABACABDDDDD

(for $k=2$:
 A#: 6
 D#: 5)

15-826 Copyright (c) 2019 C. Faloutsos 50

50

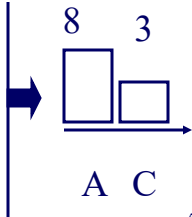
CarnegieMellon

Hot-list queries

- Given a stream of product ids (with duplicates)
- Compute
 - the k most frequent products,
 - and their counts
- with a SINGLE PASS and $O(k)$ memory

A A B A C A B C A A D E A C A

$k=2$



A C

15-826 Copyright (c) 2019 C. Faloutsos 51

51

CarnegieMellon

Applications?

15-826 Copyright (c) 2019 C. Faloutsos 52

52

CarnegieMellon

Applications?

- Best selling products
- most common words
- most busy IP destinations/sources (DoS attacks)
- summarization / synopses of datasets
- high-end histograms for DBMS query optimization

15-826

Copyright (c) 2019 C. Faloutsos

53

53

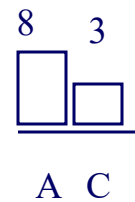
CarnegieMellon

Hot-list queries

- Given a stream of product ids (with duplicates)
- Compute
 - the k most frequent products,
 - and their counts
- with a SINGLE PASS and $O(k)$ memory

A A B A C A B C A A D E A C A

Exact: impossible

Thus: **approximate** $k=2$ 

15-826

Copyright (c) 2019 C. Faloutsos

54

54

CarnegieMellon

Hot-list queries - idea


- Keep the (approx.) k best so far, plus counts
- for a new item, if it is in the hot list
 - increment its count

A A B A C A B C A A D E A C A

↑

$k=2$

2 1



A B

15-826
Copyright (c) 2019 C. Faloutsos
55

55

CarnegieMellon

Hot-list queries - idea

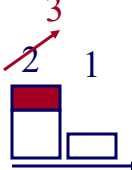
- Keep the (approx.) k best so far, plus counts
- for a new item, if it is in the hot list
 - increment its count

A A B A C A B C A A D E A C A

↑

$k=2$

3 1



A B

15-826
Copyright (c) 2019 C. Faloutsos
56

56

CarnegieMellon

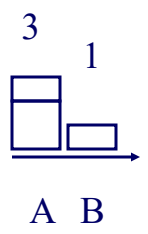
Hot-list queries - idea

- Keep the (approx.) k best so far, plus counts
- for a new item, if it is in the hot list
 - increment its count
 - else ??

A A B A C A B C A A D E A C A

↑

$k=2$



15-826 Copyright (c) 2019 C. Faloutsos 57

57

CarnegieMellon

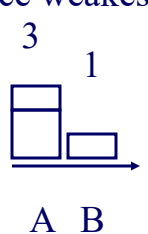
Hot-list queries - idea

- Keep the (approx.) k best so far, plus counts
- for a new item, if it is in the hot list
 - increment its count
 - else TOSS a coin, and possibly displace weakest

A A B A C A B C A A D E A C A

↑

$k=2$



15-826 Copyright (c) 2019 C. Faloutsos 58

58

CarnegieMellon

Hot-list queries - idea

- Biased coin - what are the Head/Tail prob.?

A A B A C A B C A A D E A C A

↑

k=2

6

A B

15-826 Copyright (c) 2019 C. Faloutsos 59

59

CarnegieMellon

Hot-list queries - idea

- Biased coin - what are the Head/Tail prob.?
- A: depends on count(weakest)

A A B A C A B C A A D E A C A

↑

k=2

6

A B

15-826 Copyright (c) 2019 C. Faloutsos 60

60

CarnegieMellon

Hot-list queries - idea

- Biased coin - what are the Head/Tail prob.?
- A: depends on count(weakest)
- and the new item ('D'), if it wins, it gets the **count of the item it displaced.**

15-826 Copyright (c) 2019 C. Faloutsos 61

61

CarnegieMellon

Hot-list queries - idea

- See [Gibbons+Matias 98] for proofs

15-826 Copyright (c) 2019 C. Faloutsos 62

62

CarnegieMellon

Outline


- Flajolet-Martin (and Cohen) – vocabulary size (Problem #1)
- Application: Approximate Neighborhood function (ANF)
- other, powerful approximate counting tools
 - Problem #2,
 - **Problem #3**

15-826 Copyright (c) 2019 C. Faloutsos 63

63

CarnegieMellon

Problem #3




- Given two documents
- compute quickly their similarity ($\frac{\text{\#common words}}{\text{\#total-words}}$) == Jaccard coefficient

15-826 Copyright (c) 2019 C. Faloutsos 64

64

CarnegieMellon



Problem #3'

- Given a query document q
- and many other documents
- compute quickly the k nearest neighbors of q , using the Jaccard coefficient

D1: {A, B, C} q: {A, C, D, W}
D2: {A, D, F, G}

...

15-826 Copyright (c) 2019 C. Faloutsos 65

65

CarnegieMellon

Applications?

15-826 Copyright (c) 2019 C. Faloutsos 66

66

Applications?

- Set comparisons eg.,
 - snail-mail address (set of trigrams)
- search engines - ‘similar pages’
- social networks: people with many joint friends (facebook recommendations)

Problem #3’

- Given a query document q
- and many other documents
- compute quickly the k nearest neighbors of q , using the Jaccard coefficient

- Q: how to extract a fixed set of numerical features, to index on?

CarnegieMellon

Answer

- Approximation / hashing - Cohen:

15-826 Copyright (c) 2019 C. Faloutsos 69

69

CarnegieMellon

Basic idea (Cohen)

large bit string

the
the
cat

For each document and for a given h.f. return the position of first '1'

Repeat for k h.f. -> each document becomes k numbers

15-826 Copyright (c) 2019 C. Faloutsos 70

70

CarnegieMellon

Idea

- Doc1: n_1, n_2, \dots, n_k
- Doc2: n_1', n_2', \dots, n_k'

15-826 Copyright (c) 2019 C. Faloutsos 71

71

CarnegieMellon

Idea

- Doc1: n_1, n_2, \dots, n_k
- Doc2: n_1', n_2', \dots, n_k'

1 m

- say they agree on m values

15-826 Copyright (c) 2019 C. Faloutsos 72

72

CarnegieMellon

Idea

- Doc1: n_1, n_2, \dots, n_k
- Doc2: n_1', n_2', \dots, n_k'

- say they agree on m values,
- then


Jaccard(Doc1, Doc2) $\sim m/k$

15-826 Copyright (c) 2019 C. Faloutsos 73

73

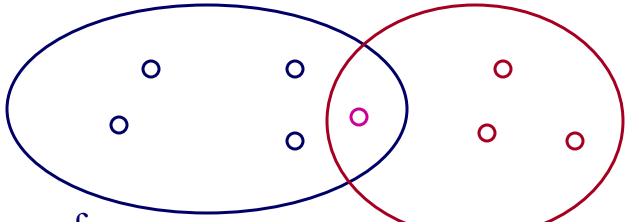
CarnegieMellon

Intuition behind proof



Andrew Tomkins

- Venn diagram



voc. terms of
Doc.#1

15-826

voc. terms of
Doc.#2

Copyright (c) 2019 C. Faloutsos

74

74

CarnegieMellon

Intuition behind proof

- Venn diagram

voc. terms of Doc.#1

voc. terms of Doc.#2

15-826 Copyright (c) 2019 C. Faloutsos 75

75

CarnegieMellon

Intuition behind proof

- Venn diagram - let w be the voc. word with the overall smallest hash value, for h.f.#1

voc. terms of Doc.#1

voc. terms of Doc.#2

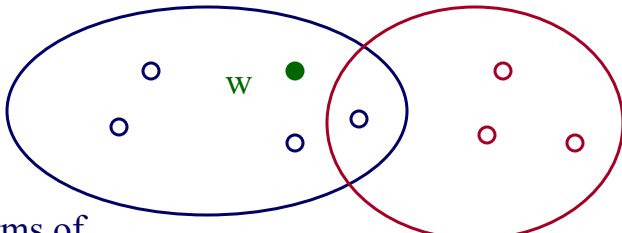
15-826 Copyright (c) 2019 C. Faloutsos 76

76

CarnegieMellon

Intuition behind proof

- Prob. that w is smallest on both is exactly Jaccard: $\#common / \#union$



voc. terms of Doc.#1

voc. terms of Doc.#2

15-826 Copyright (c) 2019 C. Faloutsos 77

77

CarnegieMellon

Conclusions

- Approximations can achieve the impossible!
- MF and ANF for neighborhood function
- hot-lists
- Jaccard coeff. / 'similar pages'

15-826 Copyright (c) 2019 C. Faloutsos 78

78

References

E. Cohen. Size-estimation framework with applications to transitive closure and reachability. *Journal of Computer and System Sciences*, 55(3):441-453, December 1997.

<http://www.research.att.com/~edith/Papers/tcest.ps.Z>

Phillip B. Gibbons, Yossi Matias, *New sampling-based summary statistics for improving approximate query answers*, ACM SIGMOD, 1998 Seattle, Washington, pp 331 - 342

References (cont' d)

Aristides Gionis, Dimitrios Gunopulos, Nikos Koudas, *Efficient and Tunable Similar Set Retrieval*, ACM SIGMOD 2001, Santa Barbara, California

M. Faloutsos, P. Faloutsos, and C. Faloutsos. *On power-law relationships for the internet topology*. SIGCOMM, 1999.

References (cont' d)

- P. Flajolet and G. N. Martin. Probabilistic counting algorithms for data base applications. *Journal of Computer and System Sciences*, 31:182-209, 1985.
- C. R. Palmer, P. B. Gibbons and C. Faloutsos. *Fast approximation of the "neighborhood" function for massive graphs*. KDD 2002

References (cont' d)

- C. R. Palmer, G. Siganos, M. Faloutsos, P. B. Gibbons and C. Faloutsos. *The connectivity and fault-tolerance of the internet topology*. NRDM 2001.