An Introduction to ANALYTIC COMBINATORICS

Computer Science 488 Robert Sedgewick Jérémie Lumbroso

Course overview

Analysis of algorithms

- Methods and models for the analysis of algorithms.
- Basis for a scientific approach.
- Mathematical methods from classical analysis.
- Combinatorial structures and associated algorithms.

Analytic combinatorics

- Study of properties of large combinatorial structures.
- A foundation for analysis of algorithms, but widely applicable.
- Symbolic method for encapsulating precise description.
- Complex analysis to extract useful information.





Is this course for me?

Sure, if you can answer "yes" to these questions.

- Do you like to program ?
- Do you like math ?
- Have you read Algorithms?
- Would you like to be able to read Knuth's books?





• and Flajolet's papers?

Q. Why study the analysis of algorithms and analytic combinatorics?

A. For many of the same reasons we study *algorithms* (next) !



3

Their impact is broad and far-reaching.

Internet. Web search, packet routing, file sharing, ... Biology. Human genome project, protein folding, ... Computer design. Circuit layout, file system, compilers, ... Multimedia. Movies, video games, virtual reality, ... Security. Cell phones, e-commerce, voting machines, ... Social networks. Recommendations, news feeds, advertisements, ... Physics. N-body simulation, particle collision simulation, ... Big data. Deep learning, autonomous vehicles, ...







Old roots, new opportunities.

- Analysis of algorithms dates at least to Euclid.
- Practiced by Turing and von Neumann in 1940s.
- Mostly developed by Knuth starting in 1960s.
- Steady evolution for decades.
- Analytic combinatorics dates to Euler and earlier.
- Mostly developed by Flajolet starting in 1980s.
- Many algorithms are waiting to be understood.
- Many theorems are waiting to be discovered.

"father of analysis of algorithms"



Don Knuth

"father of analytic combinatorics"



Philippe Flajolet



"If I have seen further, it is by standing on the shoulders of giants."

– Isaac Newton

To solve problems that could not otherwise be addressed.

Example: Cardinality estimation (stay tuned).

pool-71-104-94-246.lsanca.dsl-w.verizon.net 117.222.48.163 pool-71-104-94-246.lsanca.dsl-w.verizon.net 1.23.193.58 188.134.45.71 1.23.193.58 gsearch.CS.Princeton.EDU pool-71-104-94-246.lsanca.dsl-w.verizon.net 81.95.186.98.freenet.com.ua 81.95.186.98.freenet.com.ua 81.95.186.98.freenet.com.ua CPE-121-218-151-176.lnse3.cht.bigpond.net.au 117.211.88.36 msnbot-131-253-46-251.search.msn.com msnbot-131-253-46-251.search.msn.com pool-71-104-94-246.lsanca.dsl-w.verizon.net gsearch.CS.Princeton.EDU CPE001cdfbc55ac-CM0011ae926e6c.cpe.net.cable.rogers.com CPE001cdfbc55ac-CM0011ae926e6c.cpe.net.cable.rogers.com 118-171-27-8.dynamic.hinet.net cpe-76-170-182-222.socal.res.rr.com

— How many of these are different?

6

For intellectual stimulation.



"The point of mathematics is that in it we have always got rid of the particular instance, ... no mathematical truths apply merely to fish, or merely to stones, or merely to colours. So long as you are dealing with pure mathematics, you are in the realm of complete and absolute abstraction. ... Mathematics is thought moving in the sphere of complete abstraction from any particular instance of what it is talking about.

- Alfred North Whitehead



Abstract Thought 379, by Theo Dapore

"Here's to pure mathematics—may it never be of any use to anybody."

– attributed to G. H. Hardy



"PEOPLE WHO ANALYZE ALGORITHMS have double happiness. First of all they experience the sheer beauty of elegant mathematical patterns that surround elegant computational procedures. Then they receive a practical payoff when their theories make it possible to get other jobs done more quickly and more economically."



– Don Knuth

They may unlock the secrets of life and of the universe.



"Pure mathematics is, in its way, the poetry of logical ideas. One seeks the most general ideas of operation which will bring together in simple, logical and unified form the largest possible circle of formal relationships. In this effort toward logical beauty spiritual formulas are discovered necessary for the deeper penetration into the laws of nature."



– Albert Einstein



Some compelling reasons

- Their impact is broad and far-reaching.
- Old roots, new opportunities.
- To solve problems that could not otherwise be addressed.
- For intellectual stimulation.
- They may unlock the secrets of life and of the universe.
- For fun and profit.



Today. A case in point.

AofA/AC context



A Case Study: Cardinality Estimation

• Exact cardinality count

- Probabilistic counting
- Stochastic averaging
- Refinements

with special thanks to Jérémie Lumbroso

Cardinality counting

Q. In a given stream of data values, how many different values are present?

Reference application. How many unique visitors in a web log?

log.07.f3.txt

109.108.229.102 pool-71-104-94-246.lsanca.dsl-w.verizon.net 117.222.48.163 pool-71-104-94-246.lsanca.dsl-w.verizon.net 1.23.193.58 188.134.45.71 1.23.193.58 gsearch.CS.Princeton.EDU pool-71-104-94-246.lsanca.dsl-w.verizon.net 81.95.186.98.freenet.com.ua 81.95.186.98.freenet.com.ua 81.95.186.98.freenet.com.ua CPE-121-218-151-176.lnse3.cht.bigpond.net.au

6 million strings



A standard "Interview Question".

13

Wrong answer: Check every value

Check every value

- Save all the values in an array.
- Check all previous values for duplicates.
- Count a value only if no previous duplicate.

```
int[] a = StdIn.readAllLines();
int count = 1;
for (int i = 1; i < a.length; i++)
{
    for (int j = 0; j <= i)
        if (a[j] == a[i]) break;
        if (j != i) count++;
}
StdOut.print(count + " different values");
```



Q. Why is this the wrong answer?

A. QUADRATIC running time, therefore not feasible for real-world applications.

Standard answer I: Sort, then count

Sort, then count

- Save all the values in an array.
- Sort the array.
- Equal values are together in the sorted input.
- Count the first occurrence of each value.



Standard answer I: Sort, then count

Sort, then count

- Save all the values in an array.
- Sort the array.
- Equal values are together in the sorted input.
- Count the first occurrence of each value.

```
int[] a = StdIn.readAllLines();
Arrays.sort(a);
int distinct = 1;
for (int i = 1; i < a.length; i++)
    if (a[i] != a[i-1]) distinct++;
StdOut.print(distinct + " different values");
```

Used by programmers "in the wild" for decades





Programming Exam 1 COS 126 2015

Part 1. Yo on standar	ar task is to write programs that find the number of distinct values among the integers d input, assuming that the input is nonempty <i>and in sorteil order</i> .
Your task.	Add code to this template (the file count 1. java that you have downloaded):
pu	blic class Countl
(<pre>gublic static void main(string[] args) { int count = 1; int distinct = 1;</pre>
	// YOUR CODE HERE
F	3
Your code values amo	must print the number of integers on standard input <i>and</i> the number of distinct on those integers.
Example. six distinc	Test your program with the file testcountitiny.txt, which has 18 integers having values (1, 2, 4, 5, 6, and 9). Your program must behave as follows:
	nore testCountiting.txt
	iava Countl < testCountltiny.txt

Aside: Existence table

IF the values are positive integers less than U:

Use an *existence table*

- Create an array b[] of boolean values.
- For value i, set b[i] to true.
- Count the number of true values in b[].

smal	l exa	mple	2																
15	9	9	4	10	9	11	12	10	14	12	11	15	6	11	9	8	5	10	2
exist	ence	tabl	e (U	= 16) 4	5	6	7	8	q	10	11	12	13	14	15				
Ū	T	T	5	Т	Т	T	,	T	T	T	T	T	15	T	T				

Aside: Existence table

```
StdOut.print(count + " different values");
```

X Not applicable to reference application (long strings).



Programming Exam 1 COS 126 2015 Part 2

Pro	gramming Exam: Count Distinct Values
Part 1. Your t on standard in	isk is to write programs that find the number of distinct values among the integers put, assuming that the input is nonempty <i>and in sorted order</i> .
Your task. Ad	d code to this template (the file count1. java that you have downloaded):
publi	c class Countl
(P	ablic static void main(String[] args)
	<pre>int count = 1; int distinct = 1;</pre>
	// YOUR CODE HERE
}	
Your code mu values among	ast print the number of integers on standard input and the number of distinct those integers.
Example. Test six distinct val	your program with the file testcountitiny.txt, which has 18 integers having ass (1, 2, 4, 5, 6, and 9). Your program must behave as follows:
t mor	e testCountitiny.txt 1 2 2 2 2 4 4 4 5 5 6 6 9 9
8 jav	a Countl < testCountiting.txt

Standard answer II: Use a hash table

Hashing with separate chaining

- Create a table of size *M*.
- Transform each value into a "random" table index.
- Make linked lists for colliding values.
- Ignore values already in the table.

example: multiply by a prime, then take remainder after dividing by *M*.



Exact cardinality count using a hash table

Hashing with separate chaining

- Create a table of size *M*.
- Transform each value into a "random" table index.
- Make linked lists for colliding values.
- Ignore values already in the table.

Widely used and well studied textbook method.



```
Exact cardinality count in Java
```

- Input is an "iterable"
- HashSet implements a hash table
- add() adds new value (noop if already there)
- size() gives number of distinct values added

```
public static long count(Iterable<String> stream)
{
    HashSet<String> hset = new HashSet<String>();
    for (String x : stream)
        hset.add(x);
    return hset.size();
}
```

Mathematical analysis of exact cardinality count with hashing

Theorem. If the hash function uniformly and independently distributes the keys in the table, the expected time and space cost is LINEAR.



Q. Do the hash functions that we use *uniformly and independently distribute keys in the table*?

A. Not likely.

*

Scientific validation of exact cardinality count with hashing

Hypothesis. Time and space cost is linear for the hash functions we use and the data we have.

Quick experiment. Doubling the problem size should double the running time.

Driver to read N strings and count distinct values

```
public static void main(String[] args)
{
    int N = Integer.parseInt(args[0]);
    StringStream stream = new StringStream(N);
    long start = System.currentTimeMillis();
    StdOut.println(count(stream));
    long now = System.currentTimeMillis();
    double time = (now - start) / 1000.0;
    StdOut.println(time + " seconds");
}
```





Q. Is hashing effective?

A. Yes. Validated in countless applications for *over half a century*.

22

Summary of cardinality count algorithms

	time bound	memory bound
Wrong answer	N ²	Ν
Sort and count	N log N	Ν
Existence table	Ν	U
Hash table	N *	Ν



* Theoretical AofA. Hashing solution is *quadratic* in the worst case.

Theoretical AofA. If (uniform hashing assumption) then *hashing solution is linear* (expected).

Scientific AofA. Hypothesis that *hashing solution is linear* has been validated for decades.

Q. End of story?

A. Not at all !!

Cardinality Estimation

- Warmup: exact cardinality count
- Probabilistic counting
- Stochastic averaging
- Refinements

Problem: exact cardinality count requires linear space

Q. I can't use a hash table. The stream is much too big to fit all values in memory. Now what?

A. Bad news: You cannot get an exact count.

A. Good news: You *can* get an accurate estimate (stay tuned).

109.108.229.102 pool-71-104-94-246.lsanca.dsl-w.verizon.net 117.222.48.163 pool-71-104-94-246.lsanca.dsl-w.verizon.net 1.23.193.58 188.134.45.71 1.23.193.58 gsearch.CS.Princeton.EDU pool-71-104-94-246.lsanca.dsl-w.verizon.net 81.95.186.98.freenet.com.ua 81.95.186.98.freenet.com.ua 81.95.186.98.freenet.com.ua CPE-121-218-151-176.lnse3.cht.bigpond.net.au



Typical modern applications

- Can only afford to process each value *once*.
- *Estimate* of count still very useful.
- and cannot spend much time on any value

Cardinality estimation

is a fundamental problem with many applications *where memory is limited*.

Q. *About* how many different values appear in a given stream?

Constraints

- Make one pass through the stream.
- Use as few operations per value as possible
- Use *as little memory* as possible.
- Produce as accurate an estimate as possible.





To fix ideas on scope: Think of *millions* of streams each having *trillions* of values.

Probabilistic counting with stochastic averaging (PCSA)

Flajolet and Martin, Probabilistic Counting Algorithms for Data Base Applications FOCS 1983, JCSS 1985.



Philippe Flajolet 1948-2011

Contributions

- Introduced problem
- Idea of streaming algorithm
- Idea of "small" sketch of "big" data
- Detailed analysis that yields tight bounds on accuracy
- Full validation of mathematical results with experimentation
- Practical algorithm that has remained effective for decades



Bottom line: Quintessential example of the effectiveness of scientific approach to algorithm design.



Probabilistic counting starting point: two integer functions

Definition. r(x) is the **number of** trailing 1s in the binary representation of x. \leftarrow position of rightmost 0

(x)	15	14	12	12	11	10	0	Q	7	6	\int_{Σ}) ⊿	2	2	1		0	$r(\mathbf{v})$	R(v)) R (v)a	
	IJ	14	15	12		10	9	0	1	0	5	4	2	2			0	<i>(</i> (<u></u>)			~]2	
	1	0	1	1	1	1	0	1	1	1	1	1	0	1	. C)	1	1	2		10	
	1	0	1	0	1	0	1	0	1	0	0	0	1	1	. 1	_	0	0	1		1	
	0	1	1	0	1	0	0	1	0	1	0	1	1	1	. 1	-	1	5	32	100	000	
																_						
		0) 1	1	0	1	0	0	1	0	1	0	1	1	1	1	1		x			
		1	. 0	0	1	0	1	1	0	1	0	1	0	0	0	0	0	~	X		3 instri	uctions
oute.		0) 1	1	0	1	0	0	1	0	1	1	0	0	0	0	0	<i>X</i> ·	+1	\Rightarrow	on a t	ypical puter
		0	0	0	0	0	U	0	0	0	0	1	U	0	0	0	0	~X & (<u>x + 1</u>)		comp	
	(x) oute.	(x) 15 1 1 0	(x) 15 14 1 0 1 0 0 1 0 1 0 1 0 1 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0	$\begin{array}{c} (x) \\ 15 & 14 & 13 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{array}$ bute. $\begin{array}{c} 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{array}$	$\begin{array}{c} (x) \\ 15 & 14 & 13 & 12 \\ 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{array}$	$\begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{array}$	(x)1514131211101011111110101010011011010011011011100110110110110100000000	(x)1514131211109101111101010101010110101000110101010110101010110101010110101000000000	(x)15141312111098101111011011110110101010011010101011010010010100011010110110100000000000	(x)15141312111098710111110111010111011101010101110110101010110101101011010110101101011010000000000	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 &$	$\begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 &$	$ \begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 &$	$ \begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 &$	$ \begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 &$	$ \begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 &$	$ \begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 1 &$	$ \begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 & r(x) \\ 1 & 0 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 1 &$	$\begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 & r(x) & R(x) \\ 1 & 0 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 1 &$	$\begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 & r(x) & R(x) & R(x) \\ 1 & 0 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 1 &$	$\begin{array}{c} (x) \\ 15 & 14 & 13 & 12 & 11 & 10 & 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 & r(x) & R(x) & R(x)_2 \\ 1 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 &$

Bit-whacking magic:

r(x) is also "easy" to compute (don't ask).

Bottom line: r(x) and R(x) can be computed with just a few machine instructions.

Probabilistic counting (Flajolet and Martin, 1983)

Maintain a single-word *sketch* that summarizes a data stream $x_0, x_1, ..., x_N, ...$

- For each x_N in the stream, update sketch by *bitwise or* with $R(x_N)$.
- Use position of rightmost 0 in sketch to estimate lg N.



Probabilistic counting trace

X	r(x)	R(x)	sketch
$011000100110001110100111101110 {\color{black} 1}{\color{black} 1}$	2	100	00000000000000000000000000000000000000
$0110011100100011000111110000010 {\color{black} 1}$	1	10	0000000000000000000000000000000001 1 0
$000100010001110001101101100 {\color{black} 1}{\color{black} 1}$	2	100	00000000000000000000000000000000000000
$010001000111011100000001110 {\color{red} 11111}$	5	100000	00000000000000000000000000000000000000
011010000101100010111000100100	0	1	000000000000000000000000000010011 1
$00110111101100000001010010101010{\color{black}1}$	1	10	000000000000000000000000000000000001 1 1
0011010001100011101010111111100	0	1	000000000000000000000000000010011 1
$000110000100001001011100110 {\color{blue}{111}}$	3	1000	00000000000000000000000000000000000000
$0001100110011001110010000 {\color{red} 11111} \\ \\$	6	1000000	00000000000000000000000000000000000000
01000101110001001010110011111100	0	1	00000000000000000000000000000000000000

 $R(sketch) = 10000_2$ = 16

31

Probabilistic counting (first try)

```
public long R(long x)
{ return ~x & (x+1); }
public long estimate(Iterable<String> stream)
{
    long sketch;
    for (s : stream)
        sketch = sketch | R(s.hashCode());
    return R(sketch);
}
```

Maintain a *sketch* of the data

• A single word

• OR of all values of R(x) in the stream Estimate is smallest value not seen.

Early example of "a simple algorithm whose analysis isn't"

Q. (Martin) Seems a bit low. How much?

A. (unsatisfying) Obtain empirically.

Probabilistic counting (Flajolet and Martin)

```
public long R(long x)
{ return ~x & (x+1); }
public long estimate(Iterable<String> stream)
{
    long sketch;
    for (s : stream)
        sketch = sketch | R(s.hashCode());
    return R(sketch)/.77351;
}
```

Maintain a *sketch* of the data

• A single word

• OR of all values of R(x) in the stream Estimate is smallest value not seen.

with correction for bias

Early example of "a simple algorithm whose analysis isn't"

Q. Value of correction factor? How accurate?

A. (Flajolet) Do the math!



Mathematical analysis of probabilistic counting

Theorem. The expected number of trailing 1s in the PC sketch is

 $lg(\phi N) + P(lg N) + o(1)$ where $\phi \doteq .77351$

and P is an oscillating function of lg N of very small amplitude.

Proof (omitted).

1980s: Flajolet tour de force

1990s: trie parameter

21st century: standard analytic combinatorics



Kirschenhofer, Prodinger, and Szpankowski Analysis of a splitting process arising in probabilistic counting and other related algorithms ICALP 1992.

In other words. In PC code, R(sketch)/.77351 is an *unbiased statistical estimator* of N.

Validation of probabilistic counting

Hypothesis. Expected value returned is *N for random values from a large range*.

Quick experiment. 100,000 31-bit random values (20 trials)



Flajolet and Martin: Result is "typically one binary order of magnitude off."

```
Of course! (Always returns a power of 2 divided by .77351.)
```

Need to incorporate more experiments for more accuracy.

16384/.77351 = 21181 32768/.77351 = 42362 65536/.77351 = 84725

....

Cardinality Estimation

- Rules of the game
- Probabilistic counting
- Stochastic averaging
- Refinements



PCSA trace

use initial m bits for second hash			M =	= 4	
x	R(x)	sketch[0]	sketch[1]	sketch[2]	sketch[3]
10 100111101110 11	100	00000000000000000	00000000000000000	0000000000000 1 00	000000000000000000000000000000000000000
00 0111110000010 1	10	0000000000000010	000000000000000000	000000000000100	000000000000000000000000000000000000000
01 101101101100 11	100	000000000000010	0000000000000 1 00	000000000000100	000000000000000000000000000000000000000
00 000001110 11111	100000	000000000 1 00010	000000000000100	000000000000100	000000000000000000000000000000000000000
01 01110001000100	1	000000000100010	00000000000010 1	000000000000100	000000000000000000000000000000000000000
00 0010100101010 1	10	000000000100010	000000000000101	000000000000100	000000000000000000000000000000000000000
10 1010111111100	1	000000000100010	000000000000101	00000000000010 1	000000000000000000000000000000000000000
00 01011100110 111	1000	00000000010 1 010	000000000000101	000000000000101	000000000000000000000000000000000000000
11 10010000 111111	1000000	00000000101010	000000000000101	000000000000101	00000000 1 000000
10 1011001111110 1	10	00000000101010	000000000000101	000000000001 1 1	000000001000000
00 01110100110100	1	0000000010101	000000000000101	000000000000111	
		000000000101011	0000000000000101	000000000000111	000000001000000
r (sketch[])		2	1	3	0

Probabilistic counting with stochastic averaging in Java

```
public static long estimate(Iterable<Long> stream, int M)
{
    long[] sketch = new long[M];
    for (long x : stream)
    {
        int k = hash2(x, M);
        sketch[k] = sketch[k] | R(x);
    }
    int sum = 0;
    for (int k = 0; k < M; k++)
        sum += r(sketch[k]);
    double mean = 1.0 * sum / M;
    double phi = .77351;
    return (int) (M * Math.pow(2, mean)/phi);
}</pre>
```

Idea. Stochastic averaging

- Use second hash to split into
 M = 2^m independent streams
- Use PC on each stream, yielding 2^m sketches .
- Compute *mean* = average # trailing 1 bits in the sketches.
- Return 2^{mean}/.77351.

Flajolet-Martin 1983

Theoretical analysis of PCSA

Definition. The *relative accuracy* is the standard deviation of the estimate divided by the actual value.

LEMMA 4. Setting $\beta = 2^{1/q}$, with $q \ge 1$, one has for fixed q

 $\mathbb{E}[\beta^{R_n}] = n^{1/q}(d_q + P_q(\log_2 n)) + o(n^{1/q}),$

where

Theorem (paraphrased to fit context of this talk).

Under the uniform hashing assumption, PCSA

- Uses 64M bits.
- Produces estimate with a relative accuracy close to $0.78/\sqrt{M}$.

$$1 - \left(1 - \frac{1}{2^k}\right)^n - \left(1 - \frac{1}{2^{k-1}}\right)^n + \left(1 - \frac{1}{2^k} - \frac{1}{2^{k-1}}\right)^n$$

a quantity which is

$$1 - e^{-n/2^k + O(n/2^{2k})} - e^{-n/2^{k-1} + O(n/2^{2k})} + e^{-3n/2^{k-1} + O(n/2^{2k})}$$

or $O(n/2^{2k})$, which in the given range of values of k is $O(n^{-3/2}4^{-\delta})$. Thus

$$\sum_{k > (5/4) \log_2 n} 2^k p_{n,k} = O\left(n^{5/4 - 3/2} \sum_{\delta \ge 0} 4^{-\delta} 2^{\delta}\right) = O(n^{-1/4}),$$

and the same bound applies if 2 is replaced by β in the above sum.

We now consider the error that comes from the replacement of the $p_{n,k}$ by their asymptotic equivalent for "small" k. From the bounds of Theorem 2, one finds

$$\sum_{k \leq (5/4)\log_2 n} \beta^k \left[p_{n,k} - \psi\left(\frac{n}{2^k}\right) + \psi\left(\frac{n}{2^{k+1}}\right) \right] = O\left(\frac{n^{5/4q}}{n^{0.49}}\right) = O(n^{0.76/q}), \quad (29)$$

LEMMA 5. If n elements are distributed into m cells (m fixed), where the probability that any element goes to a given cell has probability 1/m, then the probability that at least one of the cells has a number of elements N satisfying

 $|N-n/m| > \sqrt{n} \log n$

h > 0.

-1/m; let N_1 be the number of elements that fall into istribution

$$\mathbf{r}(N_1 = k) = \binom{n}{k} p^k q^{n-k},\tag{30}$$

and taking logarithms of (30), for $k = pn + \delta$ and $\delta \ll n$, one finds

$$\Pr(N_1 = pn + \delta) = \exp\left(-\frac{\delta^2 + O(\delta)}{2npq} + O\left(\frac{\delta^3}{n^2}\right)\right).$$

If $\delta = \sqrt{n} \log n$, the probability (30) is exponentially small. We conclude the proof by observing that the binomial distribution is unimodal and

8)
$$\Pr\left[\bigcup_{1 \le j \le m} \left| N_j - \frac{n}{m} \right| > \sqrt{n} \log n \right] < m \Pr\left[\left| N_1 - \frac{n}{m} \right| > \sqrt{n} \log n \right].$$

We can now conclude the proof of the first part of Theorem 4. Let S denote the sum $R^{\langle 1 \rangle} + R^{\langle 2 \rangle} + \cdots + R^{\langle m \rangle}$. We have

$$\Pr(S=k) = \sum_{\substack{n_1+n_2+\cdots+n_m=n\\k_1+k_2+\cdots+k_m=k}} \frac{1}{m^n} \binom{n}{n_1, n_2, \dots, n_m} p_{n_1, k_1} p_{n_2, k_2} \cdots p_{n_m, k_m}.$$
 (31)

40

Validation of PCSA analysis

Hypothesis. Value returned is accurate to $0.78/\sqrt{M}$ for random values from a large range.

Experiment. 1,000,000 31-bit random values, *M* = 1024 (10 trials)

% java PCSA 1000000 31 1024 10
964416
997616
959857
1024303
972940
985534
998291
996266
959208
1015329



Space-accuracy tradeoff for probabilistic counting with stochastic averaging



Bottom line.

- Attain 10% relative accuracy with a sketch consisting of 64 words.
- Attain 2.4% relative accuracy with a sketch consisting of 1024 words.

Scientific validation of PCSA

Hypothesis. Accuracy is as specified for the hash functions we use and the data we have.

Validation (Flajolet and Martin, 1985). Extensive reproducible scientific experiments (!)

Validation (RS, this morning).

% java PCSA 6000000 1024 < log.07.f3.txt 1106474

<1% larger than actual value

log.07.f3.txt

109.108.229.102 pool-71-104-94-246.lsanca.dsl-w.verizon.net 117.222.48.163 pool-71-104-94-246.lsanca.dsl-w.verizon.net 1.23.193.58 188.134.45.71 1.23.193.58 gsearch.CS.Princeton.EDU pool-71-104-94-246.lsanca.dsl-w.verizon.net 81.95.186.98.freenet.com.ua 81.95.186.98.freenet.com.ua 81.95.186.98.freenet.com.ua CPE-121-218-151-176.lnse3.cht.bigpond.net.au

Q. Is PCSA effective?

A. ABSOLUTELY!

43

Summary: PCSA (Flajolet-Martin, 1983)

is a *demonstrably* effective approach to cardinality estimation

Q. About how many different values are present in a given stream?

PCSA

- Makes one pass through the stream.
- Uses a few machine instructions per value
- Uses *M* words to achieve relative accuracy $0.78/\sqrt{M}$

Results validated through extensive experimentation.

JOURNAL OF COMPUTER AND SYSTEM SCIENCES 31, 182-200 (1984) Probabilistic Counting Algorithms for Data Base Applications PHILIPPE FLANOLET INRIA, Recquencourt, 78133 Le Chesnay, France AND G. NIGEL MARTIN IBM Development Laboratory, Hursley Park, Winchester, Hampakire SO212JN, United Kingdom Received June 13, 1984; revised Anril 3, 1985 This paper introduces a class of probabilistic counting algorithms with which one can estimate the number of denistic identities in a large collection of data (typically a large file bandless of the standard of the standard standard standard standard with a standard standard with an elevation to a separations per denome scannel. The standards may ava-band on statistical observations made one bits of hashed values of mercots. They are by com-stantion statistical observations made one bits of hashed values of mercots. They are by com-stantion statistical observations without any deguadation of performances and prove epicially work in the context of data bases (quot permission). Clint scanner ress, as, 1. INTRODUCTION As data base systems allow the user to specify more and more complex queries, the need arises for efficient processing methods. A complex query can however generally be evaluated in a number of different manners, and the overall perfor-mance of a data base system depends rather crucially on the selection of appropriate decomposition strategies in each particular case. Even a problem as trivial as computing the intersection of two collections of data. A and B lends itself to a number of different treatments (see, e.g., [7]): ²⁰B' 1. Sort A, search each element of B in A and retain it if it appears in A; 2. sort A, sort B, then perform a merge-like operation to determine the inter-3. eliminate duplicates in A and/or B using hashing or hash filters, then perform Algorithm 1 or 2. Each of these evaluation strategy will have a cost essentially determined by the number of records a, b in A and B, and the number of *distinct* elements α, β in A and B, and for typical sorting methods, the costs are: 182 0022-0000/85 \$3.00

Open questions

- Better space-accuracy tradeoffs?
- Support other operations?

"IT IS QUITE CLEAR that other observable regularities on hashed values of records could have been used...

Cardinality Estimation

- Rules of the game
- Probabilistic counting
- Stochastic averaging
- Refinements

logs and loglogs

To improve space-time tradeoffs, we need to *carefully count bits*.

Relevant quantities

- N is the number of items in the data stream.
- $\lg N$ is the number of bits needed to represent numbers less than N in binary.
- Ig Ig N is the number of bits needed to represent numbers less than Ig N in binary.

For most applications

- N is less than 2⁶⁴.
- Ig N is less than 64.
- lg lg *N* is less than 7.

Typical PCSA implementations

- Could use *M* lg *N* bits, in theory.
- Use 64-bit words to take advantage of machine-language efficiencies.
- Use (therefore) 64*64 = 4096 bits with M = 64 (for 10% accuracy with $N < 2^{64}$).



We can do better: HyperLogLog algorithm (2007)



Flajolet-Fusy-Gandouet-Meunier 2007

Theorem (paraphrased to fit context of this talk).

Under the uniform hashing assumption, HyperLogLog

- Uses M log log N bits.
- Achieves relative accuracy close to $1.02/\sqrt{M}$.

Space-accuracy tradeoff for HyperLogLog



Typical PCSA implementations

- Could use *M* lg *N* bits, in theory.
- Use 64-bit words to take advantage of machine-language efficiencies.
- Use (therefore) 64*64 = 4096 bits with M = 64 (for 10% accuracy with $N < 2^{64}$).

Typical Hyperloglog implementations

- Could use *M* lg lg *N* bits, in theory.
- Use 8-bit bytes to take advantage of machine-language efficiencies.
- Use (therefore) 64*8 = 512 bits with M = 64 (for 10% accuracy with $N < 2^{64}$).



Right answer: Hyperloglog

Divide into *M* streams (stochastic averaging)

- Keep track of min(# trailing 1s).
- Use harmonic mean.

```
public static long estimate(Iterable<Long> stream, int M)
{
    byte[] bytes = new byte[M];
    for (long x : stream)
    {
        int k = hash2(x, M);
        if (bytes[k] < Bits.r(x)) bytes[k] = Bits.r(x);
    }
    double sum = 0.0;
    for (int k = 0; k < M; k++)
        sum += Math.pow(2, -1.0 - bytes[k]);
    return (int) (alpha * M * M / sum);
}</pre>
```



Validation of Hyperloglog



HyperLogLog in Practice: Algorithmic Engineering of a State of The Art Cardinality Estimation Algorithm. Extending Database Technology/International Conference on Database Theory 2013.

Summary for cardinality estimation algorithms

	time bound	memory bound (bits)	# bits for 10% accuracy for 1 billion inputs
Wrong answer	N ²	N lg N	64 billion
Sort and count	N log N	N lg N	64 billion
Existence table	Ν	U	1 billion
Hash table	N*	N lg N	64 billion
PCSA	N *	M lg N	4096
HyperLogLog	N*	M lglg N	512
HyperBitBit ?	N *	2 <i>M</i> + lglg <i>N</i>	134

Q. End of story?

A. Not quite. (RS+JL are working on an algorithm with 10% accuracy with 134 bits for $N < 2^{64}$.)

A Case Study: Cardinality Estimation

- Exact cardinality count
- Probabilistic counting
- Stochastic averaging
- Refinements

Hyperloglog is a case in point.

- Impact is broad and far-reaching.
- Old roots, new opportunities.
- Allows solution of otherwise unsolvable problems.
- Intellectually stimulating.
- Implementations teach programming proficiency.
- May unlock the secrets of life and of the universe.
- Useful for fun and profit.



One more thing

This was the last live lecture for COS 488

- We are online, with studio-produced videos.
- Same (scalable) model as for COS 126 and COS 226.





Yes, on Wednesdays

- Check your e-mail.
- Watch Lecture 1 online.
- Details and discussion Wednesday.

20th Century



21st Century

