

# Data Structures and Algorithms

Spring 2014

Cesar Agustin Garcia Vazquez

April 8, 2014

## 1 Problem Set 06

1. Suppose that Harvard ID numbers were issued randomly, with replacement. That is, your Harvard ID would consist of just 8 randomly generated digits, and no check was made to ensure that the same number was not issued twice. You might use the last four digits of your ID number as a password. Try to give exact numerical answers for the following questions:
  - a) How many people would you need to have in a room before it was more likely than not that two had the same last four digits?
  - b) How many numbers could be issued before it would be more likely than not that there is a duplicate Harvard ID number?
  - c) What would the answers for the above questions be if there were 12 digit ID numbers?

**Answer:** If we consider the birthday problem, the second person has the probability of collision equals to  $1/n$  because only one day over  $n$  has been already used. In this case, instead of  $1/n$  it is going to be  $(m+1)/n$  where  $m$  is all the possible appearance of the last 4 digits number. If we have  $n$  from 0 to 99,999 and we take only the last 4 digits, we can see that any password could be picked from 10 different ids. If we have 100,000,000 ids (0 to 99,999,999), then  $m$  becomes 10,000, i.e., the number of possible ids  $10^8$  minus the number of elements that we use as passwords  $10^4$ . So in this case, the probability of no collisions is given by equation (1).

$$\begin{aligned} P(B_k) &= \left(\frac{n-m}{n}\right) \cdot \left(\frac{n-m-1}{n}\right) \cdot \left(\frac{n-m-2}{n}\right) \cdots \left(\frac{n-m-k+1}{n}\right) \\ &= \left(1 - \frac{m}{n}\right) \cdot \left(1 - \frac{m+1}{n}\right) \cdot \left(1 - \frac{m+2}{n}\right) \cdots \left(1 - \frac{m+k-1}{n}\right) \quad (1) \end{aligned}$$

If we consider the inequality  $e^{-x} \geq 1 - x$ , then we reach the result in equation (2).

$$\begin{aligned}
P(B_k) &= \left(1 - \frac{m}{n}\right) \cdot \left(1 - \frac{m+1}{n}\right) \cdot \left(1 - \frac{m+2}{n}\right) \cdots \left(1 - \frac{m+k-1}{n}\right) \\
&\leq e^{-\frac{m}{n}} \cdot e^{-\frac{m+1}{n}} \cdot e^{-\frac{m+2}{n}} \cdots e^{-\frac{m+k-1}{n}} \\
&= e^{-\frac{m}{n} - \frac{m+1}{n} - \frac{m+2}{n} - \cdots - \frac{m+k-1}{n}} \\
&= e^{-\sum_{i=0}^{k-1} \frac{m+i}{n}} \\
&= e^{-\sum_{i=0}^{k-1} \left(\frac{m}{n} + \frac{i}{n}\right)} \\
&= e^{-\left(\sum_{i=0}^{k-1} \frac{m}{n} + \sum_{i=0}^{k-1} \frac{i}{n}\right)} \\
&= e^{-\left(\frac{m}{n} \sum_{i=0}^{k-1} 1 + \frac{1}{n} \sum_{i=0}^{k-1} i\right)} \\
&= e^{-\left(\frac{m}{n} k + \frac{1}{n} \frac{k(k-1)}{2}\right)} \\
&= e^{-\left(\frac{mk}{n} + \frac{k(k-1)}{2n}\right)} \\
&= e^{-\left(\frac{2mk+k(k-1)}{2n}\right)}
\end{aligned} \tag{2}$$

So we want to know  $k$  such that  $P(B_k) \leq e^{-\left(\frac{2mk+k(k-1)}{2n}\right)} < 1/2$ , which can be reached equation (3).

$$\begin{aligned}
e^{-\left(\frac{2mk+k(k-1)}{2n}\right)} &< \frac{1}{2} \\
-\left(\frac{2mk+k(k-1)}{2n}\right) &< \log\left(\frac{1}{2}\right) \\
-\left(\frac{2mk+k(k-1)}{2n}\right) &< \log(2^{-1}) \\
-\left(\frac{2mk+k(k-1)}{2n}\right) &< -\log(2) \\
\frac{2mk+k(k-1)}{2n} &\geq \log(2) \\
2mk+k(k-1) &\geq 2n \log(2)
\end{aligned} \tag{3}$$

From equation (3), we know that  $n = 100,000,000$  and  $m = 10,000$ , we reach in equation in (4), and by solving it using [sagemath.org](http://sagemath.org), we reach to 2 roots and we focus only in the real root, which means that  $k = 5447.8116$ , i.e., that is, with 5,447 people in a room, the probability of password collision is less than 0.5

$$\begin{aligned}
2mk+k(k-1) &\geq 2n \log(2) \\
2(10,000)k+k(k-1) &\geq 2(100,000,000) \log(2) \\
20,000k+k^2-k-200,000,000 \log(2) &\geq 0 \\
k^2+19,999k-200,000,000(0.693) &\geq 0 \\
k^2+19,999k-200,000,000(0.693) &\geq 0
\end{aligned}$$

$$k^2 + 19,999k - 138,629,436.112 \geq 0 \quad (4)$$

For the second question, we just have to do  $m = 1$ , and then we reach inequality (5), and then solving for  $k$ , we have that the number of ids that could be issued before a collision are 11,774.

$$\begin{aligned} 2mk + k(k-1) &\geq 2n \log(2) \\ 2k + k(k-1) &\geq 2(100,000,000) \log(2) \\ 2k + k^2 - k - 200,000,000 \log(2) &\geq 0 \\ k^2 + k - 200,000,000(0.693) &\geq 0 \\ k^2 + k - 200,000,000(0.693) &\geq 0 \\ k^2 + k - 138,629,436.112 &\geq 0 \end{aligned} \quad (5)$$

For the third question,  $m = 10,000$  but  $n = 1,000,000,000,000$ , so, we reach the inequality (6), and then solving for  $k$  we have that the number of people in a room, for the probability of password collision to be less than 0.5, is 1,167,452.

$$\begin{aligned} 2mk + k(k-1) &\geq 2n \log(2) \\ 2(10,000)k + k(k-1) &\geq 2(1,000,000,000,000) \log(2) \\ 20,000k + k^2 - k - 2,000,000,000,000 \log(2) &\geq 0 \\ k^2 + 19,999k - 2,000,000,000,000(0.693) &\geq 0 \\ k^2 + 19,999k - 1,386,294,361,120 &\geq 0 \end{aligned} \quad (6)$$

The number of ids is also given by inequality 7, and then solving for  $k$  we have that the number of IDs that can be issued, with probability of collision less than 0.5, is 1,177,409.

$$\begin{aligned} 2mk + k(k-1) &\geq 2n \log(2) \\ 2(10,000)k + k(k-1) &\geq 2(1,000,000,000,000) \log(2) \\ 20,000k + k^2 - k - 2,000,000,000,000 \log(2) &\geq 0 \\ k^2 + k - 2,000,000,000,000(0.693) &\geq 0 \\ k^2 + k - 1,386,294,361,120 &\geq 0 \end{aligned} \quad (7)$$

2. Suppose each person gets a random hash value from the range  $[1, \dots, n]$ . (For the case of birthdays,  $n$  would be 365.) Show that for some constant  $c_1$ , where there are at least  $c_1\sqrt{n}$  people in a room, the probability that no two have the same hash value is at most  $1/e$ . Similarly, show that for some constant  $c_2$  (and sufficiently large  $n$ ), when there are at most  $c_2\sqrt{n}$  people in the room, the probability that no two

have the same hash value is at least  $1/2$ . Make these constants as close to optimal as possible.

Hint: you may use the fact in equation (8)

$$e^{-x} \geq 1 - x \quad (8)$$

and in equation (9).

$$e^{-x-x^2} \leq 1 - x \text{ for } x \leq \frac{1}{2} \quad (9)$$

You may feel free to find and use better bounds.

**Answer:** For this problem, we are going to consider the general case of the birthday paradox problem, where the probability of no collisions when we have  $k$  people over  $n$  numbers is given by equation (10).

$$\begin{aligned} P(B_k) &= \left(\frac{n-0}{n}\right) \cdot \left(\frac{n-1}{n}\right) \cdot \left(\frac{n-2}{n}\right) \cdots \left(\frac{n-k+1}{n}\right) \\ &= 1 \cdot \left(1 - \frac{1}{n}\right) \cdot \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{k-1}{n}\right) \end{aligned} \quad (10)$$

If we consider the inequality (8), then we can reach the result in equation (11).

$$\begin{aligned} P(B_k) &\leq e^{-1/n} \cdot e^{-2/n} \cdots e^{-(k-1)/n} \\ &= e^{-\sum_{i=1}^{k-1} \frac{i}{n}} \\ &= e^{-\frac{k(k-1)}{2n}} \end{aligned} \quad (11)$$

For the first case, we want that  $P(B_k) \leq 1/e$ , so if we consider that  $e^{-\frac{k(k-1)}{2n}} \approx e^{-\frac{k^2}{2n}}$ , then we reach (12), which means that  $c_1 = \sqrt{2} \approx 1.414$

$$\begin{aligned} e^{-\frac{k^2}{2n}} &\leq \frac{1}{e} \\ -\frac{k^2}{2n} &\leq \log\left(\frac{1}{e}\right) \\ -k^2 &\leq 2n \log(e^{-1}) \\ -k^2 &\leq -2n \\ k^2 &> 2n \\ k &> \sqrt{2n} \\ k &> \sqrt{2}\sqrt{n} \end{aligned} \quad (12)$$

For the second case, we want that  $P(B_k) \geq 1/2$ , so if we consider that  $e^{-\frac{k(k-1)}{2n}} \geq e^{-\frac{k^2}{2n}}$ , then we reach (13), which means that  $c_2 = \sqrt{2\log(2)} \approx 1.177$

$$\begin{aligned}
P(B_k) &\geq \frac{1}{2} \\
e^{-\frac{k(k-1)}{2n}} &\geq \frac{1}{2} \\
e^{-\frac{k(k-1)}{2n}} &\geq e^{-\frac{k^2}{2n}} \\
e^{-\frac{k^2}{2n}} &\geq \frac{1}{2} \\
-\frac{k^2}{2n} &\geq \log\left(\frac{1}{2}\right) \\
-k^2 &\geq 2n \log(2^{-1}) \\
-k^2 &\geq -2n \log(2) \\
k^2 &< 2 \log(2)n \\
k &< \sqrt{2 \log(2)n} \\
k &< \sqrt{2 \log(2)}\sqrt{n}
\end{aligned} \tag{13}$$

3. For the document similarity scheme described in class, it would be better to store fewer bytes per document. Here is one way to do this, that uses just 48 bytes per document: take an original sketch of a document, using 84 different permutations. Divide the 84 permutations into 6 groups of 14. Re-hash each group of 14 values to get 6 new 64 bit values. Call this the *super-sketch*. Note that for each of the 6 values in the super-sketch, two documents will agree on a value when they agree on all 14 of the corresponding values in the sketch. Why does it make sense to simply assume that this is the only time a match will occur?

Consider the probability that two documents with resemblance  $r$  agree on two or more of the six sketches. Write equations that give this probability and graph the probability of a function of  $r$ . Explain and discuss your results.

What happens if instead of using a 64 bit hash value for each group in the super-sketch, we only use a 16 bit hash? An 8 bit hash?

4. Prove that 636,127 is composite by finding an appropriate witness. Be sure to give ample evidence showing that your witness is in fact witnesses. (Note: do not use a factor as a witness!) Sure, these numbers are small enough that you can exhaustively find a factor; that is not the point. A factor is not a witness, according to our definition.) Hint: you will want to write some code. You will preferably use a package that deals with big integers appropriately, as you may want to use some of this code for the next problem (RSA). We don't need a code listing for this problem, a short summary of the output should suffice.

The number 294, 409 is a Carmichael number. Prove that it is composite by finding

a witness. Briefly explain why Fermat's little theorem won't help.

**Answer:** For the number 636,127, we first generate a random number  $w$ , between  $(1, 636, 127)$  and then we just applied Fermat's little Theorem. The trick here is to correctly compute  $w^{636,126} \bmod 636, 127$ , which can be easily done using the java class *java.math.BigInteger*. If the remainder is different than 1, that means that the number is composite and then we stop and log the witness and the remainder obtained.

The number 360,850 acts as a witness that 646,127 composite, since it is  $360, 850 \equiv 36, 085 \bmod 636, 127$ .

To find a witness for the Carmichael number, Fermat's Little Theorem does not work directly with any number, since it will return the remainder as 1 if the possible witness share a factor with 294,409. This might lead to computing a lot of unnecessary computations if we don't consider this condition.

So the tweak is that after getting randomly the possible witness, we compute the greater common divisor, and if it is 1, then we drop that possible witness and we retrieve another one.

After adding the dropping condition, the program finds the witness 88,987 which returned the remainder 16,133, so we have that  $88, 987 \equiv 16, 133 \bmod 294, 409$ .

5. My RSA public key is (46947848749720430529628739081, 37267486263679235062064536973). Convert the message

Give me an A

into a number, using ASCII in the natural way. (So for "A b": in ASCII, A = 65, space = 32, and b = 98; translating each number into 8 bits gives "A b" = 010000010010000001100010 in binary.) Encode the message as though you were sending it to me using my RSA key, and write for me the corresponding encoded message in decimal.

**Answer:** First we convert every character into its corresponding ASCII number and then into 8 bits. This is shown in Table 1.

If we append each binary string after another, we get the string in equation (14).

$$\begin{array}{l} 010001110110100101110110011001010010000001101101 \\ 011001010010000001100001011011100010000001000001 \end{array} \quad (14)$$

The binary number in equation (14) of 96 bits (8 bits per 12 characters), represents the number  $x = 22,100,932,013,077,800,977,026,654,273$ , which can be obtained by using the java class *java.math.BigInteger*. If we consider the public key  $(n, e)$  given, then message to be sent is in the range  $[1, \dots, n]$ . So we just need to compute  $e(x) = x^e \bmod n$ , which yields  $e(x) = 27, 016, 764, 340, 118, 192, 395, 712, 492, 378$ .

Table 1: Values of each character in the message

Character	ASCII	Binary
G	71	0100 0111
i	105	0110 1001
v	118	0111 0110
e	101	0110 0101
	32	0010 0000
m	109	0110 1101
e	101	0110 0101
	32	0010 0000
a	97	0110 0001
n	110	0110 1110
	32	0010 0000
A	65	0100 0001

The above result was computed using *java.math.BigInteger#modPow* which internally breaks the operations and then it uses the Chinese Remainder Theorem to combine results.