

## Midterm Exam

- Don't Panic.
- The midterm contains six problems (and one just for fun). You have 100 minutes to earn 100 points.
- The midterm contains 18 pages, including this one and 4 pages of scratch paper.
- The midterm is closed book. You may bring one double-sided sheet of A4 paper to the midterm. (You may not bring any magnification equipment!) You may not use a calculator, your mobile phone, or any other electronic device.
- Write your solutions in the space provided. If you need more space, please use the scratch paper at the end of the midterm. Do not put part of the answer to one problem on a page for another problem.
- Read through the problems before starting. Do not spend too much time on any one problem.
- Show your work. Partial credit will be given. You will be graded not only on the correctness of your answer, but also on the clarity with which you express it. Be neat.
- You may use any algorithm given in class without restating it—simply give the name of the algorithm and the running time. If you change the algorithm in any way, however, you must provide complete details.
- Good luck!

Problem #	Name	Possible Points	Achieved Points
1	A Few Random Variables	12	
2	Independence Day	14	
3	Frenemy Fred	14	
4	Right in the Middle	20	
5	Building Strong Bonds	20	
6	Network Robustness	20	
<b>Total:</b>		100	

Student Number: \_\_\_\_\_

**Problem 1. A Few Random Variables**

Assume you have a collection of indicator random variables  $x_1, x_2, \dots, x_n$ . For each  $x_i$ , the  $\Pr[x_i = 1] = p_i$  (and otherwise  $x_i = 0$ ). The variables are *not independent*, and you do not know anything about how they are correlated. Which of the following statements are *always* true? For each part, explain why in less than ten words.

$$\Pr \left[ \sum_{i=1}^n x_i \geq 1 \right] \leq \sum_{j=1}^n p_j$$
TRUE      FALSE

**Explanation:**

$$\Pr \left[ \sum_{i=1}^n x_i \leq 0 \right] \leq \prod_{j=1}^n (1 - p_j)$$
TRUE      FALSE

**Explanation:**

$$\Pr \left[ \sum_{i=1}^n x_i \geq \sum_{j=1}^n \frac{p_j}{2} \right] \leq \frac{1}{2}$$
TRUE      FALSE

**Explanation:**

$$\Pr \left[ \sum_{i=1}^n x_i \geq (1 + \epsilon) \sum_{j=1}^n p_j \right] \leq e^{-(\epsilon^2/3) \sum_{j=1}^n p_j}$$
TRUE      FALSE

**Explanation:**

**Problem 2. Independence Day**

Recall that the L0-sampler is constructed from a collection of  $\log n$  distinct  $s$ -sparse samplers. To decide which indices are mapped to which samplers, we use  $\log n$  hash functions  $h_0, h_1, \dots, h_{\log n - 1}$  where  $h_i(x) = 1$  with probability  $1/2^i$  (and 0 otherwise). What happens if these hash functions are not independent?

Consider, instead, using a hash function  $h : [1, n] \rightarrow [1, n]$  that maps to a random integer from 1 to  $n$ . I.e., for every index  $k$  and every integer  $\ell \in [1, n]$  (independently for each  $k$ ):

$$\Pr[h(k) = \ell] = 1/n.$$

Then we define:  $h_i(k) = 1$  if  $h(k) \leq n/2^i$ . In this case, each hash function  $h_i$  has the proper probability distribution, but they are not independent.

If we use the  $h_i$  defined in this manner, does the L0-sampler still work as expected? If so, why? If not, why not? (Circle your answer and then explain.)

**WORKS**

**FAILS**

**Explanation:**

**Problem 3. Frenemy Fred.**

To pass the time, you create an internet service<sup>1</sup> that generates a stream of integers. You have hidden a secret pattern in the integers, and you are waiting for someone to find it! You notice that your frenemy Fred is monitoring the stream of integers, trying to figure out your code. Bored, you hack into his computer and discover that Fred is about to start running a Flajolet-Martin sketch to count the number of distinct items in your stream for the month of November. And you find the hash function that he will be using! Here is an excerpt from Fred's hash function:

1	2	3	4	5	6	7	8	9	10
0.221	0.467	0.619	0.689	0.413	0.263	0.716	0.292	0.890	0.831
11	12	13	14	15	16	17	18	19	20
0.026	0.110	0.179	0.761	0.474	0.926	0.289	0.333	0.037	0.550

(For example,  $h(13) = 0.179$ .)

You decide to mess with Fred a bit, and update your sequence generation routine to ensure that Fred's algorithm fails. And you want it to fail *as badly as possible*.

**Give a stream of integers of length 10 in which Fred's implementation of FM will return an answer that is at least 5 off:**

--	--	--	--	--	--	--	--	--	--

**What will Fred's implementation of FM return for your sequence?** (Your answer here can be approximate, e.g., round to the nearest integer.)

<sup>1</sup>You were lucky to get the URL [www.CodedStreamIsSequenceForUnderstandingNumbers.com](http://www.CodedStreamIsSequenceForUnderstandingNumbers.com)!

**Problem 4. Right in the Middle**

Your job is to monitor a stream of numbers:  $x_1, x_2, x_3, \dots$  (Perhaps these numbers are stock prices for Acme Corporation.) You do not know, in advance, how many numbers will be in the stream. Your task is to design an algorithm that uses as little space as possible and returns an approximate median for the values in the stream. To be more specific, given some error parameter  $\epsilon$ , your algorithm should return a value  $v$  that, with probability at least  $1 - \epsilon$ , is:

- bigger than at least  $1/3$  of the numbers in the stream, and
- smaller than at least  $1/3$  of the numbers in the stream.

(That is, it is between the 33rd and 67th percentile of the numbers.) For example, if the stream consists of the numbers: 1, 1, 2, 3, 5, 5, 6, 7, 9 (in any order), then your algorithm may return either the number 3 or the number 5.

**Problem 4.a.** Explain (succinctly) how your algorithm works.

**Problem 4.b.** How much space does your algorithm use? You may assume that each number in the stream can be stored in 1 unit of space.

Continued on the next page.

**Problem 4.c.** Prove that your algorithm is correct.

**Problem 5. Building Strong Bonds**

We are building a new social network that provides more precisely quantified relationships than Facebook. Instead of simply deciding whether someone is your friend or not, you can carefully adjust your relationship with each person. When you “up” a person, you strengthen your relationship, and when you “down” a person, you weaken your relationship. In this way, for each user, we can classify their relationships;

- If Alice gives Bob more “ups” than “downs,” then Bob is a friend of Alice.
- If Alice gives Bob more “downs” than “ups,” then Bob is an enemy of Alice.
- If Alice gives Bob the same number of “ups” and “downs,” then Alice and Bob have no relationship.

Notice that relationships are not symmetric. Alice may think Bob a friend, while Bob may think Alice an enemy!

The network produces a constant stream of “up” and “down” actions, e.g.,  $(up, u, v)$  to indicate that user  $u$  has given an “up” to user  $v$ , or  $(down, v, w)$  to indicate that user  $u$  has given a “down” to user  $v$ .

Continued on the next page.

**Problem 5.a.** Give an algorithm that will monitor the stream and decide whether or not there are any relationships in the graph. That is, if for every pair  $(u, v)$  the number of “up”s equals the number of “down”s, then your algorithm should return EMPTY. Otherwise, it should return FULL. Describe your algorithm and explain why it works.

What is the probability of error?

How much space does your algorithm take?

Continued on the next page.

**Problem 5.b.** Now describe an algorithm that will monitor the stream and decide whether there are more than 10 relationships in the graph. Show that your algorithm is correct.

What is the probability of error?

How much space does your algorithm take?

**Problem 6. Network Robustness**

Dr. Bessy Myst is an engineer tasked with keeping the company network up and running. And yet she is always worried about failures—links in the network just seem to keep failing. She has measured that each link in the network fails (over a given period of time) with probability about  $p$ , and she wants to know how likely the network is to remain connected throughout. Your job is to help Dr. Myst to estimate this probability.

More formally, Dr. Myst has a graph  $G = (V, E)$  which represents the topology of the network which contains  $n$  nodes and  $m$  edges. Your algorithm should take as input an error parameter  $\epsilon$  and a probability  $p$ . Assuming that each edge fails (i.e., is deleted) independently with probability  $p$ , define  $R(p)$  to be the *reliability* of the graph, i.e., the probability that the graph remains connected.

Your algorithm should return the probability  $R(p) \pm \epsilon$ , i.e., a value  $v$  such that:

$$R(p) - \epsilon \leq v \leq R(p) + \epsilon .$$

It should do this with probability at least  $2/3$ .

**Problem 6.a.** Describe an algorithm that outputs  $R(p) \pm \epsilon$  in polynomial time. (Your algorithm does not have to be sublinear.)

**Problem 6.b.** Prove that your algorithm is correct.

**Problem 6.c.** Dr. Myst decides to use her algorithm to determine the reliability of a streaming graph, i.e., a graph which is presented as a stream of edge additions and deletions. Explain how to modify your existing algorithm so that it can be used in the streaming context using  $o(n^2)$  space (i.e., less than  $n^2$  space, asymptotically). Assume  $p$  and  $\epsilon$  are fixed in advance, and you know the number of nodes  $n$  in the graph. Describe your algorithm and explain why it works.

Give the space needed by your algorithm as a function of  $n$  and  $\epsilon$ :

**Problem 6.d.** The CEO of Dr. Myst's company is excited to see her results. But she is unhappy that the approximation is additive. She asks Dr. Myst instead to devise an algorithm to calculate  $R(p)(1 \pm \epsilon)$ , i.e., a value  $v$  such that:

$$R(p)(1 - \epsilon) \leq v \leq R(p)(1 + \epsilon) .$$

Dr. Myst explains that this is much more difficult. Explain why. (Perhaps give an example where  $R(p)$  is very hard to approximate in this way. Perhaps explain where the analysis of your algorithm would fail. There are several different ways you might explain the problem.)

Note that there do exist polynomial-time approximation algorithms for  $(1 - R(p))$ , so it is not impossible. But there are reasons why it is harder.

**Scratch Paper**

**Scratch Paper**

**Scratch Paper**

**Scratch Paper**