

Problem Set 5

Due: Wednesday, February 15, by 11:59pm

Instructions

Solutions format and late policy. See PSet 1 for further details. The same requirements and policies still apply. Also follow the typesetting instructions from the prior PSets.

Collaboration policy. The written problems on this pset may be done with a **single partner**. In this case, only one person will submit the written part on Gradescope and add their partner as a collaborator. You must do the coding part (Task 6a) on your own.

Solutions submission. You must submit your solution via Gradescope. In particular:

- For the solutions to Task 1-5, as well as parts b) and c) of Task 6, submit under “PSet 5 [Written]” a **single** PDF file containing the solution to all tasks in the homework (for you and your partner). Each numbered task should be solved on its own page (or pages). Follow the prompt on Gradescope to link tasks to your pages. Do not write your names on the individual pages – Gradescope will handle that.
- For the programming part (Task 6 part a)), submit your code under “PSet 5 [Coding]” as a file called bloom.py.

Task 1 – On day or off day

[10 pts]

A student is getting ready to take an important oral examination and is concerned about the possibility of having an “on” day or an “off” day. If the student has an on day, each of the examiners will pass him independently with probability 0.8, whereas if he has an off day, this probability will be reduced to 0.4. Suppose that the student will pass the exam if a majority of the examiners pass him. If the student feels that he is twice as likely to have an off day as he is to have an on day, should he request an examination with 3 examiners or with 5 examiners?

Task 2 – Geometric and Poisson

[18 pts]

Let X be Geometric with parameter p , let Y be Poisson with parameter λ and $Z = \max(X, Y)$. Assume that X and Y are independent. For each of the following problems, your final answers should not have summations. You may want to use the Taylor series expansion of e^x , that is, that for any x ,

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

- a) (3 points) What is $\mathbb{P}(X > k)$? [Think about it from first principles]
- b) (6 points) Compute $\mathbb{P}(X > Y)$. Hint: Use the law of total probability to obtain that

$$\mathbb{P}(X > Y) = \sum_{k=0}^{\infty} \mathbb{P}(X > Y | Y = k) \cdot \mathbb{P}(Y = k).$$

- c) (4 points) Compute $\mathbb{P}(Z \geq X)$.
- d) (5 points) Compute $\mathbb{P}(Z \leq Y)$.

Task 3 – Partnerships

[16 pts]

In a class of N students, there are $M = \binom{N}{2}$ possible “partnerships” (where a partnership is just a pair of people that will work on a problem set together).

Answer each of the following questions. Make sure that each of your answers is **not** in the form of a summation for this problem. *In each case include the expectation and variance of N as part of your answer.*

- a) What is the expected value of M if N equals some fixed positive integer c with probability 1? (Your answer will be a function of c .)
- b) What is the expected value of M if N has a Poisson distribution with parameter λ ? (Your answer will be a function of λ .)
- c) What is the expected value of M if N has a geometric distribution with parameter p ? (Your answer will be a function of p .)
- d) What is the expected value of M if $N = 10X + 7$, where X is a Bernoulli random variable with parameter p ? (Your answer will be a function of p .)

Task 4 – Roll away

[12 pts]

Suppose that a fair 6-sided die is rolled repeatedly, with each roll independent of the others. Let Z be the number of rolls until (and including) the first time either a 2 or a 3 is rolled, and let W be the number of 6's rolled until the first 2 or 3 is rolled. So, for example if the sequence of die values until the first 2 or 3 is 1,5,4,4,5,6,1,6,2, then Z is 9 and W is 2.

Define

$$p(j) := \begin{cases} \mathbb{P}(W = j \mid Z = i) & j \in \{0, 1, \dots, i - 1\} \\ 0 & \text{otherwise} \end{cases}$$

Show that $p(j)$ is the probability mass function of a binomially distributed random variable and determine its parameters n and p .

Task 5 – Sample Sampling Algorithm

[18 pts]

Consider the following algorithm for generating a random sample S of size n from the set of integers $\{1, 2, \dots, N\}$, where $0 < n < N$.

```
Sample( $N, n$ ):  
   $S \leftarrow \emptyset$  //  $S$  is a set of distinct integers, initially an empty set  
  while  $|S| < n$  do  
     $x \leftarrow \text{RollDie}(N)$  //  $x$  is the outcome of rolling a fair  $N$ -sided die  
     $S \leftarrow S \cup \{x\}$  // if  $x$  is already in  $S$  it doesn't change  
  return  $S$ 
```

Let I be the number of die rolls until S is returned. Also, let I_i be the random variable which describes the number of rolls it takes from the time the set S has $i - 1$ values to the first time a new value is added after that (i.e., the set S has i values).

- a) What type of random variable from our zoo is I_i and what is/are the relevant parameter(s) for that random variable?
- b) What is I in terms of the random variables I_i ? Calculate $\mathbb{E}[I]$, expressing the result as a summation that depends on both N and n .

c) What is $\text{Var}(I)$? You can leave your answer in summation form.

Task 6 – Bloom filters [Coding+Written]

[10+7 pts]

Google Chrome has a huge database of malicious URLs, but it takes a long time to do a database lookup (think of this as a typical [Set](#)). They want to have a quick check in the web browser itself, so a space-efficient data structure must be used. A **Bloom filter** is a **probabilistic data structure** which only supports the following two operations:

- `add(x)`: Add an element x to the structure.
- `contains(x)`: Check if an element x is in the structure. If either returns “definitely not in the set” or “could be in the set”.

It does **not** support the following two operations:

- delete an element from the structure.
- return an element that is in the structure.

The idea is that we can check our Bloom filter to see if a URL is in the set. The Bloom filter is always correct in saying a URL definitely isn't in the set, but may have false positives – it may say that a URL is in the set when it isn't. Only in these rare cases does Chrome have to perform an expensive database lookup to know for sure.

Suppose that we have k **bit arrays** t_1, \dots, t_k each of length m (all entries are 0 or 1), so the total space required is only km bits or $km/8$ bytes (as a byte is 8 bits). Suppose that the universe of URL's is the set \mathcal{U} (think of this as all strings with less than 100 characters), and we have k **independent and uniform** hash functions $h_1, \dots, h_k : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$. That is, for an element x and hash function h_i , pretend $h_i(x)$ is a **discrete** $\text{Unif}[0, m-1]$ random variable. Suppose that we implement the `add` and `contains` function as follows:

Algorithm 1 Bloom Filter Operations

```
1: function INITIALIZE(k,m)
2:   for  $i = 1, \dots, k$ : do
3:      $t_i =$  new bit array of  $m$  0's
4: function ADD(x)
5:   for  $i = 1, \dots, k$ : do
6:      $t_i[h_i(x)] = 1$ 
7: function CONTAINS(x)
   return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

Refer to Section 9.4 of the textbook and the relevant lecture for more details on Bloom filters.

a) Implement the functions `add` and `contains` in the `BloomFilter` class of `bloom.py`.

To solve this task, we have set up a corresponding edstem lesson [here](#). Press the Mark button above the terminal to run the unit tests we have written for you. Passing these unit tests is not enough. We have written a number of different tests for the Gradescope autograder. Your score on Gradescope will be your actual score - you have unlimited attempts to submit.

b) Let's compare this approach to using a typical [Set](#) data structure. Google wants to store 1 million URLs, with each URL taking (on average) 23 bytes.

- How much space (in MB, 1 MB = 1 million bytes) is required if we store all the elements in a set?
- How much space (in MB) is required if we store all the elements in a Bloom filter with $k = 10$ hash functions and $m = 800,000$ buckets? Recall that 1 byte = 8 bits.

- c) Let's analyze the time improvement as well. Let's say an average Chrome user attempts to visit 36,500 URLs in a year, only 1,000 of which are actually malicious. Suppose it takes half a second for Chrome to make a call to the database (the [Set](#)), and only 1 millisecond for Chrome to check containment in the Bloom filter. Suppose the false positive rate on the Bloom filter is 4%; that is, if a website is not malicious, the Bloom filter will incorrectly report it as malicious with probability 0.04. What is the time (in seconds) taken if we only use the database, and what is the *expected* time taken (in seconds) to check all 36,500 strings if we used the Bloom filter + database combination described earlier?