

# Asmt 1: Experimenting with Statistical Principles

---

Turn in (a pdf) through Canvas by 5pm:  
Wednesday, January 28

## Overview

In this assignment you will experiment with random variation over discrete events.

It will be very helpful to use the analytical results and the experimental results to help verify the other is correct. If they do not align, you are probably doing something wrong (this is a very powerful and important thing to do whenever working with real data).

*As usual, it is highly recommended that you use LaTeX for this assignment. If you do not, you may lose points if your assignment is difficult to read or hard to follow. Find a sample form in this directory: <http://www.cs.utah.edu/~jeffp/teaching/latex/>*

## 1 Birthday Paradox (30 points)

Consider a domain of size  $n = 3000$ .

**A: (5 points)** Generate random numbers in the domain  $[n]$  until two have the same value. How many random trials did this take? We will use  $k$  to represent this value.

**B: (10 points)** Repeat the experiment  $m = 300$  times, and record for each how many random trials this took. Plot this data as a *cumulative density plot* where the  $x$ -axis records the number of trials required  $k$ , and the  $y$ -axis records the fraction of experiments that succeeded (a collision) after  $k$  trials. The plot should show a curve that starts at a  $y$  value of 0, and increases as  $k$  increases, and eventually reaches a  $y$  value of 1.

**C: (5 points)** Empirically estimate the expected the number of  $k$  random trials in order to have a collision. That is, add up all values  $k$ , and divide by  $m$ .

**D: (10 points)** Describe how you implemented this experiment and how long it took for  $m = 300$  trials.

Show a plot of the run time as you gradually increase the parameters  $n$  and  $m$ . (For at least 3 fixed values of  $m$  between 300 and 10,000, plot the time as a function of  $n$ .) You should be able to reach values of  $n = 1,000,000$  and  $m = 10,000$ .

## 2 Coupon Collectors (30 points)

Consider a domain  $[n]$  of size  $n = 100$ .

**A: (5 points)** Generate random numbers in the domain  $[n]$  until every value  $i \in [n]$  has had one random number equal to  $i$ . How many random trials did this take? We will use  $k$  to represent this value.

**B: (10 points)** Repeat step A for  $m = 300$  times, and for each repetition record the value  $k$  of how many random trials we required to collect all values  $i \in [n]$ . Make a cumulative density plot as in 1.B.

**C: (5 points)** Use the above results to calculate the empirical expected value of  $k$ .

**D: (10 points)** Describe how you implemented this experiment and how long it took for  $n = 100$  and  $m = 300$  trials.

Show a plot of the run time as you gradually increase the parameters  $n$  and  $m$ . (For at least 3 fixed values of  $m$  between 300 and 5,000, plot the time as a function of  $n$ .) You should be able to reach  $n = 20,000$  and  $m = 5,000$ .

### 3 Comparing Experiments to Analysis (24 points)

**A: (12 points)** Calculate analytically (using formulas from the notes) the number of random trials needed so there is a collision with probability at least 0.5 when the domain size is  $n = 3000$ . (Show your work.)

How does this compare to your results from Q1.C?

**B: (12 points)** Calculate analytically (using formulas from the notes) the expected number of random trials before all elements are witnessed in a domain of size  $n = 100$ ? (Show your work.)

How does this compare to your results from Q2.C?

### 4 Random Numbers (16 points)

Consider when the only random function you have is one that chooses a bit at random. In particular `rand-bit()` returns 0 or 1 at uniformly random.

**A: (6 points)** How can you use this to create a random integer number between 1 and  $n = 1024$ ?

**B: (5 points)** Describe a Las Vegas randomized algorithm (“Las Vegas” means: it may run for ever, but you can bound how long you expect it to take, and when it finishes you know it is done) for when  $n = 1000$ .

**C: (5 points)** How many calls does this take in terms of  $n$  (say I were to increase the value  $n$ , how does the number of calls to `rand-bit()` change)?

Keep in mind that in many settings generating a random bit is much more expensive than a standard CPU operation (like an addition, if statement, or a memory reference), so it is critical to minimize them.

### 5 BONUS (2 points)

Consider a domain size  $n$  and let  $k$  be the number of random trials run. Let  $f_i$  denote the number of trials that have value  $i$ . Note that for each  $i \in [n]$  we have  $\mathbf{E}[f_i] = k/n$ . Let  $\mu = \max_{i \in [n]} f_i/k$ .

Consider some parameter  $\delta \in (0, 1)$ . How large does  $k$  need to be for  $\Pr[|\mu - 1/n| \geq 0.01] \leq \delta$ ? That is, how large does  $k$  need to be for *all* counts to be within 1% of the average with probability  $\delta$ ?

How does this change if we want  $\Pr[|\mu - 1/n| \geq 0.001] \leq \delta$  (for 0.1% accuracy)?

*(Make sure to show your work)*