# Semester Project

## Harrison Zheng

harrison.zheng1@marist.edu

December 15, 2021

## 1 Assignment Description

The semester project asked students to develop a simulation of group/pooled testing. The testing protocol works by dividing a population into smaller groups and collecting a sample (i.e. saliva) from each person in a group. The samples are then combined into one and tested as one. If the test result is negative, then no one in the group is infected. However, if the test result is positive, then the group is divided into smaller groups with a test performed on each subgroup and, if necessary, every individual in each subgroup until it can be determined who the infected person(s) is. If the infection rate in the population is low enough, this protocol can significantly reduce the number of tests needed to keep the population safe. Marist College utilized group testing throughout the 2020-2021 school year to track and limit the spread of COVID-19. For their projects, students ran their testing simulation on populations of 1,000, 10,000, 100,000, and 1,000,000 people with an infection rate of 2% and put them into groups of 8. Then, they needed to output the results of their simulation.

## 2 My Implementation

In the first run of my simulation, I start with a population size of 1000 and I create an instance of the Population class with this number (Listing 2.1). The Population class holds an array of instances of the Person class (Listing 2.3). An instance of the Person class represents a person that will be screened for COVID-19 and currently, it only describes whether the person is infected (Listing 2.2). The program then calls the infectPeople function, which infects a number of people at the front of the population array and then uses the Knuth shuffle from previous assignments to spread out the infected people. This ensures that multiple groups will have at least one infected individual when we do testing. Next, the program calls the testPeople function, which tests the population in groups of 8 at a time (case 1). If it finds an infection in the first test, it divides the group into two subgroups of 4 and performs a test on each group (case 2). If it finds that only one subgroup has an infected individual, then it only individually tests the people in that subgroup. However, if it finds that both subgroups have an infected individual, it individually tests all 8 people (case 3). The program keeps track of the number of tests performed and the number of occurences of each case as it goes through the testPeople function. Finally, the program prints out the total number of tests needed to screen the population and some other information. It repeats these steps for population sizes 10,000, 100,000, and 1,000,000.

```
1  import java.text.DecimalFormat;
2
3  public class Simulation {
4      public static int POPULATION_SIZE = 0;
5      public static int GROUP_SIZE = 8;
6      public static int GROUP_COUNT = 0;
7      public static double INFECTION_RATE = .02;
8      public static int testCount = 0;
9      public static int case1Count = 0;
10     public static int case2Count = 0;
11     public static int case3Count = 0;
12
13     public static void main(String[] args) {
14         //Run my pooled testing simulation on populations
15         //of various sizes ranging from 1000-1M people
16         runSimulation(1000);
17         runSimulation(10000);
18         runSimulation(100000);
19         runSimulation(1000000);
20     }
21
22     //Run simulation with given population size
23     public static void runSimulation(int popSize) {
24         POPULATION_SIZE = popSize;
25
26         Population simulatedPop = new Population(POPULATION_SIZE);
27         simulatedPop.wholePop = infectPeople(simulatedPop.wholePop);
28         testPeople(simulatedPop.wholePop);
29         printOutput();
30
31         testCount = 0;
32         case1Count = 0;
33         case2Count = 0;
34         case3Count = 0;
35     }
36
37     //Infect people in the population based on the infection rate
38     public static Person[] infectPeople(Person[] population) {
39         //Infect the first 20 people in the population array
40         int peopleToInfect = (int)Math.round(INFECTION_RATE * population.length);
41         for (int i = 0; i < peopleToInfect; i++) {
42             population[i].isInfected = 1;
43         }
44
45         //Randomly distribute infected people throughout
46         //the population with a Knuth shuffle
47         for (int i = 0; i < population.length; i++) {
48             //Create a "random" value
49             int rand = i + (int) (Math.random() * (population.length - i));
50
```

```
51          //Swap the value at this index with the value at
52          //the index of the randomly generated number
53          Person temp = population[rand];
54          population[rand] = population[i];
55          population[i] = temp;
56      }
57      return population;
58  }
59
60  public static void testPeople(Person[] population) {
61      //Place people into groups of 8 (125 groups in our case)
62      GROUP_COUNT = POPULATION_SIZE / GROUP_SIZE;
63
64      //Perform necessary tests on each group
65      for (int i = 0; i < GROUP_COUNT; i++) {
66          //Create a new testing pool
67          Person[] currentPool = new Person[GROUP_SIZE];
68
69          for (int j = i * GROUP_SIZE; j < ((i * GROUP_SIZE) + GROUP_SIZE); j++) {
70              currentPool[j%GROUP_SIZE] = population[j];
71          }
72
73          //Test the intial group of 8
74          testCount++;
75          if (checkForInfection(currentPool)) {
76              //If an infection was found,
77              //divide pool into two equal size groups
78              Person[] newPool1 = new Person[GROUP_SIZE / 2];
79              for (int j = 0; j < (currentPool.length / 2); j++) {
80                  newPool1[j] = currentPool[j];
81              }
82              Person[] newPool2 = new Person[GROUP_SIZE / 2];
83              for (int j = (currentPool.length / 2); j < currentPool.length; j++) {
84                  newPool2[j%(GROUP_SIZE / 2)] = currentPool[j];
85              }
86
87              //Test the two subgroups
88              if ((checkForInfection(newPool1) && !checkForInfection(newPool2)) || (!
                    ↪ checkForInfection(newPool1) && checkForInfection(newPool2))) {
89                  //Case 2: If only one subgroup showed infection,
90                  //then we used 6 more tests (1 per subgroup + 4 for people in infected
                        ↪ group)
91                  testCount += 6;
92                  case2Count++;
93              }
94              else {
95                  //Case 3: If both subgroups showed infection,
96                  //then we used 10 more tests (1 per subgroup + 8 for people in infected
                        ↪ groups)
97                  testCount += 10;
98                  case3Count++;
99              }
```

```java
100              }
101              else {
102                  //Only 1 test was needed because nobody was infected
103                  case1Count++;
104              }
105          }
106      }
107
108      public static boolean checkForInfection(Person[] pool) {
109          int giantTube = 0;
110
111          //Combine the isInfected value of everyone in the pool
112          //as you would combine the biomarkers in real life
113          for (int i = 0; i < pool.length; i++) {
114              giantTube += pool[i].isInfected;
115          }
116
117          //If the combined value is still 0, then nobody was infected
118          if (giantTube == 0) {
119              return false;
120          }
121          //If the combined value is > 0, then at least one person was infected
122          else {
123              return true;
124          }
125      }
126
127      public static void printOutput() {
128          DecimalFormat df = new DecimalFormat("#.##");
129          //Print the expected instances and tests
130          System.out.println("Expected values for a population of size " + POPULATION_SIZE);
131          System.out.println("Case (1): " + GROUP_COUNT + "x0.8500 = " + Double.valueOf(df.
                ↪ format(GROUP_COUNT * 0.8500)) + " instances requiring " + Math.round(
                ↪ GROUP_COUNT * 0.8500) + " tests");
132          System.out.println("Case (2): " + GROUP_COUNT + " x 0.1496 = " + Double.valueOf(df
                ↪ .format(GROUP_COUNT * 0.1496)) + " instances requiring " + Math.round((
                ↪ GROUP_COUNT * 0.1496) * 7) + " tests");
133          System.out.println("Case (3): " + GROUP_COUNT + " x 0.0004 = " + Double.valueOf(df
                ↪ .format(GROUP_COUNT * 0.0004)) + " instances requiring " + Math.round((
                ↪ GROUP_COUNT * 0.0004) * 11) + " tests");
134          System.out.println
                ↪ ("----------------------------------------------------------------------------")
                ↪ ;
135          int expectedTests = (int) (Math.round(GROUP_COUNT * 0.8500) + Math.round((
                ↪ GROUP_COUNT * 0.1496) * 7) + Math.round( (GROUP_COUNT * 0.0004) * 11));
136          System.out.println(expectedTests + " tests are expected to screen a population of
                ↪ " + POPULATION_SIZE + " people for a disease with 2% infection rate.\n");
137
138          //Print the actual numbers generated by my simulation
139          System.out.println("Actual results for a population of size " + POPULATION_SIZE);
140          System.out.println("Case (1): " + case1Count + " instances requiring " +
                ↪ case1Count + " tests");
```

```
141        System.out.println("Case (2): " + case2Count + " instances requiring " + Math.
              ↪ round(case2Count * 7) + " tests");
142        System.out.println("Case (3): " + case3Count + " instances requiring " + Math.
              ↪ round(case3Count * 11) + " tests");
143        System.out.println
              ↪ ("--------------------------------------------------------------------------")
              ↪ ;
144        System.out.println(testCount + " tests were actually used to screen a population
              ↪ of " + POPULATION_SIZE + " people for a disease with 2% infection rate.\n");
145    }
146 }
```

## 2.2 Person Class

```
1 public class Person {
2     int isInfected;
3
4     //Initialize each Person as not infected
5     public Person() {
6         this.isInfected = 0;
7     }
8 }
```

## 2.3 Population Class

```
1 public class Population {
2     //Will consist of 0s and 1s
3     Person[] wholePop;
4
5     //Initialize population array
6     public Population(int populationSize) {
7         this.wholePop = new Person[populationSize];
8         for (int i = 0; i < populationSize; i++) {
9             this.wholePop[i] = new Person();
10        }
11    }
12 }
```

# 3 Results of Simulation on Various Population Sizes

The results for all runs of the simulation are shown below. The decimal numbers used to calculate the expected instances of each case are based on the likelihood of each scenario occurring with a 2% infection rate and they come from Dr. Labouseur's article on pooled testing. The results suggest that there will most often be a difference in the number of tests we expect to use and the number that we actually use due to randomizing the infections. The actual and expected results are much closer for populations of 1,000 and 10,000 than they are for populations of 100,000 and 1,000,000 since the randomness of infections will have a greater effect as the population gets bigger. One should also note that each time population size increases by a factor of 10, so does the total number of tests needed (approximately).

## 3.1 Results for a Population of 1,000

**Expected Results**
Case (1): 125 x 0.8500 = 106.25 instances requiring 106 tests
Case (2): 125 x 0.1496 = 18.7 instances requiring 131 tests
Case (3): 125 x 0.0004 = 0.05 instances requiring 1 tests

————————————————————————————————————————————————————

238 tests are expected to screen a population of 1000 people for a disease with 2% infection rate.


**Actual Results**
Case (1): 105 instances requiring 105 tests
Case (2): 20 instances requiring 140 tests
Case (3): 0 instances requiring 0 tests

————————————————————————————————————————————————————

245 tests were actually used to screen a population of 1000 people for a disease with 2% infection rate.


## 3.2 Results for a Population of 10,000:

**Expected Results**
Case (1): 1250 x 0.8500 = 1062.5 instances requiring 1063 tests
Case (2): 1250 x 0.1496 = 187.0 instances requiring 1309 tests
Case (3): 1250 x 0.0004 = 0.5 instances requiring 6 tests

————————————————————————————————————————————————————

2378 tests are expected to screen a population of 10000 people for a disease with 2% infection rate.


**Actual Results**
Case (1): 1069 instances requiring 1069 tests
Case (2): 169 instances requiring 1183 tests
Case (3): 12 instances requiring 132 tests

————————————————————————————————————————————————————

2384 tests were actually used to screen a population of 10000 people for a disease with 2% infection rate.


## 3.3 Results for a Population of 100,000

**Expected Results**
Case (1): 12500 x 0.8500 = 10625.0 instances requiring 10625 tests
Case (2): 12500 x 0.1496 = 1870.0 instances requiring 13090 tests
Case (3): 12500 x 0.0004 = 5.0 instances requiring 55 tests

————————————————————————————————————————————————————

23770 tests are expected to screen a population of 100000 people for a disease with 2% infection rate.


**Actual Results**
Case (1): 10627 instances requiring 10627 tests
Case (2): 1801 instances requiring 12607 tests
Case (3): 72 instances requiring 792 tests

————————————————————————————————————————————————————

24026 tests were actually used to screen a population of 100000 people for a disease with 2% infection rate.

## 3.4 Results for a Population of 1,000,000

**Expected Results**
Case (1): 125000 x 0.8500 = 106250.0 instances requiring 106250 tests
Case (2): 125000 x 0.1496 = 18700.0 instances requiring 130900 tests
Case (3): 125000 x 0.0004 = 50.0 instances requiring 550 tests

————————————————————————————————————————————————————

237700 tests are expected to screen a population of 1000000 people for a disease with 2% infection rate.


**Actual Results**
Case (1): 106334 instances requiring 106334 tests
Case (2): 17924 instances requiring 125468 tests
Case (3): 742 instances requiring 8162 tests

————————————————————————————————————————————————————

239964 tests were actually used to screen a population of 1000000 people for a disease with 2% infection rate.


# 4 Binomial vs. Hypergeometric Distribution

In the field of statistics, a binomial distribution is a probability distribution of the number of "successes" over a sequence of trials. To calculate the binomial probability of successes at a certain point in the distribution, use the formula

$$P(x) = \binom{n}{x} p^x q^{n-x}$$

where x is the total number of successes, $p^x$ is the probability of a success on a single trial, $q^{n-x}$ is the probability of a failure on a single trial, and n is the number of trials.

A hypergeometric distribution is very similar to a binomial distribution, except it describes the number of successes over a sequence of trials from a finite population without replacement. This means every time a trial is performed on an individual from the population, that individual is removed from the population for the next trial. To calculate the hypergeometric probability of successes, use the formula

$$P(x) = \frac{\binom{k}{n}\binom{N-k}{n-x}}{\binom{N}{n}}$$

where x is the number of successes, k is the number of failures, n is the sample size, and N is the total number of individuals in the population.

Although we used binomial probabilities to calculate the expected number of instances of each case in our simulations, hypergeometric probabilities actually would have provided us with better estimates, especially for the larger populations. The table in section 4.1 shows the results I found using hypergeometric probabilities. For each population size, I took the number of testing groups (ex: 1250) and multiplied it by the hypergeometric probability (ex: 0.8507) of finding $x$ infected people in a group (x is 0 for case 1, 1 for case 2, and 2+ for case 3), which I found using the formula above. This produced estimates of the number of instances we would get for each case in a simulation and upon comparison, I found that many of them are closer to the actual results from the simulation runs than our original estimates. For example, using hypergeometric probabilities produced an estimate of 12.87 instances for case 3 in a population of 10,000, which is a lot closer to the actual result of 12 than our original estimate of 0.5. This is likely due to the fact that a binomial distribution has a fixed number of independent trials and thus does not model our simulations as well as a hypergeometric distribution, which has a set number of dependent trials. When the simulation is testing a group and it finds that somebody is infected, it divides the group into 2 subgroups and then

conducts a pooled test on each one. If it finds that there is an infection in one subgroup but not the other, it removes that half of the group from any further testing. Additionally, once the simulation is done testing one group, it moves on to the next eight individuals in the array that holds the population and does not go back so any individuals that have been tested are effectively removed from the screening population. This way of sampling without replacing the contents of the sample makes a hypergeometric distribution the best way to model the simulation. When the population is small (i.e. 1000 people), which distribution we use to make estimates does not matter much as the numbers come out roughly the same. As the population gets bigger, though, the effect of removing individuals from the screening population after each pooled test becomes more pronounced and estimates from a hypergeometric distribution become more accurate than those from a binomial distribution.

## 4.1 Expected Results Based On Hypergeometric Distribution

|          | 1,000                     | 10,000                      | 100,000                      | 1,000,000                      |
|----------|---------------------------|-----------------------------|------------------------------|--------------------------------|
| Case (1) | 125 x 0.8503 = 106.28     | 1250 x 0.8507 = 1063.29     | 12500 x 0.8508 = 10634.48    | 125000 x 0.8508 = 106345.32    |
| Case (2) | 125 x 0.1398 = 17.48      | 1250 x 0.1390 = 173.74      | 12500 x 0.1390 = 1736.37     | 125000 x 0.1390 = 17362.62     |
| Case (3) | 125 x 0.0099 = 1.24       | 1250 x 0.0103 = 12.87       | 12500 x 0.0103 = 129.16      | 125000 x 0.0103 = 1292.06      |

## 5 Potential Ways to Improve Simulation

One way that I would improve my simulation is adding a mechanism that creates false positives and false negatives in the test results. For our simulations, we were asked to assume that the tests are 100% accurate. In reality, however, no test used to screen for a disease is 100% accurate. I would implement testing inaccuracies in my simulation by having it give each person in the population two attributes; their test result and their actual infection status. First, it would change 2% of the population's infection status to 1, as usual. Then during the testing phase, in random fashion, it would give the majority of the healthy people a test result of negative and a small percentage (based on real-world false positive rates) of them a test result of positive to serve as false positives. It would also give, in random fashion, most infected people a test result of positive and a small amount (based on real-world false negative rates) of them a test result of negative to serve as false negatives. After testing, my program would go through the population and compare each person's test result with their actual infection status while keeping track of disparities. This is similar to how testing centers need to perform a second test on people to confirm whether their previous result was erroneous. My program would then print out the number of false positives and false negatives in the tests along with the other results from the simulation.

Another way that I would improve my simulation is by making it so that the user can enter values for population size, group size, and infection rate when the program starts up. The program would then check to make sure that the user entered values that are within an acceptable range and run the simulation as well as output results based on those values. Allowing for user input would make my simulation more dynamic and useful for approximating test results in a real-world community where it is unlikely that the population size is an exact power of 10.