# Simulating the Luria-Delbrück experiment in `R`

## Jun Ishigohoka

```r
1  T <- 14 # Number of generations
2  n_0 <- 100 # Number of cells to in a tube
3  mu <- 1e-3 # Mutation rate
4  r <- 50 # Number of plates (A) or tubes (B)
5  n_sample <- 10000 # Number of cells to plate
```

## Implementation of the spontaneous mutation model

Let's write steps in R.

### Step 1. Simulation of passaging cells in a tube from a flask

Let `parents` a vector of $n_0$ zeros, representing $n_0$ wild type cells. For now, we (as in the original paper) assume that there are no standing genetic variation in these cells. Later you are free to relax this assumption and investigate how the presence of standing genetic variation in this step might affect the result.

```r
1  parents <- rep(0, n_0)
```

### Step 2. Growth

Let cells replicate $T$ times and mutate. Try with small $T$

```r
1  T_test <- 5
2  # Loop over T generations
3  for(t in 1:T_test){
4          # Initialise a vector children, into which we put genotype of daughter cells
5          children <- c()
6          # Loop over parent cells
7          # In each iteration, `cell` is the genotype of focal cell
8          for(cell in parents){
9                  # Loop over two daughter cells
10                 for(i in 1:2){
11                         # Bernoulli sampling of muation
12                         # Mutation is represented as addition of 1
13                         # The genotype of the focal daughter cell is appended
14                         children <- c(children, cell + rbinom(1, 1, mu))
15                 }
16         }
17         # vector parents is updated
18         parents <- children
19 }
20
21 # frequency distribution of cells freq(genotype == x)
22 table(parents)
```

```
## parents
##    0    1
## 3189   11
```

Entries greater than 1 are converted into 1 because cells whose ancestors experienced at least one muation are resistant.

```
1  parents <- as.integer(parents > 0)
2  table(parents)
```

```
## parents
##    0    1
## 3189   11
```

## Initial optimisation of steps 1 and 2

The `for` loop in lines 7-15 in the above code block can be substituted with a single binomial sampling of `n = 2 * length(parents)` ($= n_0 \times 2^T$). We can directly update `parents` instead of making an intermediate vector `children`. And without children naming of parents wouldn't make sense, so let us call the vector `genotypes`.

```
1  genotypes <- rep(0, n_0)
2
3  for(t in 1:T){
4        genotypes <- rep(genotypes, 2) + rbinom(n_0 * 2^t, 1,  mu)
5  }
6
7  genotypes <- as.integer(genotypes > 0)
8  table(genotypes)
```

```
## genotypes
##       0       1
## 1610230   28170
```

The number of resistant cells after $T$ generations is the sum of the vector.

```
1  sum(genotypes)
```

```
## [1] 28170
```

## Make a function `sim_tube` for steps 1 and 2

Let's make a function for steps 1 and 2.

```
1  sim_tube <- function(n_gens = T, mut_rate = mu, ncells_init = n_0){
2        genotypes <- rep(0, ncells_init)
3        for(t in 1:n_gens){
4              genotypes <- rep(genotypes, 2) + rbinom(ncells_init * 2^t, 1,  mut_rate)
5        }
6        genotypes <- as.integer(genotypes > 0)
7        return(genotypes)
8  }
```

Let's simulate a tube and count the number of resistant cells

```
1  tube <- sim_tube(T, mu, n_0)
2  sum(tube)
```
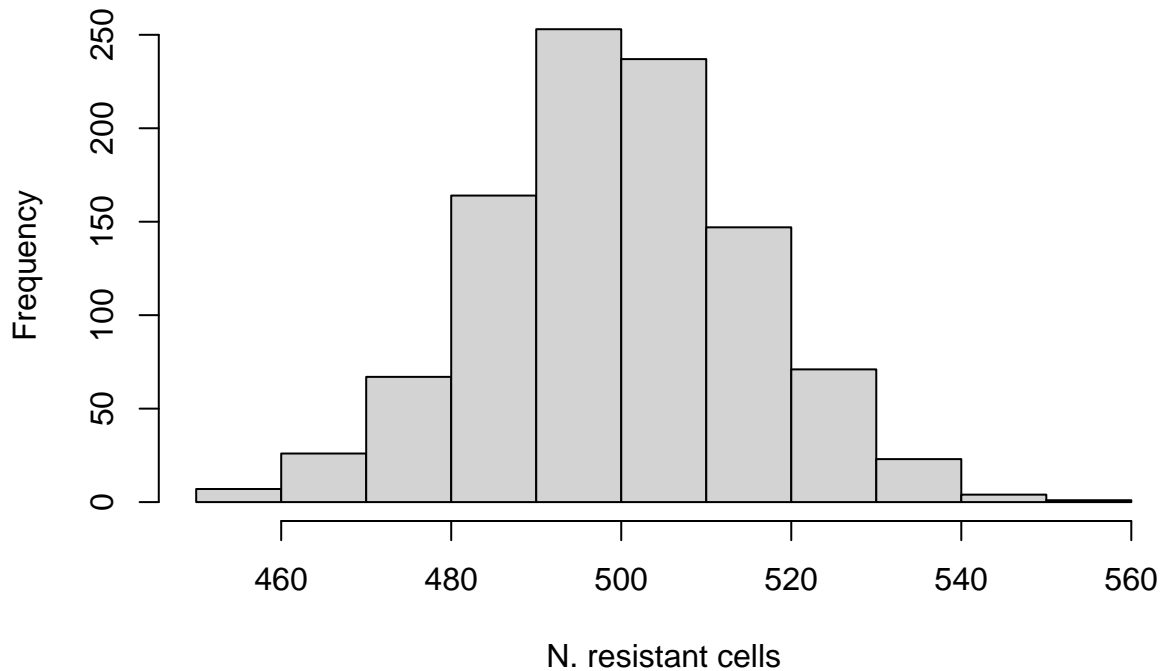
```
## [1] 26396
```

Let's see if this function actually does a right job. If we let $T = 1$, then the number of resistant cells should be around $2n_0\mu$. So, when $\mu = 0.5$ and $n_0 = 500$, there should be around 500 resistant cells. Let's simulate it 1,000 times and check the distribution of the number of resistant cells.

```r
hist(
    sapply(1:1000,
        function(x){
            return(sum(sim_tube(1, 0.5, 500)))
        }
    ),
    xlab = "N. resistant cells",
    main = ""
)
```
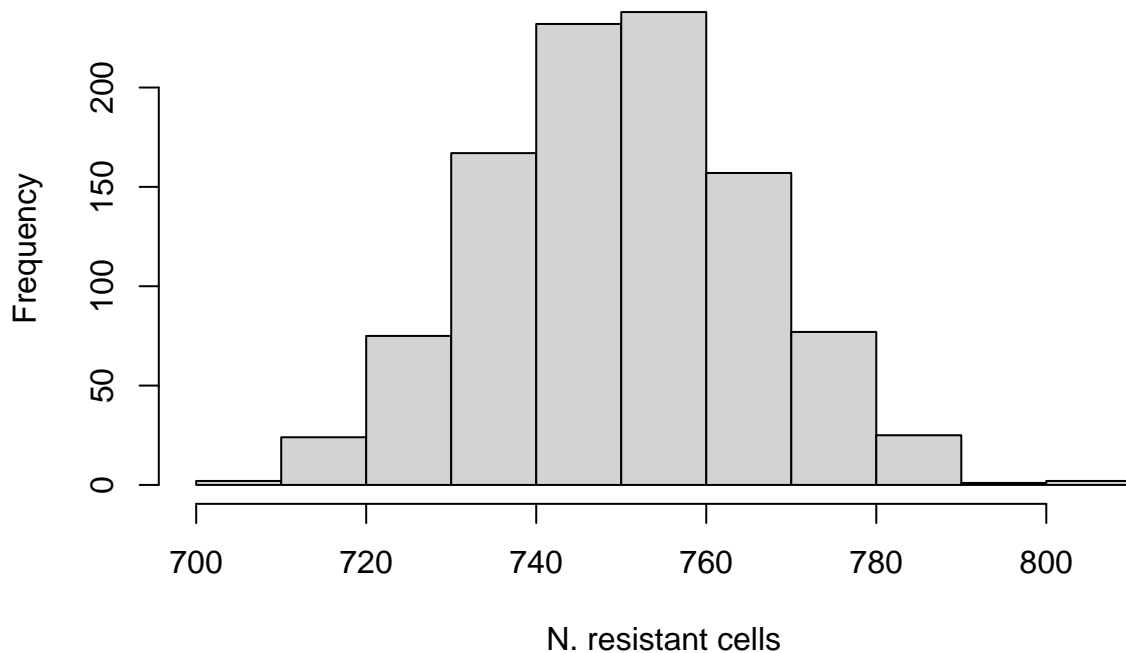


If we let $T = 2$, then the number of resistant cells should be around $n_0 \times 2^2 \times (1 - \mu^2)$. So, when $\mu = 0.5$ and $n_0 = 250$, there should be around 750 resistant cells. Let's simulate it 1,000 times and check the distribution of the number of resistant cells.

```r
hist(
    sapply(1:1000,
        function(x){
            return(sum(sim_tube(2, 0.5, 250)))
        }
    ),
    xlab = "N. resistant cells",
    main = ""
)
```

The simulator `sim_tube` seems to be working properly.

### Simulate resistant cells in experiment A

In experiment A, $n_0$ cells are let grow over $T$ generations (to $n_0 \times 2^T$ cells). The cells are plated to $r$ plates from this tube. So, we will simulate one tube before plating. The genotypes of cells are recorded in a vector `tube_a` of length $n_0 \times 2^T$.

```
tube_a <- sim_tube(n_gens = T, mut_rate = mu, ncells_init = n_0)
head(tube_a)
```

```
## [1] 0 0 0 0 0 0
```

```
length(tube_a)
```

```
## [1] 1638400
```

### Simulate resistant cells in experiment B

In experiment B, there are $r$ tubes, each of which accommodates $n_0$ cells and let them grow over $T$ generations. So, we will simulate $r$ tubes. The genotypes are recorded in a matrix `tubes_b` with $r$ rows and $n_0 \times 2^T$ columns.

```
tubes_b <- t(
            sapply(1:r,
                function(x){
                        sim_tube(n_gens = T, mut_rate = mu, ncells_init = n_0)
                }
            )
)
```

## Step 3. Simulation of plating

In experiment A, we plate $n_{sample}$ cells from 1 tube to $r$ plates. In experiment B, we plate $n_{sample}$ cells from each of $r$ tubes to each of $r$ plates, respectively.

4

**Experiment A**

In experiment A, given a vector of genotypes of cells in a tube, we want a matrix with $r$ rows and $n_{sample}$ columns representing genotypes of cells. One can shuffle the input vector representing genotypes of cells in a tube and take the first $n_sample$ cells for the first plate, second $n_sample$ cells for the second plate, and so on.

Let's start with plating 2 cells from 6 cells in a tube. Genotypes of the plated 2 cells can be stored in a vector.

```
1  n_plates <- 1
2  n_cells_plate <- 2
3
4  tube_test <- c(1, 1, 0, 0, 0, 0)
5  # shuffle
6  tube_test <- sample(tube)
7  plate_test <-tube_test[1:n_cells_plate]
```

Then let's try plating 2 cells from 10 cells in each of 3 plates. The output genotypes are stored in a matrix with 3 rows and 2 columns.

```
1  n_plates <- 3
2  n_cells_plate <- 2
3
4  tube_test <- c(1, 1, 0, 0, 0, 0, 1, 0, 1, 0)
5  # shuffle
6  tube_test <- sample(tube)
7  # A vector to store cells to plate in n_plates plates
8  plates_test <-tube_test[1:(n_cells_plate * n_plates)]
9  # Make it into a matrix
10 plates_test <- matrix(plates_test, nrow = n_plates, byrow = T)
```

Based on the above exercise we can write a function `sim_plate`.

```
1  sim_plate <- function(tube, n_plates, ncells_plate){
2         tube_shuf <- sample(tube)
3         plates <-tube_shuf[1:(ncells_plate * n_plates)]
4         plates <- matrix(plates, nrow = n_plates, byrow = T)
5         return(plates)
6  }
```

By running this function on our simulated `tube_a`, we can obtain `plates_a`, a matrix of genotypes.

```
1  plates_a <- sim_plate(tube = tube_a, n_plates = r, ncells_plate = n_sample)
```

**Experiment B**

In experiment B, given a matrix with $r$ rows representing genotypes of cells in $r$ tube, we want a matrix with $r$ rows and $n_{sample}$ columns representing genotypes of cells. Conveniently, for each tube we can use `sim_plate` function with `n_plates = 1`.

For example, for the first tube,

```
1  plate_b_1 <- sim_plate(tube = tubes_b[1,], n_plates = 1, ncells_plate = n_sample)
```

We can apply this to each row of matrix `tubes_b` (and transpose the output) to obtain a matrix of genotypes `plates_b`.

```
1  plates_b <- t(
2              sapply(1:r,
3                  function(x){
4                      return(sim_plate(tube = tubes_b[x,],
```

```
5                                            n_plates = 1,
6                                            ncells_plate = n_sample))
7                       }
8                 )
9  )
```

## Step 4. Compute mean and variance

Both `plates_a` and `plates_b` are matrices whose row corresponds to a replicate plate. The entry of the matrix is genotype, where 0 is wild type and 1 is resistant mutant. Therefore the sum of each row is the number of resistant colonies in each plate. We can summarise the results in two values: mean and variance of the number of colonies per plate.

```
1  result <- c(mean_a = mean(rowSums(plates_a)),
2             mean_b = mean(rowSums(plates_b)),
3             var_a = var(rowSums(plates_a)),
4             var_b = var(rowSums(plates_b)))
5  result
```

```
##     mean_a    mean_b     var_a     var_b
##   171.1800  142.7000  190.9261 1263.0714
```

# Implementation of induced mutation model

In both experiments A and B, cells should mutate at the same rate (mutation rate) after plating. So, this is binomial sampling with rate parameter of $\mu$ and size parameter of $n_sample$

```
1  sim_ld_ind <- function(n_plates, n_sample, mu){
2        plates_a <- rbinom(n_plates, n_sample, mu)
3        plates_b <- rbinom(n_plates, n_sample, mu)
4        res <- c(mean_a = mean(plates_a),
5               mean_b = mean(plates_b),
6               var_a = var(plates_a),
7               var_b = var(plates_b)
8        )
9        return(res)
10 }
```

```
1  sim_ld_ind(n_plates = r, n_sample = n_sample, mu = mu)
```

```
##     mean_a    mean_b     var_a     var_b
##   9.920000  9.480000  7.585306 12.458776
```

We can run it multiple times using `sapply`

```
1  t(
2    sapply(1:5,
3         function(x){
4               return(sim_ld_ind(r, n_sample, mu))
5         }
6    )
7  )
```

```
##      mean_a mean_b    var_a     var_b
## [1,]   9.30  10.02  8.091837 12.958776
## [2,]   9.56  10.02  7.394286  9.775102
```

```
## [3,]  10.00  10.12 13.510204  8.720000
## [4,]   8.98   9.94  8.264898  9.567755
## [5,]  10.20  10.44 10.693878 10.169796
```

# Questions

Let's address questions with our simulators. Note that our simulators are not very well optimised, and some questions may be still iplausible (e.g. large $T$, large $\mu$). We will further optimise the simulators in the next section, but you are already ready to play around with what you have to get an intuition of the experiment.

1. aaa
2. aaa
3. aaa

# Further optimisations of the spontaneous mutation simulator

In the above simulators `sim_tube` and `sim_plate`, we recorded the genotypes of all cells in tubes and plates until the final step. This is very inefficient for our purpose. The information of index of the genotype vector is not used: Even if someone shuffled our vector at any step of our simulation, we would not suffer. The only information we need is actually the numbers of wild type cells adn resistant cells, instead of genotype of millions of cells.

Below, we will try to improve the scripts to make it more scalable and faster.

## Step 1. Simulation of passaging cells in a tube from a flask

Instead of recording genotypes of $n_0$ cells, we can have two objects to keep the number of wild type and resistant cells.

```
1  n_wt <- n_0
2  n_re <- 0
```

## Step 2. Growth

In each generation, daughter cells of resistant cells are all resistant. Some of daughter cells of wild type cells mutate to resistant, and the number of such cells follow a binomial distribution with the size parameter of 2 x n. wild type parent cells.

```
1  for(t in 1:T){
2          n_re <- 2 * n_re + rbinom(1, 2 * n_wt, mu)
3          n_wt <- n_0 * 2^t - n_re
4  }
5
6  print(c(n_wt = n_wt, n_re = n_re))
```

```
##    n_wt    n_re
## 1619291   19109
```

## Make a function `sim_tube_count`

```
1  sim_tube_count <- function(n_gens = T, mut_rate = mu, ncells_init = n_0, ncells_res_init = 0){
2          n_wt <- ncells_init - ncells_res_init
3          n_re <- ncells_res_init
4          for(t in 1:n_gens){
5                  n_re <- 2 * n_re + rbinom(1, 2 * n_wt, mut_rate)
```

```
6              n_wt <- ncells_init * 2^t - n_re
7          }
8          return(n_re)
9  }
```

Note that I added an optional argument `ncells_res_init`, the number of resistant cells in the initial passage, reflecting standing variation (default: 0).

Now, this `sim_tube_count` returns an integer, the number of resistant cells. By using this function, we can simulate the number of resistant cells in experiment A.

```
1  tube_count_a <- sim_tube_count(n_gens = T, mut_rate = mu, ncells_init = n_0)
```

By applying this function, we can simulate the number of resistant cells in experiment B. The counts are stored in a vector `tubes_count_b`.

```
1  tubes_count_b <- sapply(1:r,
2                      function(x){
3                          sim_tube_count(n_gens = T, mut_rate = mu, ncells_init = n_0)
4                      }
5  )
```

## Step 3. Simulation of plating

Before plating, we have `n_re` resistant cells and `n_wt` wild type cells. In experiment A, we sample `n_sample` cells without replacement sequentially over r times. Each plating is equivalent to taking some balls from a box with some red and blue balls without replacement, and the number of balls with one colour follows a hypergeometric distribution.

```
1  sim_plate_count <- function(n_re, n_wt, n_plates, n_sample){
2          # Number of resistant cells in plates
3          plates_count <- c()
4          # Loop over n_plates plates
5          for(i in 1:n_plates){
6                  plates_count <- c(plates_count, rhyper(1, n_re, n_wt, n_sample))
7                  n_re <- n_re - plates_count[i]
8                  n_wt <- n_wt - (n_sample - plates_count[i])
9          }
10         return(plates_count)
11 }
```

For experiment A, we use this function to obtain the number of resistant cells on $r$ plates.

```
1  plates_count_a <- sim_plate_count(n_re = tube_count_a,
2                              n_wt = n_0 * 2^T - tube_count_a,
3                              n_plates = r,
4                              n_sample = n_sample
5  )
```

For experiment B, we can apply this function with `n_plates = 1` over $r$ times.

```
1  plates_count_b <- sapply(tubes_count_b,
2                      function(x){
3                          sim_plate_count(n_re = x,
4                                      n_wt = n_0 * 2^T - x,
5                                      n_plates = 1,
6                                      n_sample = n_sample
7                          )
```

```
8                                    }
9    )
```

## Step 4. Compute mean and variance

```
1    result <- c(mean_a = mean(plates_count_a),
2              mean_b = mean(plates_count_b),
3              var_a = var(plates_count_a),
4              var_b = var(plates_count_b))
5
6    result
```

```
##     mean_a    mean_b     var_a      var_b
##   151.6200   139.6000   109.7506  1255.8367
```

## Make a function to do all. . .

```
1    sim_ld_spo <- function(n_gens, mut_rate , ncells_init , n_sample , n_plates,  ncells_res_init = 0){
2            tube_count_a <- sim_tube_count(n_gens = n_gens, mut_rate = mut_rate, ncells_init = ncells_init)
3            tubes_count_b <- sapply(1:n_plates,
4                                function(x){
5                                        sim_tube_count(n_gens = n_gens, mut_rate = mut_rate, ncells_ini
6                                }
7            )
8            plates_count_a <- sim_plate_count(n_re = tube_count_a,
9                                            n_wt = ncells_init * 2^n_gens - tube_count_a,
10                                           n_plates = n_plates,
11                                           n_sample = n_sample
12           )
13           plates_count_b <- sapply(tubes_count_b,
14                                function(x){
15                                        sim_plate_count(n_re = x,
16                                                      n_wt = ncells_init * 2^n_gens - x,
17                                                      n_plates = 1,
18                                                      n_sample = n_sample
19                                        )
20                                }
21           )
22           result <- c(mean_a = mean(plates_count_a),
23                     mean_b = mean(plates_count_b),
24                     var_a = var(plates_count_a),
25                     var_b = var(plates_count_b))
26           return(result)
27   }
```

We can run this multiple times using `sapply`

```
1    t(sapply(1:5,
2          function(x){
3          sim_ld_spo(n_gens = 30,
4                  mut_rate = 1e-4,
5                  ncells_init = n_0,
6                  n_sample = n_sample,
7                  n_plates = r,
```

```
8                    ncells_res_init = 0)
9            }
10   )
11   )
```

```
##       mean_a mean_b    var_a     var_b
## [1,]  27.54  27.92 26.94735  72.40163
## [2,]  29.96  29.74 25.10041 137.42082
## [3,]  24.04  30.78 30.40653 145.23633
## [4,]  28.36  28.44 27.58204  94.82286
## [5,]  26.18  29.44 19.17102 147.51673
```

## Solutions to some of the questions

```
1    T <- 15
2    mu <- 1e-4
3
4
5    #for(i in 1:length(seq(0,1,0.2))){
6    #        p <- seq(0,1,0.2)[i]
7    #        t(sapply(1:10,
8    #                 function(x){
9    #                  sim_ld_spo(n_gens = T,
10   #                         mut_rate = mu,
11   #                         ncells_init = n_0,
12   #                         n_sample = (n_0 * 2^T)/100,
13   #                         n_plates = r,
14   #                         ncells_res_init = n_0 * p)
15   #                 }
16   #        )
17   #        )
18   #}
```