

# Task 05: Revisiting a Lab

## I. General

**Type up all of your work in a text editor.** Basically, you should NEVER type things directly into the R terminal. Type them into a text editor, then either run them or copy/paste them into R.

Before you begin, make a new folder in Tasks called Task.05, and save an empty file named task03.r in that folder.

When you're **done** with this assignment, turn it in by (1) saving your text document, (2) opening your Terminal or GitBash, (3) navigating to the appropriate directory using `cd` and (4) typing:

```
git add -A (enter)
git commit -m "yourname Task 04" (enter)
git push -u origin master (enter)
```

## II. Background

This R activity has three separate parts. First, you're going to do some simple simulation experiments using the `simPop` function. These are called *in silico* experiments, as they are done on a computer. You'll make a (virtual) population of little hermaphroditic creatures that mate and then die (non-overlapping generations; individuals can mate more than once though!). You'll track their fitnesses and look at how allele frequencies change as a function of population size, heritability of allele b, and selection acting on allele b.

Next, you'll get a touch more realistic. Instead of experimenting on simulated populations, you'll be reanalyzing data from a BIOL 112 lab in which students simulated evolution physically by "eating" asexually reproducing sticks. Unlike the simulation, you do not know the relative strength of selection here (as it was produced by students), and you cannot vary the population size (the experiment is over, these are just the data).

Finally, you'll combine the two approaches. You'll simulate additional trials of the student-led experiment, and compare your simulated results to the observed results in an attempt to learn something about how students behave during this lab.

## III. Simulation

We're going to be doing simulations in this handout. We touched on them briefly last week, but this week is going to be heavily driven by them. A simulation is when, on your computer, you invent a fake world with a set of rules that you decide on. Then you let that world go and grow for a period of time (defined in those fixed rules), and have it tell you how things went.

The `simPop` function we're going to use invents a world of hermaphrodites that have exactly two alleles and reproduce sexually, but whose population size never varies (that is, each generation

there are always the same number of individuals). This population progresses from some start point for some number of generations, and then you receive a record of the frequency of both the two alleles and all three possible genotypes for each generation.

You can control the simulation by setting the arguments in the `simPop` function to different values. Adjusting the `Popsiz` argument changes the population size, for instance, while the `nGenerations` argument sets how long the simulated world should run. `initial_p` controls the starting frequency of the a allele, while `h` and `s` control the heritability and selection coefficient for the b allele.

As with last week, this handout mostly just shows you the basic usage of `simPop`, but as you can see, you can use it to study. I strongly encourage you to do so. Look at how effective selection is when the allele is recessive versus overdominant, look at how drift and selection interact. Does the most fit allele always go to fixation? What combinations of drift and heritability prevent that? There are many nuances to our material that you can explore with these functions. I've also provided a convenient plotting tool.

```
source("http://jonsmitchell.com/code/fxn05.R")

# run a simple simulation
Pop1 <- simPop(Popsiz = 50, nGenerations = 100, initial_p = 0.5, h =
  1, s = 0)

plot(1:nrow(Pop1), Pop1[,1], ylim=c(0, 1), type = "l", xlab="generation
  ", ylab="allele freq.", lwd=2)
lines(1:nrow(Pop1), Pop1[,2], lwd=2, col='red')
legend("topleft", legend = c("a", "b"), col = c("black", "red"), lwd =
  2, bty="n")

# I also set up a plotting wrapper function. This works mostly the same
# as normal simPop
# except you can set the number of runs (nruns) to do multiple
# simulations at once.
# It will plot the relative fitnesses of the three genotypes (left) and
# the frequencies of
# allele a for each of the nruns populations with a different color per
# simulation

plotFit( nruns = 10, n = 50, ngens = 100, init_p = 0.5, h = 1, s = 0 )
```

## IV. Days Long Past

This week, we're going to return to a lab you did in the past if you took BIOL 112 here at Tech. The third lab in that class was a Natural Selection experiment involving toothpicks. If you didn't take 112 here, or if you just want a memory refresher, students formed groups and spread 60 colored toothpicks (10 toothpicks of each of 6 colors) onto a colored paper background. One student "ate" the toothpicks one at a time until only 20 toothpicks remained. Then, the 20 "surviving" toothpicks "reproduced" to create two additional copies of their color back into the pool, resulting in 60 toothpicks again. Predation and reproduction were then repeated twice. At the end, the  $\chi^2$  statistic was calculated ( $\frac{(expected-observed)^2}{expected} = \chi^2$ ).

For reference, the 112 lab handout is posted on eCampus for you to re-read!

The datasheet you'll be analyzing is the "whole class" data file from the 2012 to 2020 labs combined. Each row represents the final population size for each of the six colors, for one group, on one background, from one year. In that freshmen-level class, you learned a lot of foundational things (that you must remember for this class! It's a prerequisite!) but, as it is a freshmen-level class, there were some necessary simplifications. We're going to explore one of those by reanalyzing the data from that lab.

## V. What is a test statistic?

A *test statistic* is a way to quantify how different a set of observations is from your expectation. For instance, in a t-Test, you compare the difference in the means between two groups (this is the "t statistic"). You *expect* the t-statistic to be zero if there's no difference between groups. But random variation in samples means that the test statistic may not be exactly zero even if there is no meaningful difference between groups. So the further from zero your observed value is, *clearer it is* that the two groups are different.

Almost all "statistical tests" you know/at one point learned are based on this idea. You calculate some measurement of your data relative to some expectation (usually zero) and get a number that describes how far your observed data are from that expectation. Then you compare your calculated number to a known distribution (usually in the form of a table) to get the "p-value". The notion is that the p-value tells you the probability your data differ from your expectation due solely to sampling errors. However, this is essentially never true, as p-values are accurate if and only if all of the assumptions of your test statistic are met and your expectation isn't unwarranted.

If the above paragraph seems confusing (or if you've had a statistics class and found it confusing), don't fret: it seems confusing because it *is* confusing. The whole idea of using "standardized" test statistics only makes sense in a pre-computer world where you needed to look things up in published tables. Because you couldn't have an infinite number of published tables, people would only have a few (e.g., a table for the t, F, or  $\chi^2$ -statistics) and then would design experiments *to fit with the tables they had*.

You can visualize this in R

```
# Let's assume that we expect to see four equal categories
Expectation <- c(10, 10, 10, 10)

# Now let's pretend what we actually see is four uneven categories
Observed <- c(15, 15, 5, 5)

# We can calculate the Chi-squared statistic using these numbers
Chisq <- sum( ( ( Expectation - Observed ) ^ 2 ) / Expectation )

# And we can visualize what this means pretty readily
barplot(rbind(Expectation, Observed), beside = T, main = bquote(chi^2 ~
  "=" ~.(Chisq)), legend.text=c("expected", "observed"))
```

Try setting the Observed vector to c(5, 0, 0, 35), and (2, 3, 10, 30), and other combinations that you can think of. What's the  $\chi^2$  value when you observe all 10's? What about when all 40 observations are only in one category? How does the  $\chi^2$  relate to the *evenness* of the bars?

Nowadays, we have computers. We can calculate tables trivially, if we want, or we can simulate whole alternate realities in which experiments were conducted differently, or assumptions were

126 violated and compare our observations to those alternate worlds.

127 This week's task is a little funny. We're going to try to learn about evolutionary processes  
128 by comparing how actual freshmen perform a physical experiment to how hypothetical perfectly-  
129 behaved freshmen perform a simulated experiment. We'll do this by looking at calculated  $\chi^2$  values  
130 from the 112 lab and comparing them to simulations.

131 We're going to take the observed data. We're going to calculate the test statistics and ask  
132 how often we find "significant" results. Then we're going to simulate an alternative scenario and  
133 compare our observed data to the simulation. We'll use that comparison to ask whether finding  
134 a "significant" result really means what we think it means, and see if we can figure out what  
135 evolutionary processes are most responsible for the observed changes.

136 To get started, open an empty .R file, as always, set your working directory, and begin adding  
137 in the following code:

```
138 # Read in the data from past years. I'd encourage you to look at this
139 data. See Task 02a to remind yourself how to do that!
140 results <- read.csv("http://jonsmitchell.com/data/biol112labresults.csv",
141 stringsAsFactors=F)
142
143 # Simple subsetting of data to have only the right pieces
144 counts <- results[,c("yellow", "red", "green", "blue", "black", "tan")]
145
146 # Let's list out the background colors (they differ from the toothpicks
147 !)
148 backgrounds <- c("White", "Red", "Yellow", "Green", "Blue", "Black")
149
150 # Let's set some (slightly nicer) than default colors. This line is
151 optional.
152 backgroundCol <- c("white", "#d53e4f", "#fee08b", "#abdda4", "#3288bd", "
153 black")
154
155 # Now, let's calculate the Chi-squared statistic for the first row of
156 counts.
157 calcChi(counts[1,])
158
159 # Now let's calculate Chi-squared for ALL of the rows at once!
160 Chisqs <- apply(counts, 1, calcChi)
161
162 # You've all seen Chi-squared values before, but interpreting what
163 they actually mean, visualizing what the chi-squared number is
164 telling you, might be a bit tricky. Use the following function
165 multiple times.
166 plotChis(counts)
167
168 # Look at the plots. Look at the chi squared values for each panel, and
169 the heights of the bars. How "even" are the bars when the chi
170 squared is high? How "even" are they when it is low? What does
171 plotChis() show you about how to interpret the number a Chi-squared
172 test spits out?
```

```

173
174 # What's the average Chi-squared? How would you interpret that given
175   the (many) plots you should've just looked at?
176 Avg <- mean(Chisqs)
177
178 # How does the average chi-squared compare to the critical value in the
179   packet (11.70)?
180
181 # Does the chi-squared differ by background?
182 backgroundAvs <- tapply(Chisqs, results[,3], mean)
183
184 # The critical value here is 11.70. That is, if the Chi-squared number
185   is greater than 11.70, then the p-value is <0.05.
186 # When the p-value is <0.05, we say that the difference between
187   observation and expectation is statistically clear
188 # So let's find out how many of our Chi-squareds show a clear
189   difference!
190 propSig <- length( which( Chisqs > 11.70) ) / length(Chisqs)
191 percSig <- round(100 * propSig)
192
193 # percSig tells you the percent of trials that had a "significant" p-
194   value. Does that number surprise you?
195 # As a bit of foreshadowing, do you think the only thing driving that
196   very high number is natural selection?
197
198 # Percents are good, but plots are better. Let's make a plot of Chi-
199   squared results.
200 par(las = 1, mar = c(4, 4, 1, 1), mgp = c(2, 0.5, 0), tck = -0.01, cex.
201   axis=1)
202 hist(Chisqs, main="", xlab="chi-squared values", ylab="frequency")
203
204 # Let's see if there're any background-specific patterns
205 # We'll set up an empty plot to put data into
206 par(las = 1, mar = c(4, 4, 1, 1), mgp = c(2, 0.5, 0), tck = -0.01, cex.
207   axis=1)
208 plot(1, 1, xlim=c(0, 400), ylim=c(1, 8.5), xlab="", ylab="", type="n",
209   yaxt="n")
210
211 # Because this plot is weird, we're going to setup the y-axis and x-
212   axis label separately
213 axis(2, at = 1:length(backgrounds), labels = backgrounds)
214 mtext(side=1, expression(chi^2), cex=1.75, line=2.5)
215
216 # Now, FOR each background, we'll add a histogram of the data.
217 counter <- 1
218 for (i in backgrounds) {
219   Data <- Chisqs[which(results[,3] == i)]
220   addHist(Y=counter, Dat=Data, Color=backgroundCol[counter])

```

```

221   counter <- counter + 1
222 }
223
224 # Now let's add a line representing the critical value
225 abline( v = 11.70, lty=2, lwd = 2, col='black')
226
227 # The more to the right of that line a distribution is, the more often
228   trials on that background were "significant"
229 # Do you see any meaningful differences between backgrounds?

```

## 230 VI. Running a simulation

231 In the functions you downloaded from the website at the top, there's a function called **simDraws**.  
 232 This function will simulate the toothpick experiment, but will do so "blindly". Each time, 40 of  
 233 the 60 toothpicks are drawn *at random* without regard for color or background.

234 Basically, the **simDraws** simulates what would happen if the predator student had their eyes  
 235 closed throughout the "eating" process.

```

236 # Let's simulate running this experiment without natural selection ten
237   thousand times
238 Simulation <- simDraws(10000)
239
240 # Now let's add our Chi-squareds from the 10,000 simulations to our
241   plot
242 addHist(Y=7, Dat=Simulation, Color="lightgray")
243 mtext(side=2, at=7, line=0, "simulated")
244
245 # And compare it to the line, too
246 abline(v = 11.70, lty=2, lwd=2)

```

247 Look at the relationship between the critical value and the simulated distribution. What per-  
 248 centage of the time does the selection-free simulation find a "significant" (greater than 11.70) result?  
 249 Remember: a significant result means that there is a less than 5% chance that you'd get propor-  
 250 tions as different as what you observed if you were just randomly drawing toothpicks. How can the  
 251 selection-free simulation produce final counts so different from the initial counts?

252 Now, it'd also be interesting to simulate actual selection. That is, instead of simulating the  
 253 procedure as if the student playing predator had their eyes closed, instead simulate the experiment  
 254 as if the student playing predator *followed the instructions perfectly*.

255 The trick is, we don't actually know what that would look like! But just because we don't know  
 256 what it would look like doesn't mean we can't simulate it. In fact, not knowing how it would work  
 257 is actually where the power of simulations is greatest. We can simulate a *range* of scenarios, and  
 258 compare those simulated scenarios to the observed data!

259 So that's what we'll do. The idea is that, on a given background, some toothpicks will be more  
 260 likely to "survive" and "reproduce" than others. So we can set a "reproduce probability" (fitness;  
 261  $w$ ) for each toothpick type.

```

262 # no fitness differences. The Fit value for each of the six toothpick
263   colors are all 1.
264 Fit <- c(1, 1, 1, 1, 1, 1)

```

```

265 names(Fit) <- 1:6
266 Simulation2 <- simDraws(1e4, w = Fit)
267 addHist(Y=8, Dat=Simulation2, Color=rgb(0,0,0,0.25))
268
269 # one tooth pick type selected against
270 Fit <- c(0.1, 1, 1, 1, 1, 1)
271 names(Fit) <- 1:6
272 Simulation3 <- simDraws(1e4, w = Fit)
273 addHist(Y=8, Dat=Simulation3, Color=rgb(0,0,0,0.25))
274
275 # three tooth pick types selected against
276 Fit <- c(0.5, 0.6, 0.7, 1, 1, 1)
277 names(Fit) <- 1:6
278 Simulation4 <- simDraws(1e4, w = Fit)
279 addHist(Y=8, Dat=Simulation4, Color=rgb(0,0,0,0.25))
280
281 # five selected against
282 Fit <- c(0.1, 0.2, 0.3, 0.4, 0.5, 1)
283 names(Fit) <- 1:6
284 Simulation5 <- simDraws(1e4, w = Fit)
285 addHist(Y=8, Dat=Simulation5, Color=rgb(0,0,0,0.25))
286
287 # insane selection
288 Fit <- c(0.1, 0.1, 0.1, 0.1, 0.1, 1)
289 names(Fit) <- 1:6
290 Simulation6 <- simDraws(1e4, w = Fit)
291 addHist(Y=8, Dat=Simulation6, Color=rgb(0,0,0,0.25))
292 mtext(side=2, at=8, line=0, "sel. sim.")

```

293 Notice that none of these match the empirical distributions. The student-generated experiments  
294 have a peak that looks kind of like the no-selection simulations, but they also have a series of extra  
295 peaks out on the extreme end that look like the high-selection simulations. We can plot what a  
296 *mixture* of our simulations would look like really easily:

```

297 Simulation7 <- c(Simulation2, Simulation3, Simulation4, Simulation5,
298 Simulation6)
299 addHist(Y=8, Dat=Simulation7, Color=rgb(0,0,1,0.25))

```

300 How does the mixture compare to the student-generated data? Looking at how the mixture was  
301 generated, do most student groups show evidence of strong selection? If you had to describe the  
302 relative strength of different evolutionary processes across all the groups in all the labs across all  
303 of the years, what would you say?

## 304 VII. Inference

305 What evolutionary process(es) are at work in the lab-as-done-by-humans?

306 What evolutionary process(es) are at work in the lab-as-simulated-by-the-computer?

307 What do the graphs tell us about the *relative strength* of the evolutionary processes the BIOL-112  
308 students are simulating?

What tells you more about what processes occurring here: comparing the student numbers to a single critical value (11.70), or comparing the student numbers to the simulated numbers?

Imagine you added the possibility for a toothpick to mutate into a different type. What would that do to the  $\chi^2$  values? (hint: think about what part of this process causes such extreme distortions even without selection)

## VIII. Extra Credit

Modify the `simDraws` function to include the possibility that toothpicks could mutate during reproduction (i.e., some fixed probability  $\mu$  that an “offspring” would be a different color than their “parent”). Compare the calculated  $\chi^2$  values with mutation & drift to those with mutation, drift, & selection. The `simPop` function includes a mutation option with the same parameter name ( $\mu$  is pronounced & coded as “mu”), so you can use that as a partial guide.