

Task 04: 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_04, 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. Sampling

A sample is a subset of the total population that you take measurements/data from. So, if I wanted to measure the height of people in West Virginia and compare them to people in Virginia, it'd obviously be impossible for me to measure every single person in both states. So what I would do instead is take measurements from a small, randomly selected group of individuals in both states and compare the heights in my two small groups (samples).

If you've ever had trouble with this concept, come by my office! I have some code below to illustrate the effects of sampling, but this is a concept that's going to be central to almost everything we do for the rest of the semester, so even a little misunderstanding of it (if you allow that misunderstanding to persist) can be a problem.

```
# Make our populations
trueMean1 <- 5
trueSD1 <- 5
population1 <- rnorm(1e6, trueMean1, trueSD1)

trueMean2 <- 4
trueSD2 <- 5
population2 <- rnorm(1e6, trueMean2, trueSD2)

# Now take a sample of each population
Size <- 50
Sample1 <- sample(population1, Size)
Sample2 <- sample(population2, Size)
```

```

37
38 # Compare the samples! Are they different? Were the populations
39   different?
40 boxplot(Sample1, Sample2)

```

41 III. Basic Genetics

42 Way back in BIOL 111 and 112, you learned about the basic principles undergirding inheritance
 43 and genetics in DNA-based life forms of planet Earth. In sexually reproducing organisms (like us),
 44 multiple parents (two in mammals) contribute half of their genome to the offspring. The offspring
 45 then grows up, and when it's ready, scrambles the DNA it inherited from its parents via a process
 46 called recombination, then splits its DNA in half and passes that half off to *its* offspring.

47 We're going to model this process in R. We're going to assume there are four grandparents
 48 (paternal & maternal grandma & grandpa) and two parents (Alan and Brenda). We're going to
 49 track the genes passed on from each of the four grandparents (in the `makeFounder` function, you
 50 can manipulate the number of loci with the `len` argument).

51 We'll model 10,000 loci from each grandparent, and follow those loci into Alan & Brenda, and
 52 then into their children.

```

53 # Read in the needed functions
54 source("http://jonsmitchell.com/code/simFxn04.R")
55
56 # First, we'll make the grandparents. I would encourage you to use head
57   () and nrow() and such to examine these objects so you understand
58   them.
59 MatGrandma <- makeFounder("grandma_mom")
60 MatGrandpa <- makeFounder("grandpa_mom")
61 PatGrandma <- makeFounder("grandma_da")
62 PatGrandpa <- makeFounder("grandpa_da")
63
64 # Now, we're going to have the paternal grandparents make Alan. Again,
65   examine the object!
66 Alan <- makeBaby(PatGrandma, PatGrandpa)
67
68 # Now I would like for you to write the line of code that will make
69   Brenda.
70
71 # Once you have Alan and Brenda, they'll have their first child. This
72   child will be our focus.
73 Focus <- makeBaby(Brenda, Alan)
74
75 # Because each locus is tagged with which grandparent and parent it
76   came from, we can search to see how many genes came from any one
77   known ancestor of Focus. We'll use grep(), which is basically like
78   Ctrl+F. It will find any matches. So grandpa_mom and grandma_mom
79   both include "mom" and thus will both match to the line below. We'll
80   find how many loci Focus inherited from Brenda, and then divide
81   that by the total number of loci present. This gives us the percent

```

```

82   of genes shared between Brenda and Focus. Before you look at the
83   ToMom object, what should the number be?
84   ToMom <- length( grep("mom", Focus ) ) / length( Focus )
85
86   # We can also look at how many genes Focus shares with each of fhis
87   maternal grandparents. What should these numbers be? Do they match
88   your expectation?
89   ToMomMom <- length( grep( "grandma_mom", Focus ) ) / length( Focus )
90   ToMomDad <- length( grep( "grandpa_mom", Focus ) ) / length( Focus )
91
92   # Is Focus equally related to each maternal grandparent? What about
93   each paternal grandparent? Is this what you expected? What is the
94   average relatedness of Focus to all four grandparents?
95
96   # Let's imagine Focus gets a sibling. Yay, sibling!
97   Sibling_01 <- makeBaby(Brenda, Alan)
98
99   # How much DNA do you expect Focus to share with Sibling_01? Is that
100   the amount actually shared?
101   ToSib <- length( intersect( Focus, Sibling_01 ) ) / length( Focus )
102
103   # Now, let's imagine Brenda and Alan are extremely busy, and have 1,000
104   more children. How many genes does Focus share with each of the
105   1,000 siblings?
106   ManySiblings <- replicate( 1e3, length( intersect( Focus, makeBaby(
107     Brenda, Alan ) ) ) / length( Focus ) )
108
109   # We can summarize the data using quantiles and the mean
110   quantile(ManySiblings)
111   mean(ManySiblings)
112
113   # And we can also plot the data
114   plot(density(ManySiblings), main="", xlab="proportion shared genes")
115
116   # Please provide an explanation for why you see a range of values in
117   these analyses. Note that you can adjust the number of loci by
118   changing the "len" argument in the makeFounder() function calls
119   above.

```

120 IV. Hardy-Weinberg Equilibrium

```

121   # Given an allele frequency, p, we can calculate the expected genotype
122   frequencies
123   HWE <- function(p) {
124     aa <- p^2
125     ab <- 2 * p * (1 - p)
126     bb <- (1 - p)^2

```

```

127   return(c(aa=aa, ab=ab, bb=bb))
128 }
129 HWE(0.5)
130
131 # We can make a blank plot...
132 plot(1, 1, type="n", xlim=c(0, 1), ylim=c(0, 1), xlab="freq. allele a",
133      ylab="geno. freq")
134
135 # ...then calculate genotype frequencies for a bunch of allele
136   frequencies
137 p <- seq(from = 0, to = 1, by = 0.01)
138 GenoFreq <- t(sapply(p, HWE))
139
140 # Now we can plot our known allele frequency (p) against our expected
141   genotype frequencies (GenoFreq)
142 lines(p, GenoFreq[, "aa"], lwd=2, col="red")
143
144 # Before moving on, can you read and understand this plot? What happens
145   to the frequency of aa individuals as the frequency of the a allele
146   increases in the population? What happens to the aa frequency as a
147   decreases? Is time shown on this plot? Is geographic space? If you
148   want to do well in this class it is vital that you know exactly what
149   is (and is not) shown on a graph like this.
150
151 # Now, let's add the other genotypes
152 lines(p, GenoFreq[, "ab"], lwd=2, col="purple")
153 lines(p, GenoFreq[, "bb"], lwd=2, col="blue")
154 legend("top", legend=c("aa", "ab", "bb"), col=c("red", "purple", "blue"),
155      lty=1, lwd=2, bty="n")
156
157 # Let's use this now! Let's simulate a population, and look at how its
158   allele and genotype frequencies vary.
159 Pop <- simPop(500)
160
161 # Now let's add these points to the HWE plot you made above
162 points(Pop[, "freqa"], Pop[, "Genotypes.aa"]/500, pch=21, bg="red")
163
164 # Does the frequency of the aa genotype in your population match the
165   expectation from Hardy-Weinberg?
166
167 # What if we do another simulation with a much smaller population
168 Pop <- simPop(50)
169 points(Pop[, "freqa"], Pop[, "Genotypes.aa"]/50, pch=22, bg="red")
170
171 # What's changed and why?

```

172 V. Two-Allele Drift

```

173 # install the learnPopGen package
174 library(learnPopGen)
175
176 # Run this over and over again, watching the lines.
177 # Ne is how many individuals there are in each population
178 # nrep is how many populations you are simulating at once
179 # pause is just how fast the lines grow
180 x <- genetic.drift(Ne=200, nrep=5, pause=0.01)
181
182 # Run and rerun the line above, changing Ne, and noting the patterns
183     you see
184
185 # Let's look at how population sizes effect the time to extinction for
186     one allele
187 # First, we'll make a bunch of populations of different sizes, from 5
188     to 50 individuals
189 PopSizes <- 5:50
190
191 # Next, we'll say that there are 5 populations with each given size
192 Samples <- rep(PopSizes, 5)
193
194 # Now, we'll simulate all 230 of those populations & get the time one
195     of the two alleles went extinct
196 # Note: this will take some time (~20 sec on my computer, maybe longer
197     or shorter depending on your machine)
198 tExt <- sapply(Samples, function(x) nrow(simPop(x, 500)))
199
200 # To fit a line (Linear Model) to data in R you use the lm() function.
201 Line <- lm(tExt ~ Samples)
202
203 # To see the fit, use summary()
204 summary(Line)
205
206 # To extract the coefficients, use $coef
207 Line$coef
208
209 # To add it to the plot, use abline()
210 plot(Samples, tExt)
211 abline(Line)
212
213 # You may consider comparing the above to Line2 <- lm( tExt~Samples + 0
214     ), and ask yourself which one makes more sense given these data (
215     note: run them both and look at the output to figure out what the +0
216     does!)
217
218 # Linear models assume that all points are roughly equally far from the
219     line, and that there are as many above the line as below. What do
220     you notice, looking at this graph, about the distance between the

```

221 points and the line as the population size increases? What does that
222 mean to you?

223 **VI. Extra Credit**

224 The issue identified in the last question has a name: heteroskedasticity. Use your Google skills to
225 find a way to fit a linear model that is *robust* to heteroskedasticity and compare the estimate of
226 the slope.