

About Me

Eric C. Anderson

Contents

Who I am and where I came from	1
Research Interests	3
Influential papers	3
The mathematics behind my research	3
My computing experience	3
What I hope to get out of this class	5
Evaluating some R code	5
Citations	6

Who I am and where I came from

I grew up in a small town in Southern California. When I was 11 years old, I spent five weeks living with a family friend who ran the hyperbaric chamber adjacent to the University of Southern California Marine Lab on Catalina Island. Days spent snorkeling the kelp beds, keeping intertidal organisms in an aquarium I was allowed to fill, and watching researchers tag blue sharks convinced me that I wanted to be a marine biologist.

In high school, however, I got into rock climbing and became a little math/physics nerd. I went to Stanford University planning to major in Mathematics, but got super intimidated by all the smart, wonky people in my first honors math class. I took a year-long detour to Prescott College, where I had my first class in evolution and wondered why the hell my high-school biology class was not taught in the context of evolution.

I ended up studying human biology at Stanford, and then went to University of Washington for graduate work. I started in Fisheries, but then got involved in genetics research that needed some new statistical developments. Thus, my math nerdiness was resurrected and I did my PhD in Quantitative Ecology and Resource Management, focusing on the use of Monte Carlo methods for inference from population genetic models.

After a two-year postdoc at Berkeley, I started working for the National Marine Fisheries Service in Santa Cruz, CA. I actually did become a marine biologist, sort of! I still work for NMFS, but a year and a half ago I moved to Fort Collins with Kristen, and have an affiliate position in FWCB.

When I am not working I love getting out and being active. My top four things to do are:

1. Snorkeling in rivers and creeks, backpacking, and hiking (all with my family),
2. Playing hockey (I learned to skate after moving here. What a blast!),
3. Biking,
4. Playing with our pair of awesome, springy cats.

Here is a picture of me with daughter Zoe, looking for aquatic invertebrates above the CSU Mountain Campus.



Research Interests

I'm interested in all manner of statistical inference from genetic data. Lately I have been working on the genetic basis of run timing in Chinook salmon and other salmonids.

Influential papers

When I was a graduate student, I heard Peter Green speak about his work on reversible jump MCMC for the analysis of finite mixture models. One of the problems I was working with at the time was estimating proportions of salmon from different rivers that were being caught in the ocean—the mixed stock fishery problem. I spent a lot of time working through Richardson and Green (1997) and learned a lot about MCMC and RJMCMC in the process.

Later on, much of my work on Bayesian inference of pedigrees from genetic data (Anderson and Ng 2016) builds upon the idea of factor graphs described by Kschischang et al. (2001).

The mathematics behind my research

I have worked a lot with the coalescent process, so let's put down the expected time during which there are k extant lineages in a population of size N .

$$\mathbb{E}T_k = \frac{4N}{k(k-1)}.$$

And, while we are at, let's throw down a description of one of the update steps in the sum-product algorithm for acyclic factor graphs:

$$\mu_{f_j \rightarrow v_i}(x_i) = \sum_{x_{C \setminus i} \in \mathcal{X}_{C \setminus i}} h_j(x_{C \setminus i}, x_i) \prod_{k \in C \setminus i} \mu_{v_k \rightarrow f_j}(x_k).$$

My computing experience

I started programming in BASIC on our old Apple IIe in 1983. In high school I implemented a basic program to plot some fractal images. After that, I didn't really do any programming until grad school when I took a course in C. Here is some C code that I wrote:

```
if(RU!=NULL) {
    RepUnitZSum = (int *)calloc(RU->NumRepUnits,sizeof(int));
    RepUnitPis = DvalVector(0,RU->NumRepUnits, 0.0, 1.0, .01);
    RepUnitPofZs = (dval ***)calloc(N,sizeof(dval **));
    for(i=0;i<N;i++) {
        RepUnitPofZs[i] = DvalVector(0,RU->NumRepUnits-1, 0.0,1.0, -1); /* no histograms for th
    }
    if(B0->PiTraceInterval>0) {
        repPi_tracef = OpenTraceFile("rep_unit_pi_trace.txt", "trace file of reporting unit Pi valu
    }
    if(B0->ZSumTraceInterval>0) {
        repZSumtracef = OpenTraceFile("rep_unit_zsum_trace.txt", "trace file of reporting unit ZSum
    }
}
```

Wow! That is pretty ugly.

When I was a postdoc, John Novembre and the other members of Monty Slatkin's lab at Berkeley got me hooked on using the Unix shell, programming in bash, and writing short scripts in awk and sed. Here is a little awk script that takes the output of SGE's `qacct` command and makes a nice, tidy table of it

```
#!/usr/bin/awk -f

# an awk script.
# it expects the output of qacct like this:
# qacct -o eriq -b 09271925 -j ml

# make it executable and run it like this:
# qacct -o eriq -b 09271925 -j ml | tidy-qacct

# if you pass it a job number that was not one of your jobs it
# just skips the error message that comes up.

# note that the output of qacct is space delimited

/^=====/ {++n; next} # increment run counter, then skip these lines
/^error:/ {next} # skip it if you told it to get a wrong job number

# now, every data line it gets things. It compiles the header
# all the way through, in case some reports have more columns...
NF > 1 {
    tag = $1;
    if(!(tag in header_vals)) {
        header[++h] = tag;
        header_vals[tag]++;
    }
    $1 = ""; # remove the tag from the full line of stuff
    values[n, tag] = $0; # assign the values to the tag
}

# at the end of it all, we print the header and then all the values:
END {
    # print the header
    printf("%s", header[1]);
    for(i=2;i<=h;i++)
        printf("\t%s", header[i]);
    printf("\n");

    # cycle over individuals and fields and print em
    for(j=1;j<=n;j++) {
        printf("%s", values[j,header[1]]);
        for(i=2;i<=h;i++)
            printf("\t%s", values[j,header[i]]);
        printf("\n");
    }
}
```

I used to rather dislike the R programming language and felt it was dreadfully slow. It has gotten a lot better in the last two decades. The introduction of Hadley Wickham's tidy data analysis framework has really

improved things.

What I hope to get out of this class

I hope that I will:

- Help students understand enough about Unix and programming to lessen the pain of learning to do bioinformatics.
- Be able to advance students' own research.
- Impart to the students an appreciation of the importance of making research reproducible.

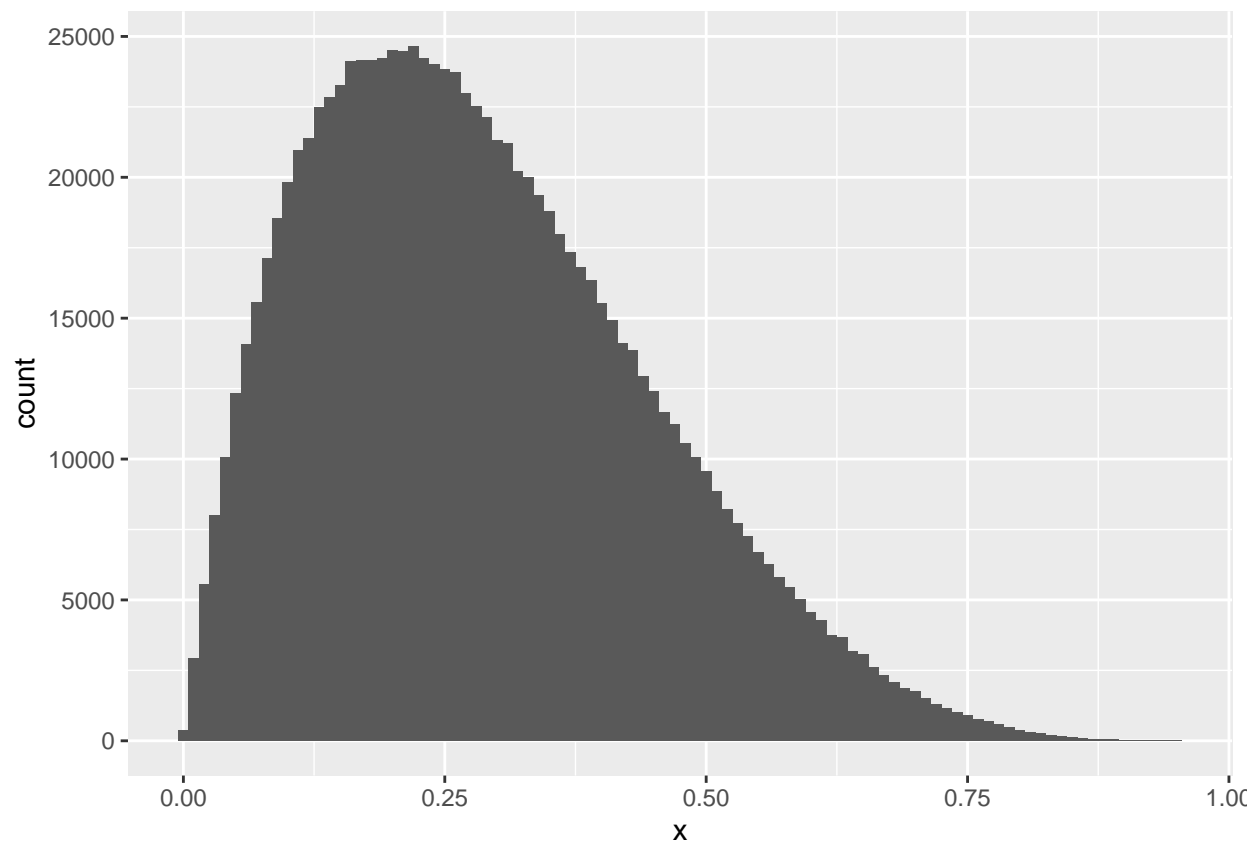
Evaluating some R code

I'm going to just simulate one million beta random variables from a Beta(2,5) distribution and plot a histogram of it.

```
library(tidyverse)

beta_rvs <- tibble(
  x = rbeta(1e06, shape1 = 2, shape2 = 5)
)

ggplot(beta_rvs, aes(x = x)) +
  geom_histogram(binwidth = 0.01)
```



Citations

Anderson, Eric C, and Thomas C Ng. 2016. “Bayesian Pedigree Inference with Small Numbers of Single Nucleotide Polymorphisms via a Factor-Graph Representation.” *Theoretical Population Biology* 107. Elsevier: 39–51.

Kschischang, Frank R, Brendan J Frey, Hans-Andrea Loeliger, and others. 2001. “Factor Graphs and the Sum-Product Algorithm.” *IEEE Transactions on Information Theory* 47 (2): 498–519.

Richardson, Sylvia, and Peter J Green. 1997. “On Bayesian Analysis of Mixtures with an Unknown Number of Components (with Discussion).” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 59 (4). Wiley Online Library: 731–92.