

BIOS 6611 Homework 6

Due Monday, October 15, 2018 **by noon** to Canvas Assignment Basket

Graft-versus-host disease (GvHD) is a secondary complication of allogeneic (from a donor) hematopoietic stem cell transplantation (the only potentially curative treatment for leukemia and many other disorders). It is characterized by the transplanted cells of the donor inappropriately attacking the tissues of the transplant recipient, or host. You're running a large cohort study of stem cell recipients in which a placebo is given to Treatment Group A and a new GvHD prophylaxis drug is given to Treatment Group B.

The risk factors for GvHD are extremely complex and still little understood. However, the characteristics of the donor are extremely important. HLA-matched siblings are considered the "gold standard" donor because they appear to be associated with substantially decreased GvHD risk (systematic reviews suggest 10-40% GvHD incidence in HLA-matched, related donors, versus 50 - 80% in other donors).

We will not perform a "complete" analysis of this data (which would include commenting on whether or not we observe a difference in GvHD risk by treatment), but focus on viewing the relationships between the prior and the posterior distributions of responses to the drugs.

- A) Load `gvhd.txt` into R, then subset the data to focus on only transplant recipients with an HLA-matched sibling donor.
- B) Calculate the proportion of recipients that got GvHD in the Treatment A group. Repeat for Treatment B.
- C) Among transplant recipients with HLA-matched donors, is there a significant association between treatment and GvHD at the 5% level of significance? Carry this test out using both a permutation test, and either an exact or asymptotic method, as appropriate. Summarize your results and comment on differences, if any, between the two methods you applied.
- D) Using the `seq()` function, create a vector called `p_grid` that has 30 evenly spaced probabilities from 0 to 1.
- E) Assume that whether or not a patient has GvHD is a binary feature modeled by a Bernoulli distribution (see Lecture 4). Using the `dbinom()` function, find the likelihood of the number of GvHD cases among subjects in Treatment A at each value in `p_grid`. You should end up with a 30-element long vector of probabilities. Save this vector as "likelihood".

- F)** Use the following code to generate a possible prior distribution of the probability of GvHD for HLA-matched, related donors (which is based on the existing literature information). What prior distribution does this represent (e.g., normal, Poisson, uniform) and what parameters does this distribution have?

```
prior_MRD <- ifelse( p_grid > 0.1 & p_grid < 0.4, 0.3, 0)
```

Note: This is an “improper” prior, i.e. it does not integrate to 1. This is ok, but you will have to account for this when you obtain the mean proportion of GvHD based on this prior in part H below.

- G)** Calculate the posterior distribution, using the following code:

```
posterior <- likelihood * prior_MRD / sum(likelihood * prior_MRD)
```

- H)** Find the means of the prior distribution and the posterior distribution numerically. Hint: Recall the definition of expected value for discrete events.
- I)** Plot the likelihood, prior, and posterior (as Y-variables) against `p_grid` (X-variable) for Treatment group A in the same figure. Make the line of each distribution a different color. Summarize what you observe.
- J) EXTRA CREDIT:** Repeat parts E through I for those who received Treatment B. Comment on how likely you think there is to be a difference between the two treatments.
- K) EXTRA CREDIT:** To see what happens with smaller sample sizes, randomly (but reproducibly!) subsample your data so that there are only 5-10 subjects for Treatment A. Remake the plot in part I (i.e. the plot with the prior, likelihood and posterior for Treatment A) with this subsample. Obtain the posterior mean based on the smaller sample. What do you *qualitatively* notice about the posterior distribution and mean with the smaller sample compared with the posterior distribution and mean from the larger sample?