

# 36-350 Lab 6 – October Surprise

The distribution of talent is a major question of statistical research and a foundation of many young statistical careers (including 75% of your instructors), and there is no shortage of the distribution of observable talent in professional sports, of interest to many statisticians (including 50% of your instructors). In today's lab we will explore the distribution of one particular talent in baseball, reaching base safely, using the Beta distribution.

The Beta is a random variable bounded between 0 and 1 and often used to model the distribution of proportions. The probability distribution function for the Beta with parameters  $\alpha$  and  $\beta$  is

$$p(x|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}$$

where  $\Gamma()$  is the Gamma function, the generalized version of the factorial. Thankfully, for this assignment, you need not know what the Gamma function is; you need only know that the mean of a Beta is  $\frac{\alpha}{\alpha+\beta}$  and its variance is  $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ .

For this assignment you will test the fit of the Beta distribution to the on-base percentages (OBPs) of hitters in the 2014 Major League Baseball season; each plate appearance (PA) results in the batter reaching base or not, and this measure is the fraction of successful attempts. This set has been pre-processed to remove those players with an insufficient number of opportunities for success.

## Part I

1. Load the file [<http://www.stat.cmu.edu/~acthomas/mlb-obp.csv>] into a variable of your choice in R. How many players have been included? What is the minimum number of plate appearances required to appear on this list? Who had the most plate appearances? What are the minimum, maximum and mean OBP?
2. Plot the data as a histogram with the option `probability=TRUE`. Add a vertical line for the mean of the distribution. Does the mean coincide with the mode of the distribution?
3. Eyeball fit. Add a `curve()` to the plot using the density function `dbeta()`. Pick parameters  $\alpha$  and  $\beta$  that matches the mean of the distribution but where their sum equals 1. Add three more `curve()`s to this plot where the sum of these parameters equals 10, 100 and 1000 respectively. Which of these is closest to the observed distribution?

## Part II

4. Method of moments fit. Find the calculation for the parameters from the mean and variance from [[http://en.wikipedia.org/wiki/Beta\\_distribution](http://en.wikipedia.org/wiki/Beta_distribution)] and solve for  $\alpha$  and  $\beta$ . Create a new density histogram and add this `curve()` to the plot. How does it agree with the data?
5. Calibration. For the previous part, find the 99 percentiles of the actual distribution using the `quantile()` function and plot them against the 99 percentiles of the beta distribution you just fit using `qbeta()`. How does the fit appear to you?
6. Optional if you have time – MLE fit. Create a function for the log-likelihood of the distribution that calculates `-sum(dbeta(your.data.here, your.alpha, your.beta, log=TRUE))` and has one argument `p=c(your.alpha, your.beta)`. Use `nlm()` to find the minimum of the negative of the log-likelihood. Take the MOM fit for your starting position. How do these values compare?