

# MATH 390.03-02 / 650 Fall 2015 Homework #3

Professor Adam Kapelner

Due 4PM in my mail slot, Friday, February 26, 2015

(this document last updated Wednesday 17<sup>th</sup> February, 2016 at 6:45pm)

## Instructions and Philosophy

The path to success in this class is to do many problems. Unlike other courses, exclusively doing reading(s) will not help. Coming to lecture is akin to watching workout videos; thinking about and solving problems on your own is the actual “working out.” Feel free to “work out” with others; **I want you to work on this in groups.**

Reading is still *required*. For this homework set, read about the beta prior, the binomial-beta bayesian formulation and the beta-binomial model. Also read ch4 in McGrayne.

The problems below are color coded: **green** problems are considered *easy* and marked “[easy]”; **yellow** problems are considered *intermediate* and marked “[harder]”; **red** problems are considered *difficult* and marked “[difficult]” and **purple** problems are extra credit. The *easy* problems are intended to be “giveaways” if you went to class. Do as much as you can of the others; I expect you to at least attempt the *difficult* problems.

Problems marked “[MA]” are for the masters students only (those enrolled in the 650 course). For those in 390, doing these questions will count as extra credit.

This homework is worth 100 points but the point distribution will not be determined until after the due date. See syllabus for the policy on late homework.

Up to 10 points are given as a bonus if the homework is typed using L<sup>A</sup>T<sub>E</sub>X. Links to installing L<sup>A</sup>T<sub>E</sub>X and program for compiling L<sup>A</sup>T<sub>E</sub>X is found on the syllabus. You are encouraged to use [overleaf.com](http://overleaf.com). If you are handing in homework this way, read the comments in the code; there are two lines to comment out and you should replace my name with yours and write your section. The easiest way to use overleaf is to copy the raw text from hwxx.tex and preamble.tex into two new overleaf tex files with the same name. If you are asked to make drawings, you can take a picture of your handwritten drawing and insert them as figures or leave space using the “\vspace” command and draw them in after printing or attach them stapled.

The document is available with spaces for you to write your answers. If not using L<sup>A</sup>T<sub>E</sub>X, print this document and write in your answers. I do not accept homeworks which are *not* on this printout. Keep this first page printed for your records.

NAME: \_\_\_\_\_

## Problem 1

These are questions about McGrayne's book, chapters 4-7.

- (a) [easy] Describe four things Bayesian modeling was applied to during WWII and identify the people who developed each application.
- (b) [harder] What do you think was the main reason Bayesian Statistics fell out of favor at the end of WWII?
- (c) [harder] Why weren't the leaders of Statistics world in the 1950's able to answer the think-tank's question about the  $\mathbb{P}$ (war in the next 5 years)?
- (d) [easy] Who was responsible for reviving the interest in Bayesian Statistics post-WWII and why?

- (e) [difficult] In 1955, there were no midair collisions of two planes. How was the actuary able to estimate that the number would be above zero?
- (f) [easy] The main attack on Bayesian Statistics has always been subjectivity. Answer the following question how Savage would have answered it: “If prior opinions can differ from one researcher to the next, what happens to scientific objectivity in data analysis?” Do you believe Savage’s idea is the way science works in the real world?
- (g) [difficult] [MA] On page 104, Sharon writes, “Bayesians would also be able to concentrate on what happened, not on what *could* have happened according to Neyman Pearson’s sampling plan”. (Note that the “Neyman Pearson’s sampling plan” is synonymous with Frequentist Statistics). Explain (1) how Bayesians concentrate on “what happened” and (2) how Frequentists concentrate on what “*could* have happened” in the context on page 104.
- (h) [easy] Who were the two tireless champions of Bayesian Statistics throughout the 50’s, 60’s and 70’s and where geographically were they located during the majority of their career?

## Problem 2

This problem is concerned with the logical definition of probability. As a review, we have:

1. **The Objective View** — This is the view that probabilities are properties of the physical world and can only be defined by either
  - (a) its **Long Run Frequency** Seeing the same event over and over again and tabulating when the event occurs will create a frequency which will asymptotically become the probability or
  - (b) its **Propensity** which means deep down inside, the physical object is wired for events in certain proportions.

Thus, events that are non physical such as the probability of Donald Trump winning the 2016 election is outside of the purview of probability. The objective view of probability is tied to the frequentist view of statistics. We also have the...

2. **The Epistemic View** — This is the view that probabilities are inherently living inside the minds of human beings who are forced to grapple with uncertainty as they see it. Laplace believed probability is an illusion because we don't have certainty about the universe. The two definitions here are that probability...
  - (a) is **Logical** which means that given the same information, everyone would come to the same conclusion.
  - (b) is **Subjective** which means that given the same information, everyone would *not* come to the same conclusion. Thus probability is defined as the degree of belief of some individual which differs from another individual.

Thus events that are non physical such as the probability of Donald Trump winning the 2016 election can now be legitimate probabilities as they can be computed.

- (a) [easy] We discussed last time that the logical definition requires the principle of indifference which goes by many different names. We will now go about showing that the principle of indifference has a tenuous foundation and thereby rendering the logical theory of probability inadequate. Thus, our conclusion will be that Bayesian Statistics runs on the Subjective definition which we may develop in a later homework. We begin with demonstrating a paradox in the logical definition for a discrete set of  $\theta_0$ .

Imagine you have a library with thousands of books but all are either red, green, yellow or purple but you don't know the proportions of the books' colors. Imagine you are blindfolded and select a random book and you are only interested if it's *red* or *not red*. According to the principle of indifference, what is your prior probability that the book is red? Remember,  $|\Theta_0| = 2$  here.

- (b) [easy] Imagine you are blindfolded and select a random book and you are interested if it's red, green, yellow or purple. According to the principle of indifference, what is your prior probability that the book is red? Remember,  $|\Theta_0| = 4$  here.
- (c) [harder] Why do (a) and (b) constitute a paradox? Does this limit the application of the principle of indifference?
- (d) [easy] Now we're going to work on trashing the principle of indifference for continuous measures. Thanks to Wikipedia... Imagine I have a cube-shaped box. The length of the side we call  $S$  and we know it's less than 1 inch in length. Thus its prior distribution under the principle of indifference should be  $S \sim U(0, 1)$ . What is the expected length of the side a priori (that is according to the prior distribution)? Please do not overthink this.
- (e) [harder] Given the same prior as previously, imagine the surface area which is calculated as  $6S^2$ . Find expectation of the surface area. You may need to look at your Math 241 notes to get the variance of the uniform r.v. Is the answer you get equal to using the prior mean length and computing the surface area based on that, i.e.  $6\mathbb{E}[S]^2$ ?

(f) [harder] We will not prove this, but it shouldn't come as a surprise that the surface area is not distributed uniformly between 0 and 6 given what you saw in your last answer. What does this mean for the principle of indifference? It only works...

(g) [easy] It may be argued that one is only indifferent to the length of the side and that it's okay this implies a non-indifference to the surface area and volume. But here's a paradox that's harder to argue with which I got from Gillies' book. Imagine you have a drink in front of you made up of some parts orange juice and some parts water. You have a prior belief that the ratio of wine/water is  $U(1/4, 4)$  using the principle of indifference. What is the probability (according to your prior belief) that the  $\mathbb{P}(\text{wine/water} \geq 2)$ ?

(h) [easy] Using the principle of indifference, what does this imply that the ratio of water/wine should be  $U(1/4, 4)$  as well? It shouldn't matter whether you pick wine divided by water or water divided by wine, right? Answer yes/no and give your gut reaction.

(i) [easy] Assuming that you wrote "yes" to the previous problem, calculate

$$\mathbb{P}\left((\text{wine/water})^{-1} \geq 2^{-1}\right) = \mathbb{P}\left(\text{water/wine} \leq \frac{1}{2}\right).$$

And: is this different to the answer you got in (e)?

- (j) [difficult] Why is this a paradox? Write a couple sentences how this should doom the principle of indifference in the case of a continuous  $\Theta_0$  space.

- (k) [harder] [MA] Now we're going to see how this fails. Under the prior belief that the ratio of wine/water is  $U(1/4, 4)$ , derive the PDF of the ratio of water/wine. Is it also  $U(1/4, 4)$ ? If you need a refresher on this stuff, see here. (Note: we're going to see Jeffrey's answer to this issue soon enough in class).

### Problem 3

A quick question on de Finetti's theorem for MA students.

- (a) [easy] [MA] If I tell you that  $X_1, \dots, X_n$  are exchangeable, what does this guarantee? What kind of model can be built? Reference de Finetti's theorem.

## Problem 4

We will now be looking at the beta-prior, binomial-likelihood Bayesian model.

- (a) [easy] Using the principle of indifference, what should the prior on  $\theta$  (the parameter for the Bernoulli model) be?
- (b) [easy] Let's say  $n = 6$  and your data is 0, 1, 0, 1, 0, 1. What is the likelihood of this event?
- (c) [easy] Does it matter the order as to which the data came in? Yes/no.
- (d) [harder] Show that the unconditional joint probability (the denominator in Bayes rule) is a beta function and specify its two arguments. We did this in class.
- (e) [harder] Put your answer from (a), (b) and (d) together to find the posterior probability of  $\theta$  given this dataset. Show that it is equal to a beta distribution and specify its parameters.



(f) [easy] Now imagine you are not indifferent and you have some idea about what  $\theta$  could be a priori and that subjective feeling can be specified as a beta distribution. (1) Draw the five basic shapes that the beta distribution can take on, (2) give an example of  $\alpha$  and  $\beta$  values that would produce these shapes and (3) write a sentence about what each one means for your prior belief. These shapes are in the notes.

(g) [harder] Imagine  $n$  data points of which you don't know the realization values. Show that  $\theta \mid X \sim \text{Beta}(\alpha + x, \beta + (n - x))$ . Note that  $x := \sum_{i=1}^n x_i$  which is the total number of successes and thereby  $n - x$  is the total nubmer of failures. The answer is in the notes but try to do it yourself.

- (h) [easy] What does it mean that the beta distribution is the “conjugate prior” for the binomial likelihood?
- (i) [harder] Stare at that distribution,  $\theta \mid X \sim \text{Beta}(\alpha + x, \beta + (n - x))$ . Some say the values of  $\alpha$  and  $\beta$  can be interpreted as follows:  $\alpha$  is considered the prior number of successes and  $\beta$  is considered the prior number of failures. Why is this a good interpretation?
- (j) [harder] By the principle of indifference, how many successes and failures is that equivalent to seeing a priori?
- (k) [easy] Why are large values of  $\alpha$  and/or  $\beta$  considered to compose a “strong” prior?
- (l) [harder] [MA] What is the weakest prior you can think of and why?

(m) [difficult] I think a priori that  $\theta$  should be expected to be 0.8 with a standard error of 0.02. Solve for the values of  $\alpha$  and  $\beta$  based on my a priori specification.

(n) [difficult] Prove that the posterior predictive distribution is  $X^* | X \sim \text{Bernoulli} \left( \frac{x+\alpha}{n+\alpha+\beta} \right)$ .  
MA students — do this yourself. Other students — use my notes and justify each step.  
I use a property of the gamma function.

- (o) [harder] The frequentist estimate of  $\theta$  is  $\hat{p} = 3/6 = 0.5$ . So a frequentist would probably use a posterior predictive distribution (if he had such a thing) as  $X^* \mid X \sim \text{Bernoulli}(0.5)$ . Why conceptually does this answer differ from your answer in (n)?

- (p) [easy] Assume the dataset in (b) where  $n = 6$ . Assume  $\theta \sim \text{Beta}(\alpha = 2, \beta = 2)$  a priori. Find the  $\hat{\theta}_{\text{MAP}}$ ,  $\hat{\theta}_{\text{MMSE}}$  and  $\hat{\theta}_{\text{MAE}}$  estimates for  $\theta$ .

For the  $\hat{\theta}_{\text{MAE}}$  estimate, you'll need to obtain a quantile of the beta distribution. Use R on your computer or online using R-Fiddle. The `qbeta` function in R finds arbitrary beta quantiles. Its first argument is the quantile desired e.g. 2.5%, the next is  $\alpha$  and the third is  $\beta$ . So to find the 97.5%ile of a  $\text{Beta}(\alpha = 2, \beta = 2)$  for example you type `qbeta(.975, 2, 2)` into the R console.

- (q) [harder] Why are all three of these estimates the same?

- (r) [easy] Write out an expression for the 95% credible region for  $\theta$ . Then solve computationally using the `qbeta` function in R.
- (s) [easy] Compute a 95% frequentist CI for  $\theta$ .
- (t) [difficult] Let  $\mu : \mathbb{R} \rightarrow \mathbb{R}^+$  be the Lebesgue measure which measures the length of a subset of  $\mathbb{R}$ . Why is  $\mu(\text{CR}) < \mu(\text{CI})$ ? That is, why is the Bayesian Confidence Interval tighter than the Frequentist Confidence Interval? Use your answers from (r) and (s).
- (u) [easy] Explain the disadvantages of the highest density region method for computing credible regions.