

# MATH 343 / 643 Homework #1

Professor Adam Kapelner

Due 11:59PM March 2 on github

(this document last updated 10:29pm on Tuesday 25<sup>th</sup> February, 2025)

## 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 as much as you can online about the topics we covered.

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.

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 7 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**. 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 general questions about Gibbs Sampling.

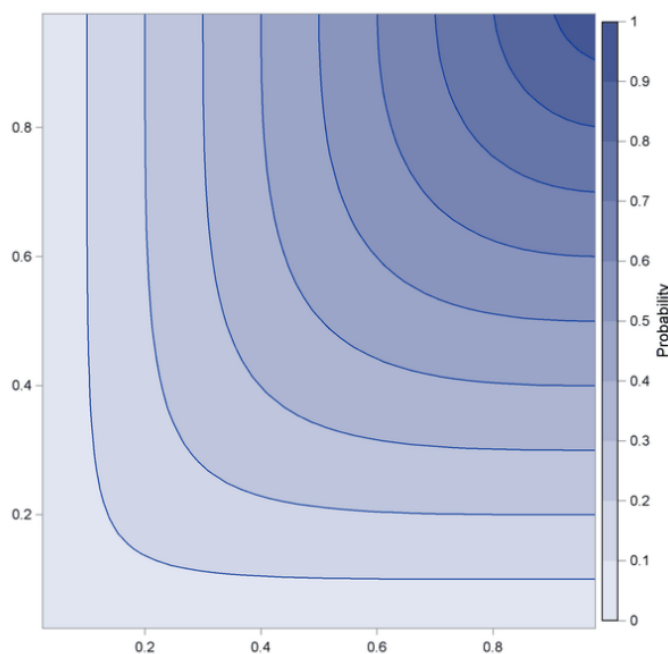
- (a) [easy] Let  $\dim[\boldsymbol{\theta}] = p$  and assume a prior  $f(\boldsymbol{\theta})$  to be continuous. Describe the steps of the systematic sweep Gibbs Sampler algorithm below that will converge to  $f(\boldsymbol{\theta}|\mathbf{X})$ . Label the steps that are necessary for the  $p$  dimensions separately e.g. Step 2.1, Step 2.2, ..., Step 2.p. You need to reference these step numbers later on in the problem.

- (b) [easy] What are all the items you need to know in order to write the code that implements a Gibbs Sampler?

- (c) [easy] Explain what burning of the Gibbs sample chain is and why it is necessary.

(d) [easy] Explain what thinning of the chain is and why it is necessary.

(e) [easy] Pretend you are estimating  $\mathbb{P}(\theta_1, \theta_2 | X)$  and the joint posterior looks like the picture below where the  $x$  axis is  $\theta_1$  and the  $y$  axis is  $\theta_2$  and darker colors indicate higher probability. Begin at  $[\theta_1, \theta_2] = [0.5, 0.5]$  and simulate 5 iterations of the systematic sweep Gibbs sampling algorithm by drawing new points on the plot.



## Problem 2

Consider a count model that has many zeroes. We choose to fit it with a hurdle model

$$X_1, \dots, X_n \stackrel{iid}{\sim} \begin{cases} 0 & \text{w.p. } \theta_1 \\ \text{ShiftedExtNegBinomial}(\theta_2, \theta_3, +1) & \text{w.p. } 1 - \theta_1 \end{cases}$$

where the shifted distribution is just the extended negative binomial distribution so that the probability of realizing a count of one is the probability of realizing a count of zero, the probability of realizing a count of two is the probability of realizing a count of one, etc. i.e.

$$\text{ShiftedExtNegBinomial}(\theta_2, \theta_3, +1) := p(x) = \frac{\Gamma(x_i - 1 + \theta_2)}{(x_i - 1)! \Gamma(\theta_2)} (1 - \theta_3)^{x_i - 1} \theta_3^{\theta_2}.$$

- (a) [harder] What is the parameter space for all three parameters of interest? This may require looking at your MATH 340 notes.
- (b) [harder] Assume a flat prior  $f(\theta_1, \theta_2, \theta_3) \propto 1$ . Find the kernel of the posterior distribution  $f(\theta_1, \theta_2, \theta_3 | \mathbf{x}, n_0, n_+)$  where  $\mathbf{x} := \{x_1, \dots, x_n\}$ , the observations. Let  $n_0$  be the number of zeroes in the dataset and  $n_+ := n - n_0$ , the number  $> 0$  in the dataset.
- (c) [easy] Find the conditional distribution  $f(\theta_1 | \mathbf{x}, n_0, n_+, \theta_2, \theta_3)$  as a brand name rv.
- (d) [easy] Find the kernel of the conditional distribution  $f(\theta_2 | \mathbf{x}, n_0, n_+, \theta_1, \theta_3)$ .

- (e) [easy] Is the conditional distribution  $f(\theta_2 | \mathbf{x}, n_0, n_+, \theta_1, \theta_3)$  a brand name rv? Yes/no
- (f) [easy] Find the conditional distribution  $f(\theta_3 | \mathbf{x}, n_0, n_+, \theta_1, \theta_2)$  as a brand name rv.
- (g) [easy] Is it possible to get inference for this model using a Gibbs Sampler? Why or why not?

### Problem 3

Consider the change point model

$$X_1, X_2, \dots, X_{\theta_3} \stackrel{iid}{\sim} \mathcal{N}(\theta_1, \sigma_1^2) \text{ independent of } X_{\theta_3+1}, X_{\theta_3+2}, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\theta_2, \sigma_2^2)$$

- (a) [harder] What is the parameter space for all five parameters of interest?
- (b) [harder] Assume a flat prior  $\theta_1, \theta_2, \theta_3$  and Jeffrey's prior for  $\sigma_1^2, \sigma_2^2$  which are assumed a priori independent of one another. Find the kernel of the posterior distribution.

(c) [harder] Find the kernels of all five conditional distributions. If they are proportional to a known distribution, name it.

(d) [harder] Find the conditional PMF of  $\theta_3$ .

(e) [easy] Is it possible to get inference for this model using a Gibbs Sampler? Why or why not?

#### **Problem 4**

Consider the discrete mixture model:

$$X_1, \dots, X_n \stackrel{iid}{\sim} \begin{cases} \text{Poisson}(\theta_0) & \text{w.p. } \rho \\ \text{Poisson}(\theta_1) & \text{w.p. } 1 - \rho \end{cases}$$

(a) [harder] What is the parameter space for all three parameters of interest?

(b) [harder] Assume a flat prior on all parameters. Find the kernel of the posterior distribution.

(c) [easy] Is this proportional to any known distribution?

(d) [harder] Is it possible to make a Gibbs Sampler to get inference here? Why or why not.

(e) [harder] Let's use data augmentation. Add  $I_1, \dots, I_n$  as parameters whose parameter space is  $\{0, 1\}$  where  $I_i = 1$  denotes that the  $i$ th observation has membership in the  $\text{Poisson}(\theta_0)$  distribution and  $I_i = 0$  denotes that the  $i$ th observation has membership in the  $\text{Poisson}(\theta_1)$  distribution. Now find the kernel of the posterior distribution.

- (f) [harder] Find the kernels of all four conditional distributions (for  $\theta_0, \theta_1, \rho, I_i$ ). If they are proportional to a known distribution, name it.

- (g) [easy] Is it possible to get inference for this model using a Gibbs Sampler after data augmentation? Why or why not?

### Problem 5

These are general questions about Permutation Testing.

- (a) [easy] What are the null and alternative hypotheses for a two-sample permutation test?
- (b) [easy] Let  $n_1$  and  $n_2$  be the sample sizes from population one and population two respectively. How many possible sample “permutations” are there? I put permutations in quotes because it’s not truly a “permutation” in the sense that you were taught in MATH 241.



(c) [easy] Give three examples of a test statistic to employ within the body of the loop of a permutation test.

(d) [difficult] Explain how you would calculate a p-value in a permutation test.

### **Problem 6**

These are general questions about the Bootstrap. Assume  $X_1, \dots, X_n \stackrel{iid}{\sim}$  some DGP.

(a) [easy] Describe the steps in the bootstrap procedure for the estimate  $\hat{\theta} := w(x_1, \dots, x_n)$  which estimates  $\theta$ .

(b) [easy] In what situations should the bootstrap be employed instead of other inferential procedures you learned about?

- (c) [difficult] Explain in what situations the bootstrap fails. Read online about this.

### Problem 7

These are questions about parametric survival using the Weibull model i.e.

$$Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Weibull}(k, \lambda) := f(y) = k\lambda^k y^{k-1} e^{-\lambda^k y^k} \mathbb{1}_{y>0}, \quad F(y) = 1 - e^{-\lambda^k y^k}, \quad S(y) = e^{-\lambda^k y^k}$$

- (a) [difficult] Assume no censoring in the data. Find closed form expressions and/or equations for the MLEs of  $k$  and  $\lambda$

- (b) [easy] Assume censoring in the data so that  $\mathbf{c}$  is the binary vector that is zero when censored and one if measured. Let  $\mathbf{y}$  be the vector of measurements or censored values if not measured. Find  $\ell(k, \lambda; \mathbf{y}, \mathbf{c})$ .

## Problem 8

These are questions about nonparametric survival inference.

- (a) [easy] Show that the empirical survival function is equal to the product limit estimator form with no censoring. Make sure to define what your notation means.
- (b) [easy] Consider the dataset  $y = \{79, 81, 92, 95, 105, 107, 122\}$  measured in days. Draw the estimate of  $S(y)$ .
- (c) [harder] Let your parameter of interest  $\theta$  be survival past 106 days. Compute a 95% CI for  $\theta$ .

(d) [harder] Test  $H_a : \theta > 0.5$ .

(e) [easy] Explain how you would use the bootstrap to find a CI for the median. Explain why the bootstrap won't be so accurate in this example.

(f) [harder] Rederive the Kaplan-Meier estimator for the survival function.

(g) [harder] Consider the dataset  $y = \{79, 81, 92+, 95, 105+, 107, 122\}$  measured in days where the “+” signs indicate censored values. Draw the Kaplan-Meier estimate of  $S(y)$  in a different color atop the estimate in (b). Try to make it to scale as best as possible.

(h) [harder] Explain how you would use the bootstrap to find a CI for the median.

(i) [easy] Write the hypotheses for the log-rank test.

(j) [easy] Write the formula for the test statistic in the log-rank test.