## Final Exam [30 pts]

Monday, May 13th, 2019 1:30-3:30 pm

PUBH 8442: Bayes Decision Theory and Data Analysis

Give your final answers in simplified, closed form wherever possible. However, partial credit will be awarded for incomplete solutions. Good luck!

1. Linear Model with Measurement Error [18 pts]

In what follows, assume that all variables are independent unless otherwise specified. Consider a linear model with measurement error in the predictors. Data  $(X_i, Y_i)$  are observed, where

$$Y_i = \beta X_i + \epsilon_i,$$
  

$$\epsilon_i \sim N(0, \sigma^2), \text{ and}$$
  

$$\tilde{X}_i \mid X_i \sim N(X_i, \tau^2) \text{ for } i = 1, \dots, n.$$

Assume  $\sigma^2$  is known. We use an improper flat prior for  $\beta : p(\beta) = 1 \forall \beta \in \mathbb{R}$ , and in parts (a-d) we use an inverse-gamma prior for  $\tau^2$ , IG(a, b):

$$p(\tau^2) = \frac{b^a}{\Gamma(a)} (\tau^2)^{-(a+1)} e^{-(b/\tau^2)}.$$

Let  $\mathbf{Y} = (Y_1, \ldots, Y_n)$ ,  $\mathbf{X} = (X_1, \ldots, X_n)$ , and  $\mathbf{\tilde{X}} = (\tilde{X}_1, \ldots, \tilde{X}_n)$ .

(a) (3 points) What is the marginal distribution of **Y** given **X**,  $\beta$  and  $\tau^2$ :  $p(\mathbf{Y} \mid \beta, \tau^2, \mathbf{X})$ ?

Each  $Y_i$  is normally distributed with mean

$$E(Y_i) = \beta E(X_i) + E(\epsilon_i) = \beta X_i$$

and variance

$$\operatorname{Var}(Y_i) = \beta^2 \operatorname{Var}(\tilde{X}_i) + \operatorname{Var}(\epsilon_i)) = \beta^2 \tau^2 + \sigma^2.$$

Their joint marginal distribution is  $\mathbf{Y} \sim \text{MVN} \left(\beta \mathbf{X}, (\beta^2 \tau^2 + \sigma^2) I_{n \times n}\right).$ 

(b) (3 points) What is the conditional posterior distribution of  $\beta$ ,  $p(\beta \mid \tilde{\mathbf{X}}, \tau^2, \mathbf{X}, \mathbf{Y})$ ? Note that  $\tilde{\mathbf{X}}\beta \sim MVN(\mathbf{Y}, \sigma^2 I_{n \times n})$ , and multiplying both sides by  $(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T$  yields

$$\beta \sim MVN((\tilde{\mathbf{X}}^T\tilde{\mathbf{X}})^{-1}\tilde{\mathbf{X}}^T\mathbf{Y}, (\tilde{\mathbf{X}}^T\tilde{\mathbf{X}})^{-1}\sigma^2)$$

(c) (3 points) What is the conditional posterior distribution of  $\tau^2$ ,  $p(\tau^2 | \tilde{\mathbf{X}}, \beta, \mathbf{X}, \mathbf{Y})$ ? Note that  $X_i - \tilde{X}_i \sim N(0, \tau^2)$ , so by the normal-inverse-gamma model, the conditional posterior for  $\tau^2$  is  $\mathrm{IG}\left(a + \frac{n}{2}, \beta + \frac{1}{2}\sum(X_i - \tilde{X}_i)^2\right)$ . This may also be derived directly, using the density functions:

$$p(\tau^{2} \mid \tilde{\mathbf{X}}, \beta, \mathbf{X}, \mathbf{Y}) \propto IG(\tau^{2}; a, b) \prod_{i=1}^{n} N(\tilde{X}_{i}; X_{i}, \tau^{2})$$

$$\propto (\tau^{2})^{-(a+1)} e^{-(b/\tau^{2})} \cdot (\tau^{2})^{-n/2} e^{-1/(2\tau^{2})\sum(X_{i} - \tilde{X}_{i})^{2}}$$

$$\propto (\tau^{2})^{-(a+n/2+1)} e^{-(1/\tau^{2})(b+(1/2)\sum(X_{i} - \tilde{X}_{i})^{2})}$$

$$\propto IG(a + \frac{n}{2}, b + \frac{1}{2}\sum(X_{i} - \tilde{X}_{i})^{2}).$$

- (d) (5 points) Describe explicitly a Gibbs sampling algorithm to simulate from the joint posterior distribution  $p(\beta, \tau^2, \tilde{\mathbf{X}} | \mathbf{X}, \mathbf{Y})$ . Initialize  $\beta^{(0)}, \tau^{2^{(0)}}$ . Then, for t = 1, ..., T:
  - i. Generate  $\tilde{\mathbf{X}}^{(t)}$  from  $p(\tilde{\mathbf{X}} \mid \beta^{(t-1)}, \tau^{2^{(t-1)}}, \mathbf{X}, \mathbf{Y})$
  - ii. Generate  $\beta^{(t)}$  from  $p(\beta \mid \tilde{\mathbf{X}}^{(t)}, \tau^{2^{(t-1)}}, \mathbf{X}, \mathbf{Y})$
  - iii. Generate  $\tau^{2^{(t)}}$  from  $p(\tau \mid \tilde{\mathbf{X}}^{(t)}, \beta(t-1), \mathbf{X}, \mathbf{Y})$

Step (2) uses the solution to part b, and step (3) uses the solution to part c. For step (1), because  $\tilde{X}_i \sim N(X_i, \tau^2)$  and  $Y_i/\beta \sim N(\tilde{X}_i, \sigma^2/\beta^2)$ , an application of the normal-normal model gives the conditional posterior  $\tilde{X}_i \sim N(\mu_i, \text{Var}_i)$ where

$$\mu_i = \frac{(\sigma^2/\beta^2)X_i + \tau^2(Y_i/\beta)}{(\sigma^2/\beta^2) + \tau^2}$$

and

$$\operatorname{Var}_{i} = \frac{(\sigma^{2}/\beta^{2})\tau^{2}}{(\sigma^{2}/\beta^{2}) + \tau^{2}}$$

for i = 1, ..., n.

(e) (4 points) Instead of using an IG(*a*, *b*) prior for  $\tau^2$ , describe a reasonable empirical Bayes approach to estimate  $\tau^2$ . Write your answer in closed form as a function of the observed data,  $\hat{\tau}^2 = f(\mathbf{X}, \mathbf{Y}, \sigma^2)$ . From part (a),  $\operatorname{var}(Y_i) = \beta^2 \tau^2 + \sigma^2$ . Substitute  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$  and let  $s^2$  be the empirical variance of the residuals  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - X_i \hat{\beta})^2$ . Then, solve  $s^2 = \hat{\beta}^2 \tau^2 + \sigma^2$  to get  $\hat{\tau}^2 = \frac{s^2 - \sigma^2}{\hat{\beta}^2}$ .

There may be sensible alternative approaches.



Figure 1

2. Importance Sampling [4 pts]

Figures 1 illustrate importance sampling from a N(0, 1) distribution. The left column gives four different importance densities (dashed line) that were used to draw importance samples from the target distribution (black line). The right column shows histograms of the sample weights generated from each scenario on the left. Match the weight histogram on the right (A,B,C, or D) that corresponds to each panel on the left (I, II, III, or IV). I: A, II: C, III: D, IV: B 3. Genomic testing [8 pts]

Consider genetic data for  $N_1$  individuals that are affected by a disease, and  $N_0$  unaffected individuals. A single neucleotide polymorphism (SNP) is recorded as present (Y = 1) or absent (Y = 0) for each individual. Thus, the data are of the form  $Y_{ij} \in \{0, 1\}$  for affected/unaffected status j = 0, 1, and individuals  $i = 1, \ldots, N_j$ . Assume

$$Y_{ij} \stackrel{\text{indep}}{\sim} \text{Bernoulli}(\theta_j)$$

so  $\theta_0$  is the probability that the SNP is present for an unaffected individual and  $\theta_1$  is the probability that the SNP is present for an affected individual. Consider two models, one where disease status has no effect  $(M_0)$  and one where it does  $(M_a)$ :

- $M_0: \theta_0 = \theta_1 = \theta$  and  $\theta \sim \text{Beta}(1, 1)$ ,
- $M_a: \theta_0 \sim \text{Beta}(1,1) \text{ and } \theta_1 \sim \text{Beta}(1,1), \text{ where } \theta_0 \text{ and } \theta_1 \text{ are independent.}$
- (a) (2 pts) What are the posterior distributions of  $\theta_0$  and  $\theta_1$  under  $M_a$ ,  $p(\theta_0 | \mathbf{Y}, M_a)$  and  $p(\theta_1 | \mathbf{Y}, M_a)$ ? By the beta-binomial model,  $p(\theta_0 | \mathbf{Y}, M_a) = \text{Beta} \left( 1 + \sum_{i=1}^{N_0} Y_{i0}, 1 + N_0 - \sum_{i=1}^{N_0} Y_{i0} \right)$ and  $p(\theta_1 | \mathbf{Y}, M_a) = \text{Beta} \left( 1 + \sum_{i=1}^{N_1} Y_{i1}, 1 + N_1 - \sum_{i=1}^{N_1} Y_{i1} \right)$ .
- (b) (2 pts) What is the posterior distribution of  $\theta_0$ , under  $M_0$ ,  $p(\theta_0 | \mathbf{Y}, M_0)$ ? By the beta-binomial model and because under  $M_0 \ \theta_0 = \theta_1 = \theta$ ,

$$p(\theta_0 \mid \mathbf{Y}, M_0) = \text{Beta}\left(\left(1 + \sum_{i=1}^{N_0} Y_{i0} + \sum_{i=1}^{N_1} Y_{i1}, 1 + N_0 + N_1 - \sum_{i=1}^{N_0} Y_{i0} - \sum_{i=1}^{N_1} Y_{i1}\right)\right)$$

(c) (4 pts) Assume  $S_0$  unaffected individuals have the SNP and  $S_1$  affected have the SNP. What is the Bayes factor for  $M_0$  over  $M_a$ ? The Bayes factor is  $\frac{p(\mathbf{Y}|M_0)}{p(\mathbf{Y}|M_a)}$ , where

$$p(\mathbf{Y} \mid M_0) = \int p(\mathbf{Y} \mid \theta) p(\theta) d\theta$$
  
=  $\int \theta^{S_0 + S_1} (1 - \theta)^{N_0 + N_1 - S_0 - S_1} d\theta$   
=  $B(S_0 + S_1 + 1, N_0 + N_1 - S_0 - S_1 + 1)$ 

and

$$p(\mathbf{Y} \mid M_a) = \int p(\mathbf{Y} \mid \theta_0) p(\theta_0) d\theta_0 \int p(\mathbf{Y} \mid \theta_1) p(\theta_1) d\theta_1$$
  
=  $\int \theta_0^{S_0} (1 - \theta_0)^{N_0 - S_0} d\theta_0 \int \theta_1^{S_1} (1 - \theta_1)^{N_1 - S_1} d\theta_1$   
=  $B(S_0 + 1, N_0 + 1) B(S_1 + 1, N_1 + 1).$