

## DS-GA 3001: Applied Statistics Practice Midterm

### Instructions:

- You have **100 minutes**, 4:55PM - 6:35PM
- The exam has 4 problems, totaling 100 points.
- Please answer each problem in the space below it.
- You are allowed to carry the textbook, your own notes and other course related material with you. Electronic devices are not allowed.
- Please read the problems carefully.
- We use boldcase letters  $\theta, \mathbf{x}, \dots$  to distinguish vectors from scalars.
- Unless otherwise specified, you are required to provide explanations of how you arrived at your answers.
- You can use previous parts of a problem even if you did not solve them.
- The problems may not be arranged in an increasing order of difficulty. If you get stuck, it might be wise to try other problems first.
- Good luck and enjoy!

Full name: \_\_\_\_\_

N number: \_\_\_\_\_

**1. Binary choice questions.** (40 points)

For each of the statements, decide if it is “True” or “False”. Provide explanations if you think it is “False”. Each question is worth 5 points.

- (a) For  $\theta \in \mathbb{R}$ , let  $y \sim p_\theta$  denote the distribution where  $y$  is uniformly distributed on the interval  $[\theta, \theta + 1]$ . This family  $(p_\theta)_{\theta \in \mathbb{R}}$  is an exponential family.

- (b) Let  $\mathbf{y} = (y_1, \dots, y_n)$  be a sample of  $n$  i.i.d. observations from  $p_\theta$ , and  $D_n(\theta_1; \theta_2)$  be the deviance between two parameters  $\theta_1, \theta_2 \in \mathbb{R}$  based on  $\mathbf{y}$ . Then  $D_n(\theta_1; \theta_2) = nD_1(\theta_1; \theta_2)$ , where  $D_1(\theta_1, \theta_2)$  is the deviance based on a single observation  $y_1$ .

- (c) Let  $P$  be an unknown continuous distribution over  $\mathbb{R}$ . Given i.i.d.  $Y_1, \dots, Y_n \sim P$ , Alice computes the following statistic

$$T = T(Y_1, \dots, Y_n) = \text{number of distinct values in } (Y_1, \dots, Y_n).$$

For example,  $T(1, 2, 3) = 3$ , and  $T(1.4, 1, 1.4) = 2$ . Alice would like to estimate the variance of  $T$  via bootstrap: she draws  $m$  bootstrap samples  $\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(m)}$ , where each sample  $\mathbf{Y}^{(j)}$  is a collection of  $n$  uniformly random draws (with replacement) from  $\{Y_1, \dots, Y_n\}$ . Alice proceeds to compute  $T^{(j)} = T(\mathbf{Y}^{(j)})$ , and uses the sample variance of  $(T^{(j)} : j = 1, \dots, m)$  to estimate the true variance of  $T$ .

Claim: for large  $(m, n)$ , this bootstrap estimate is close to the true variance of  $T$ .

- (d) In a GLM with parameter  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$ , Alice would like to test if  $\beta_1 = \beta_2 = \dots = \beta_p$ . She runs the following procedure:

- compute the unrestricted MLE  $\hat{\boldsymbol{\beta}}^{(1)}$  and the corresponding log-likelihood  $\ell_1$ ;
- compute the restricted MLE  $\hat{\boldsymbol{\beta}}^{(2)}$  subject to the constraint  $\hat{\beta}_1^{(2)} = \hat{\beta}_2^{(2)} = \dots = \hat{\beta}_p^{(2)}$ , and compute the corresponding log-likelihood  $\ell_2$ .

She then claims that under the null hypothesis  $\beta_1 = \beta_2 = \dots = \beta_p$ , asymptotically one should have  $2(\ell_1 - \ell_2) \sim \chi_{p-1}^2$ . Is this claim correct?

- (e) For model selection, intuitively speaking AIC aims to balance between two terms:
- the negative log-likelihood, which shrinks with an increasing model complexity;
  - the number of model parameters, which grows with an increasing model complexity.
- (f) Suppose  $\mathbf{D}_1$  and  $\mathbf{D}_2$  are two survival datasets for males and females, respectively. Then the following ways to plot the survival curves are equivalent:
- plot the Kaplan-Meier curves for males and females, respectively;
  - fit the Cox model on  $\mathbf{D}_1 \cup \mathbf{D}_2$  with the feature “gender”, then plot the fitted survival curves for males and females, respectively.

- (g) Recall that in the Cox model, the complete likelihood is  $L(\boldsymbol{\beta}, h)$  and the profile likelihood is  $pL(\boldsymbol{\beta}) = \max_h L(\boldsymbol{\beta}, h)$ . Bob claims that computing the profile maximum likelihood  $\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta}} pL(\boldsymbol{\beta})$  is equivalent to one single iteration of the EM algorithm, where one first computes  $\hat{h}$  and then computes  $\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta}} L(\boldsymbol{\beta}, \hat{h})$ . Is this claim correct?

- (h) For exponential families with missing data, the incomplete log-likelihood may *no longer* be concave in the natural parameter  $\boldsymbol{\theta}$ .

**2. EM algorithm with both covariates missing.** *(20 points)*

Let the sample  $(x_1, y_1), \dots, (x_{n+m}, y_{n+m})$  be i.i.d. drawn from an exponential family  $p_\theta(x, y) = \exp(\langle \theta, T(x, y) \rangle - A(\theta))h(x, y)$ . However, the  $y$ -covariate is missing in the first  $n$  observations, and the  $x$ -covariate is missing in the last  $m$  observations. In other words, our observed sample is  $(x_1, \dots, x_n, y_{n+1}, \dots, y_{n+m})$ .

- (a) Write out the incomplete log-likelihood for the observations (up to additive constants), in terms of  $A(\theta)$  and its conditional variants. *(10 points)*

- (b) Describe the EM algorithm for the MLE computation. You should give the details of both E and M steps; you need not give proofs. *(10 points)*

### 3. Mixture model with known locations. (20 points)

In a mixture model, we have a dataset  $(y_1, \dots, y_n)$  with  $y_i \sim p_{\theta_i}$ , where the unknown parameters  $\theta_1, \dots, \theta_n$  are i.i.d. drawn from an unknown distribution  $\pi$ . We can think of  $\theta_i$  as the “locations” of the mixture, and the vector  $\pi$  as the “weights”.

Throughout this problem we assume that both  $\theta_i$  and  $y_i$  take discrete values, i.e.

$$\begin{aligned}\theta_i &\in \Theta = \{\theta^1, \dots, \theta^M\}, \\ y_i &\in \mathcal{Y} = \{y^1, \dots, y^N\}.\end{aligned}$$

Therefore, we may represent  $\pi$  as a probability vector  $(\pi_1, \dots, \pi_M)$ , in the sense that  $\mathbb{P}(\theta = \theta^j) = \pi_j$  if  $\theta \sim \pi$ . We will also use the notation  $K_{j\ell} = p_{\theta^j}(y^\ell) = \mathbb{P}(y = y^\ell \mid \theta = \theta^j)$  to denote the conditional probability of observing  $y = y^\ell$  when  $\theta = \theta^j$ .

- (a) If  $\theta \sim \pi$  and  $y \sim p_\theta$ , write down the marginal distribution of  $y$  in terms of  $(\pi, K)$ .  
(5 points)

- (b) Write down the log-likelihood of the dataset  $(y_1, \dots, y_n)$ , as a function of  $\pi$ ; we assume that  $\Theta, \mathcal{Y}, K$  are known. Is the log-likelihood concave in  $\pi$ ? (5 points)

(c) Now suppose that  $\pi_j$  takes a form of a one-dimensional exponential family, i.e.

$$\pi_j = \exp(\beta T_j - A(\beta)) h_j$$

for some known  $(T_j, h_j)$  and  $A(\cdot)$ . Is your log-likelihood in (b) concave in  $\beta$ ? (5 points)

(d) Now suppose that  $\Theta$  is unknown, so the matrix  $K = K(\Theta)$  becomes a function of  $\Theta$ . Is your log-likelihood in (b) jointly concave in  $\pi$  and  $\Theta$ ? (5 points)



**4. Survival analysis.** (20 points)

- (a) Consider a survival dataset  $\{(t_i, d_i, n_i)\}_{i=1}^N$ , where  $n_i$  is the number of individuals who have survived through time  $t_i$ , and  $d_i$  is the number of deaths at time  $t_i$ . For simplicity we assume that there is no censoring.

Recall that the Kaplan-Meier estimator for the hazard rate at each time  $t_i$  is

$$\hat{h}(t_i) = \frac{d_i}{n_i}.$$

We assume that  $d_i \sim B(n_i, h(t_i))$ , where  $B(n, p)$  denotes the binomial distribution with  $n$  trials and success probability  $p$ , and  $h(t_i)$  is the true hazard at  $t_i$ . Compute  $\text{Var}(\hat{h}(t_i))$ . (5 points)

(b) Recall that Kaplan-Meier estimator for the survival function is

$$\widehat{S}(t_i) = \prod_{j:t_j \leq t_i} (1 - \widehat{h}(t_j)).$$

Suppose we know that

$$\text{Var}(\log \widehat{S}(t_i)) \approx \sum_{j:t_j \leq t_i} \left( \frac{1}{1 - \widehat{h}(t_j)} \right)^2 \text{Var}(\widehat{h}(t_j)).$$

Find the approximate variance  $\text{Var}(\widehat{S}(t_i))$  using the delta method and the plug-in approach.

You should use your result in (a), and express your final answer using  $\widehat{S}(t_i)$  and  $\{(t_i, d_i, n_i)\}_{i=1}^N$ . (5 points)

- (c) Now consider a dataset  $\{(t_i, \Delta_i, \mathbf{x}_i)\}_{i=1}^N$  with features  $\mathbf{x}_i \in \mathbb{R}^d$ , true-death indicators  $\Delta_i \in \{0, 1\}$ , and distinct stopping times  $t_i$ . The Cox model assumes that

$$h(t \mid \mathbf{x}) = e^{\boldsymbol{\beta}^\top \mathbf{x}} h(t),$$

where  $\boldsymbol{\beta}$  independent of time.

Now suppose we would like to incorporate the time dependence by  $\boldsymbol{\beta}(t) = g(t)\boldsymbol{\beta}$ , or equivalently,

$$h(t \mid \mathbf{x}) = e^{g(t)\boldsymbol{\beta}^\top \mathbf{x}} h(t).$$

Write out the partial likelihood you will use to estimate  $\boldsymbol{\beta}$ . You may assume that  $g$  is a known function. (5 points)

- (d) Propose a model for  $h(t \mid \mathbf{x})$  if the features  $\mathbf{x}$  are assumed to be time-dependent, i.e.  $\mathbf{x} = \mathbf{x}(t)$ . Do you think it is helpful to include both time-dependent features  $\mathbf{x}(t)$  and the time-dependent coefficient  $g(t)$  in (c)? (5 points)