

EGT2
ENGINEERING TRIPOS PART IIA

Wednesday 26 April 2023 14.00 to 15.40

Module 3F8

CRIB: INFERENCE

*Answer not more than **three** questions.*

All questions carry the same number of marks.

*The **approximate** percentage of marks allocated to each part of a question is indicated in the right margin.*

*Write your candidate number **not** your name on the cover sheet.*

STATIONERY REQUIREMENTS

Single-sided script paper

SPECIAL REQUIREMENTS TO BE SUPPLIED FOR THIS EXAM

CUED approved calculator allowed

Engineering Data Book

10 minutes reading time is allowed for this paper at the start of the exam.

You may not start to read the questions printed on the subsequent pages of this question paper until instructed to do so.

You may not remove any stationery from the Examination Room.

1 A data scientist uses a generative classifier to predict the binary label $y \in \{0, 1\}$ from the scalar feature $x \in \mathbb{R}$. The forms of $p(y)$ and $p(x|y)$ are known and given by

$$p(y) = \begin{cases} 0.4 & \text{if } y = 1 \\ 0.6 & \text{if } y = 0 \end{cases}, \quad p(x|y) = \mathcal{N}(x|m_y, 1) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x - m_y)^2}{2}\right\},$$

where m_0 and m_1 are two scalar mean parameters.

(a) Assuming that m_0 and m_1 are known, write down an expression for the classifier's predictive probability $p(y = 1|x)$. Write your expression using the notation $\mathcal{N}(x|m, v)$ to represent a Gaussian density with mean m and variance v evaluated at x . [25%]

(b) The data scientist introduces the following prior for m_0 and m_1 :

$$p(m_0, m_1) = \mathcal{N}(m_0|0, 1) \times \mathcal{N}(m_1|0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{m_0^2}{2}\right\} \times \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{m_1^2}{2}\right\}.$$

Given a data point (x_n, y_n) , write down an expression for the posterior distribution of m_0 and m_1 as a product of Gaussian densities using the notation $\mathcal{N}(x|m, v)$ as before. Do not calculate the value of the normalisation constant for the posterior and write Z instead. [25%]

(c) The data scientist observes a dataset $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$ containing N data points. According to the generative classifier, the log-likelihood for m_0 and m_1 is

$$L(m_0, m_1) = \sum_{n=1}^N -I[y_n = 0] \times \frac{1}{2}(x_n - m_0)^2 - \sum_{n=1}^N I[y_n = 1] \times \frac{1}{2}(x_n - m_1)^2 + \text{constant},$$

where $I[\cdot]$ is equal to 1 if its input is true and 0 otherwise. Using this result, write down an expression for the maximum a posteriori estimates of m_0 and m_1 given \mathcal{D} . [25%]

(d) The maximum a posteriori estimates of m_0 and m_1 are given by $\hat{m}_0 = 1$ and $\hat{m}_1 = -1$. The data scientist then uses the classifier to make predictions using these estimates and the following loss (negative reward) function:

$$L(y, \hat{y}) = \begin{cases} 0 & \text{if } y = 0 \text{ and } \hat{y} = 0 \\ 2 & \text{if } y = 0 \text{ and } \hat{y} = 1 \\ 1 & \text{if } y = 1 \text{ and } \hat{y} = 0 \\ 0 & \text{if } y = 1 \text{ and } \hat{y} = 1 \end{cases},$$

where y is the true class label and \hat{y} is the classifier's prediction. What is the classifier's decision boundary in this case? Hint: The decision boundary is the input x at which predicting one class or the other yields the same expected loss. [25%]

(a) Using Bayes rule we obtain

$$p(y = 1|x) = \frac{p(x|y = 1)p(y = 1)}{p(x)} = \frac{\mathcal{N}(x|m_1, 1)0.4}{\mathcal{N}(x|m_0, 1)0.6 + \mathcal{N}(x|m_1, 1)0.4}. \quad (1)$$

[25%]

(b) The posterior is given by the product of the likelihood and the prior. In particular,

$$\begin{aligned} p(m_0, m_1|x_n, y_n) &\propto p(x_n|y_n, m_{y_n})p(y_n)p(m_0)p(m_1) \\ &\propto \mathcal{N}(x_n|m_{y_n}, 1)\mathcal{N}(m_0|0, 1)\mathcal{N}(m_1|0, 1). \end{aligned} \quad (2)$$

Therefore,

$$p(m_0, m_1|x_n, y_n) = \frac{\mathcal{N}(x_n|m_{y_n}, 1)\mathcal{N}(m_0|0, 1)\mathcal{N}(m_1|0, 1)}{Z}. \quad (3)$$

[25%]

(c) The log posterior for m_0 and m_1 is given by,

$$\begin{aligned} LP(m_0, m_1) &= \sum_{n=1}^N -I[y_n = 0] \frac{1}{2}(x_n - m_0)^2 - \sum_{n=1}^N I[y_n = 1] \frac{1}{2}(x_n - m_1)^2 \\ &\quad - \frac{m_0^2}{2} - \frac{m_1^2}{2} + \text{constant}. \end{aligned} \quad (4)$$

The quantity above is maximized when the gradient with respect to m_0 and m_1 is zero.

That is,

$$\hat{m}_0 = \frac{\sum_{n=1}^N I[y_n = 0]x_n}{\sum_{n=1}^N I[y_n = 0] + 1}, \quad \hat{m}_1 = \frac{\sum_{n=1}^N I[y_n = 1]x_n}{\sum_{n=1}^N I[y_n = 1] + 1}, \quad (5)$$

[25%]

(d) The classifier's decision boundary is the input x at which predicting one class or the other yields the same loss. That is,

$$\begin{aligned} \sum_y L(y, \hat{y} = 1)p(y|x) &= \sum_y L(y, \hat{y} = 0)p(y|x) \\ L(y = 0, \hat{y} = 1)p(y = 0|x) &= L(y = 1, \hat{y} = 0)p(y = 1|x) \\ \frac{L(y = 0, \hat{y} = 1)}{L(y = 1, \hat{y} = 0)} &= \frac{p(y = 1|x)}{p(y = 0|x)} \\ 2 &= \frac{\mathcal{N}(x|\hat{m}_1, 1)0.4}{\mathcal{N}(x|\hat{m}_0, 1)0.6} \end{aligned}$$

$$3 = \exp \left\{ -0.5(x - \hat{m}_1)^2 + 0.5(x - \hat{m}_0)^2 \right\}$$
$$x \approx -0.5986.$$

[25%]

2 A product of experts model makes predictions for a quantity of interest by multiplying the predictive densities of individual experts. Consider a regression problem and a total of E experts given by linear models with additive Gaussian noise, parameter vectors $\mathbf{w}_1, \dots, \mathbf{w}_E$ and non-linear feature transformation functions $\phi_1(\cdot), \dots, \phi_E(\cdot)$, respectively. In this case, the product of experts predictions are given by

$$p(y_n|\mathbf{x}_n) \propto \prod_{e=1}^E \mathcal{N}(y_n|\mathbf{w}_e^T \phi_e(\mathbf{x}_n), \sigma_e^2),$$

where σ_e^2 is the variance of the zero-mean additive Gaussian noise assumed by expert e .

(a) It can be shown that $p(y_n|\mathbf{x}_n)$ is Gaussian:

$$p(y_n|\mathbf{x}_n) = \mathcal{N}(x|m, v) = \frac{1}{\sqrt{2\pi v}} \exp\left(-\frac{(x-m)^2}{2v}\right).$$

Give an expression for m and v in terms of the means and variances of the individual experts, that is, $\mathbf{w}_e^T \phi_e(\mathbf{x}_n)$ and σ_e^2 , $e = 1, \dots, E$. For this, use the following property of the product of Gaussian densities: $\mathcal{N}(x|m_1, v_1) \times \mathcal{N}(x|m_2, v_2) \propto \mathcal{N}(x|m_3, v_3)$, where $v_3^{-1} = v_1^{-1} + v_2^{-1}$ and $m_3/v_3 = m_1/v_1 + m_2/v_2$. Will v increase or decrease as one more expert is added into the model? [25%]

(b) Under some specific choice of transformation functions ϕ_1, \dots, ϕ_E , the training error of a product of experts model that is fitted to the data by maximum likelihood can be made arbitrarily small by just increasing the number of experts in the model. However, using a very large number of experts may not be a good idea in practice. Why? [25%]

(c) Consider now an alternative mixture of experts model that uses the same experts as the model above and has uniform mixing weights. The predictive distribution is now

$$p(y_n|\mathbf{x}_n) = \frac{1}{E} \left\{ \sum_{e=1}^E \mathcal{N}(y_n|\mathbf{w}_e^T \phi_e(\mathbf{x}_n), \sigma_e^2) \right\}.$$

What are the mean and variance of this predictive distribution? [25%]

(d) Jensen's inequality says that for a convex function φ , numbers x_1, \dots, x_n in its domain, and positive weights a_i , $\varphi(\sum_i a_i x_i / \sum_i a_i) \leq \sum_i a_i \varphi(x_i) / \sum_i a_i$. Using this and the fact that the quadratic function is convex, or otherwise, show that the variance of the predictive distribution for the previous mixture of experts model has to be greater than or equal to the minimum of $\sigma_1^2, \dots, \sigma_E^2$. Using this fact and the response to part (a) above, justify why the product of experts model may be preferred over the mixture of experts one. [25%]

(a) The mean of the e -th Gaussian in our product is $m_e = \mathbf{w}_e^T \phi_e(\mathbf{x}_n)$. The variance is $v_e = \sigma_e^2$. Let $\alpha = \sum_e v_e^{-1}$ and let $\beta = \sum_e m_e / v_e$. The product of Gaussians is then

$$\prod_{e=1}^E \mathcal{N}(y_n | m_e, v_e) \propto \mathcal{N}(y_n, \beta \alpha^{-1}, \alpha^{-1}). \quad (6)$$

The product is a Gaussian density with mean $m = \beta \alpha^{-1}$ and variance $v = \alpha^{-1}$. When you add one more expert into the mixture, the variance $v = \alpha^{-1}$ will decrease, since $\alpha = \sum_e v_e^{-1}$ will increase as we add one more positive element into this sum. [25%]

(b) The training error will be made arbitrarily small by increasing the number of experts but this could cause overfitting problems as the model fits noise in the training data and the test error increases beyond an optimal value obtained with a smaller number of experts. [25%]

(c) The mean of the predictive distribution is given by

$$\mathbf{E}_{p(y_n | \mathbf{x}_n)}[y_n] = \frac{1}{E} \sum_e \mathbf{w}_e^T \phi_e(\mathbf{x}_n). \quad (7)$$

The second moment is given by

$$\mathbf{E}_{p(y_n | \mathbf{x}_n)}[y_n^2] = \frac{1}{E} \sum_e \left((\mathbf{w}_e^T \phi_e(\mathbf{x}_n))^2 + \sigma_e^2 \right). \quad (8)$$

The variance is, therefore, given by

$$\text{VAR}_{p(y_n | \mathbf{x}_n)}[y_n] = \frac{1}{E} \sum_e \left((\mathbf{w}_e^T \phi_e(\mathbf{x}_n))^2 + \sigma_e^2 \right) - \left(\frac{1}{E} \sum_e \mathbf{w}_e^T \phi_e(\mathbf{x}_n) \right)^2. \quad (9)$$

[25%]

(d) Let's denote by m_e the mean of the predictions by expert e . Then

$$\begin{aligned} \text{VAR}_{p(y_n | \mathbf{x}_n)}^E[y_n] &= \frac{1}{E} \sum_e m_e^2 + \frac{1}{E} \sum_e \sigma_e^2 - \left(\frac{1}{E} \sum_e m_e \right)^2 \\ &\geq \frac{1}{E} \sum_e \sigma_e^2 \geq \min\{\sigma_1^2, \dots, \sigma_E^2\}, \end{aligned} \quad (10)$$

where we have used that

$$\frac{1}{E} \sum_e m_e^2 \geq \left(\frac{1}{E} \sum_e m_e \right)^2, \quad (11)$$

which is obtained by using Jensen's inequality and the fact that the quadratic function is convex. The product of experts is preferred because adding new experts always makes the predictive distribution sharper. This does not happen with the mixture of experts as the predictive variance can never be smaller than the minimum predictive variance of one of its experts.

[25%]

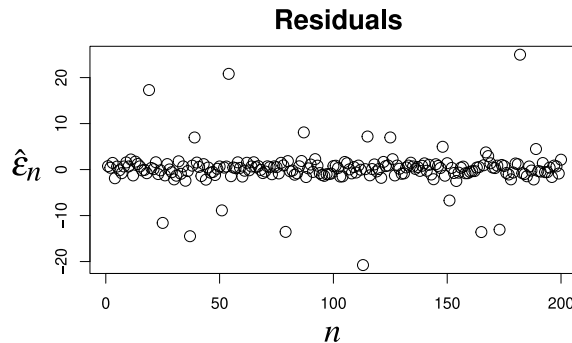
3 Consider a regression dataset $\{(x_n, y_n)\}_{n=1}^N$ and a linear regression model with coefficients $\boldsymbol{\theta} \in \mathbb{R}^d$ and Gaussian additive noise with variance σ^2 . Let $\boldsymbol{\Omega} = \{\boldsymbol{\theta}, \sigma^2\}$ be the set of model parameters. The likelihood for the n -th data point is then given by

$$p(y_n|x_n, \boldsymbol{\Omega}) = \mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma^2), \quad \text{where} \quad f(x_n; \boldsymbol{\theta}) = \sum_{i=1}^d \phi_i(x_n)\theta_i$$

with the ϕ_i being non-linear basis functions. The model is fitted with respect to both $\boldsymbol{\theta}$ and σ^2 using maximum likelihood, giving the following estimates for these parameters:

$$\hat{\boldsymbol{\theta}}_{\text{ML}} = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}, \quad \hat{\sigma}_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (y_n - f(x_n; \hat{\boldsymbol{\theta}}_{\text{ML}}))^2,$$

where $\boldsymbol{\Phi}$ is the $N \times d$ matrix given by $[\boldsymbol{\Phi}]_{ni} = \phi_i(x_n)$ and \mathbf{y} is an N -dimensional vector containing all the y_n . After fitting the model, the training residuals $\hat{\epsilon}_n = y_n - f(x_n; \hat{\boldsymbol{\theta}}_{\text{ML}})$ are calculated and it is noticed that a small proportion of them have much larger magnitude than the others as shown in the following plot:



- (a) Explain the negative effect that these large residuals will have in $\hat{\boldsymbol{\theta}}_{\text{ML}}$ and $\hat{\sigma}_{\text{ML}}^2$ as compared to the case when no such large residuals are present in the data. [25%]
- (b) To alleviate the problems caused by the large residuals, a latent binary variable h_n is introduced for each data point to distinguish between regular magnitude ($h_n = 0$) and large magnitude ($h_n = 1$) residual data points. The same f is used for both types of data points. Write an expression for the new likelihood $p(y_n|x_n, \boldsymbol{\Omega}, h_n)$ and a probability distribution for the prior $p(h_n)$. What parameters, in addition to $\boldsymbol{\theta}$, does this new model have? [25%]
- (c) The new model is fitted using the EM algorithm. State the E-step. [25%]
- (d) State the M-step for $\boldsymbol{\theta}$ when all the other model parameters are kept fixed. [25%]

(a) The large residuals may cause $\hat{\sigma}_{ML}^2$ to be too high for the majority of the data points with small residuals. At the same time, $\hat{\sigma}_{ML}^2$ may be too low for the few data points with large residuals. Regarding $\hat{\theta}_{ML}$, the large residuals may also increase the error in this estimate as these few data points will account for most of the squared errors of the model.

[25%]

(b) The new likelihood is

$$p(y_n|x_n, \mathbf{\Omega}, h_n) = (1 - h_n)\mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma_0^2) + h_n\mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma_1^2). \quad (12)$$

The prior is $p(h_n) = \pi^{1-h_n}(1 - \pi)^{h_n}$. In addition to $\boldsymbol{\theta}$, the model has now two noise variances σ_0^2 and σ_1^2 for regular and large magnitude residual data points and mixing proportion π indicating the prior probability that a point has a regular magnitude residual. That is, $\mathbf{\Omega} = \{\boldsymbol{\theta}, \sigma_0^2, \sigma_1^2, \pi\}$.

[25%]

(c) The E step calculates the posterior distribution of the binary latent variables, that is,

$$q(h_n) \propto p(y_n|x_n, \boldsymbol{\theta}, h_n)p(h_n). \quad (13)$$

In particular, we have that

$$\begin{aligned} q(h_n = 0) &= \frac{\pi\mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma_0^2)}{\pi\mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma_0^2) + (1 - \pi)\mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma_1^2)}, \\ q(h_n = 1) &= \frac{(1 - \pi)\mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma_1^2)}{\pi\mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma_0^2) + (1 - \pi)\mathcal{N}(y_n|f(x_n; \boldsymbol{\theta}), \sigma_1^2)}. \end{aligned} \quad (14)$$

[25%]

(d) The M step for $\boldsymbol{\theta}$ maximises the log of the likelihood times prior averaged with respect to the latents

$$\max_{\boldsymbol{\theta}} \sum_n \sum_{h_n} q(h_n) \log \{p(y_n|x_n, \mathbf{\Omega}, h_n)p(h_n)\}. \quad (15)$$

Taking derivative wrt to $\boldsymbol{\theta}$ and setting to zero to find the maximum

$$\sum_j (\boldsymbol{\Phi}^\top \mathbf{H}_j \mathbf{y} - \boldsymbol{\Phi}^\top \mathbf{H}_j \boldsymbol{\Phi} \boldsymbol{\theta}) / \sigma_j^2 = 0, \quad (16)$$

where \mathbf{H}_j is a diagonal matrix with entries in its diagonal given by $q(h_n = j)$ with $j = 0, 1$. The final solution is

$$\boldsymbol{\theta} = \left(\sum_j \boldsymbol{\Phi}^\top \mathbf{H}_j / \sigma_j^2 \boldsymbol{\Phi} \right)^{-1} \sum_j \boldsymbol{\Phi}^\top \mathbf{H}_j / \sigma_j^2 \mathbf{y}. \quad (17)$$

[25%]

4 A scalar linear dynamical system evolves according to

$$x_t = \lambda_1 x_{t-1} + \lambda_2 x_{t-2} + \sigma \epsilon_t, \quad (18)$$

where the ϵ_t are independent and identically distributed with $\epsilon_t \sim \mathcal{N}(0, 1)$.

(a) What is the Markov order of this system, and why? [20%]

(b) If λ_1 were set to be zero in equation (18), what would then be the relationship between states x_t at odd and even times? [20%]

(c) Re-write the system in equation (18) as a first order system in a new vector-valued variable \mathbf{z}_t such that

$$\mathbf{z}_t = \mathbf{A}\mathbf{z}_{t-1} + \mathbf{e}_t, \quad \text{where } \mathbf{e}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}).$$

In this expression, \mathbf{A} and \mathbf{Q} are matrices and \mathbf{z}_t and \mathbf{e}_t vectors. Write \mathbf{z}_t , \mathbf{A} and \mathbf{Q} in terms of quantities from equation (18). [20%]

(d) A dynamical system is unstable if its state grows without bound with time. What are the conditions of stability for the previous system in terms of the eigenvalues of \mathbf{A} ? [20%]

(e) An observation vector \mathbf{y}_t is obtained from \mathbf{z}_t as follows:

$$\mathbf{y}_t = \mathbf{B}\mathbf{z}_t + \mathbf{s}_t, \quad \text{where } \mathbf{s}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{T})$$

and \mathbf{B} and \mathbf{T} are matrices. Write down an expression for $p(\mathbf{y}_t | \mathbf{z}_{t-1})$. [20%]

(a) It's a second order system because the state depends on the two preceding states. [20%]

(b) If $\lambda_1 = 0$ then the states at odd times are completely independent of states at even times; there are effectively two independent dynamical systems. [20%]

(c) The new state variable has to be two dimensional $\mathbf{z}_t = (x_t, x_{t-1})^\top$. We have

$$\mathbf{A} = \begin{bmatrix} \lambda_1 & \lambda_2 \\ 1 & 0 \end{bmatrix}, \text{ and } \mathbf{Q} = \begin{bmatrix} \sigma^2 & 0 \\ 0 & 0 \end{bmatrix}, \quad (19)$$

which can be verified by writing out coordinate-wise. [20%]

(d) The state can only grow unboundedly if the matrix \mathbf{A} has eigenvalues with magnitude larger than 1. The condition for stability is that both eigenvalues are between -1 and 1. [20%]

(e) Because \mathbf{y}_t is obtained through a linear Gaussian model as a function of \mathbf{z}_t , we have that

$$p(\mathbf{y}_t | \mathbf{z}_{t-1}) = \int p(\mathbf{y}_t | \mathbf{z}_t) p(\mathbf{z}_t | \mathbf{z}_{t-1}) d\mathbf{z}_t. \quad (20)$$

To solve this integral we note that $\mathbf{z}_t | \mathbf{z}_{t-1} \sim \mathcal{N}(\mathbf{A}\mathbf{z}_{t-1}, \mathbf{Q})$. Since

$$\mathbf{y}_t = \mathbf{B}\mathbf{z}_t + \mathbf{s}_t, \quad \text{where } \mathbf{s}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{T}), \quad (21)$$

we then have that

$$\mathbf{y}_t | \mathbf{z}_{t-1} \sim \mathcal{N}(\mathbf{B}\mathbf{A}\mathbf{z}_{t-1}, \mathbf{B}\mathbf{Q}\mathbf{B}^\top + \mathbf{T}). \quad (22)$$

[20%]

END OF PAPER

THIS PAGE IS BLANK