

EGT2  
ENGINEERING TRIPOS PART IIA

---

Thursday 25 April 2019 14.00 to 15.40

---

**Module 3F8**

**INFERENCE**

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

*All questions carry the same number of marks.*

*The **approximate** number 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

Information 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.**

1 (a) Explain what is *maximum likelihood estimation* and how it is used to estimate parameters in a probabilistic model from data. [20%]

(b) A regression problem comprises scalar inputs  $x_n$  and scalar outputs  $y_n$  which are linearly related:  $y_n = mx_n + c + \varepsilon_n$ . The noise  $\varepsilon_n$  is Gaussian with mean 0 but with variance that depends quadratically on the input:  $p(\varepsilon_n) = \mathcal{N}(\varepsilon_n; 0, \alpha x_n^2)$ . Due to physical constraints, the inputs lie in the region  $1 < x_n < 100$ .

The offset  $c$  and noise parameter  $\alpha$  are known, but the slope  $m$  must be learned from a training dataset  $\{y_n, x_n\}_{n=1}^N$ .

(i) Compute the log-likelihood of  $m$ . [20%]

(ii) Compute the maximum likelihood estimate of  $m$ . [40%]

(iii) You are allowed to select a new input location  $x$  at which you will be provided with a corresponding output  $y$  and the new pair  $\{y, x\}$  will be added to the training data. Which value of  $x$  will be most informative about the parameter  $m$ ? Explain your reasoning. [20%]

2 A retailer has access to the collar size and waist size measurements of  $N$  customers. The retailer wants to use this information to estimate the relative size of each customer in order to recommend further products. The mean of the data is removed and the two measurements for each customer are stacked into a vector  $\mathbf{y}_n = [y_{1,n}, y_{2,n}]^\top$ . The full data set is denoted  $\{\mathbf{y}_n\}_{n=1}^N$ .

The retailer models the pair of measurements for each customer in terms of a scalar latent relative-size variable  $s_n$

$$\mathbf{y}_n = \begin{bmatrix} y_{1,n} \\ y_{2,n} \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} s_n + \sigma_y \begin{bmatrix} \varepsilon_{1,n} \\ \varepsilon_{2,n} \end{bmatrix} = \mathbf{w}s_n + \sigma_y \boldsymbol{\varepsilon}_n.$$

Here the weights  $\mathbf{w} = [w_1, w_2]^\top$  capture the dependence of the two measurements on relative-size and  $\boldsymbol{\varepsilon}_n = [\varepsilon_{1,n}, \varepsilon_{2,n}]^\top$  is independent and identically distributed Gaussian noise with mean zero and identity covariance  $p(\boldsymbol{\varepsilon}_n) = \mathcal{N}(\boldsymbol{\varepsilon}_n; \mathbf{0}, \mathbf{I})$ . Relative-size is assumed to follow a standard Gaussian distribution  $p(s_n) = \mathcal{N}(s_n; 0, 1)$ .

(a) The retailer wants to fit the parameters of the model  $\theta = \{\mathbf{w}, \sigma_y\}$  using maximum likelihood. Compute the likelihood of the parameters  $p(\{\mathbf{y}_n\}_{n=1}^N | \theta)$ . [30%]

(b) Having fit the model, the retailer would like to estimate the size of each customer. Compute the posterior distribution over the relative-size variable given the measurements  $p(s_n | \mathbf{y}_n, \theta)$ . [40%]

(c) The retailer would like to enhance the model by including two additional pieces of information. First, they know the shoe size of their customers, which takes integer values in the range  $\{2 \dots 12\}$ . Second, they know the age of their customers which affects their relative-size. Describe how the model can be altered to incorporate this information. [30%]

3 The *binary latent feature model* first draws two binary latent variables  $s_1$  and  $s_2$  from independent Bernoulli distributions. That is,  $p(s_1 = 1|\theta) = \pi_1$  and  $p(s_2 = 1|\theta) = \pi_2$ . Observations  $\mathbf{y}$ , which are real valued and  $D$  dimensional, are produced by multiplying the latent variables by the associated weights ( $\mathbf{w}_1$  and  $\mathbf{w}_2$ ), adding these contributions, and corrupting with isotropic Gaussian noise of variance  $\sigma_y^2$ . That is,  $p(\mathbf{y}|s_1, s_2, \theta) = \mathcal{N}(\mathbf{y}; \mathbf{w}_1 s_1 + \mathbf{w}_2 s_2, \mathbf{I}\sigma_y^2)$ .

Above, the model parameters have been denoted  $\theta = \{\pi_1, \pi_2, \mathbf{w}_1, \mathbf{w}_2, \sigma_y^2\}$ .

(a) Mathematically define a *mixture of Gaussians model*. Your definition should identify the *mixing proportions*, *component means* and *component variances*. [20%]

(b) Express the binary latent feature model, described at the start of this question, as a mixture of Gaussians stating the *mixing proportions*, *component means* and *component variances* in terms of the parameters  $\theta$ . [40%]

(c) A machine learner will employ the EM algorithm to learn the parameters of the *binary latent feature model* defined above.

(i) Compute the E-step update by deriving the posterior distribution over the latent variables given an observed variable,  $p(s_1, s_2|\mathbf{y}, \theta)$ . [15%]

(ii) Describe how you would mathematically derive the M-step update. Your description should outline the steps involved, but does not require detailed calculation. [25%]

4 (a) Define a *hidden Markov model* that has a discrete hidden state. Your definition should identify the *initial state probabilities*, *transition probabilities* and the *emission distribution*. [20%]

(b) A sequence model comprises a discrete latent variable  $s_t$  and a continuous observed variable  $y_t$ . The discrete state always begins with value 1, that is  $p(s_1 = 1) = 1$ . The state then evolves according to following rule

$$p(s_t = k | s_{t-1} = k') = \begin{cases} 0.9 & \text{for } k = k' + 1 \\ 0.1 & \text{for } k = 1 \\ 0 & \text{for all other } k \end{cases}.$$

The observations are generated from Gaussian distributions  $p(y_t | s_t) = \mathcal{N}(y_t; s_t^2, 0.1^2)$ .

(i) Sketch a typical sample from the model for  $T = 20$  time steps. Your sketch should show both the latent process and the corresponding observed process. Explain the salient features. [30%]

(ii) An algorithm has been used to compute the posterior distribution over the latent variable  $s_t$  from  $t$  observations  $y_{1:t}$ . The posterior is found to be

$$p(s_t = k | y_{1:t}) = \begin{cases} \frac{2}{3} & \text{for } k = 3 \\ \frac{1}{3} & \text{for } k = 4 \\ 0 & \text{for all other } k \end{cases}.$$

Compute the predictive distribution over the next observed variable i.e.  $p(y_{t+1} | y_{1:t})$ . Explain each step in your calculation. [40%]

(iii) What algorithm could be used to compute the posterior distribution  $p(s_t = k | y_{1:t})$ ? Would the standard implementation of this algorithm need to be modified to handle long sequences? Explain your reasoning. [10%]

**END OF PAPER**

**THIS PAGE IS BLANK**