

# COMPUTER SCIENCE TRIPOS Part IB – mock – Paper 6

## 1 Foundations of Data Science (DJW)

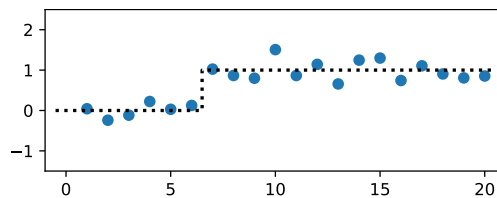
- (a) A 0/1 signal is being transmitted over a noisy channel. The transmitted signal at timeslot  $i \in \{1, \dots, n\}$  is  $x_i \in \{0, 1\}$ , and furthermore we know that this signal starts at 0 and then flips to 1, i.e. there is a parameter  $\theta \in \{1, \dots, n-1\}$  such that

$$x_i = \begin{cases} 0 & \text{for } i \leq \theta, \\ 1 & \text{for } i > \theta, \end{cases}$$

but the value of  $\theta$  is unknown. The channel is noisy, and the received signal in timeslot  $i$  is

$$Y_i \sim x_i + \text{Normal}(0, \varepsilon^2)$$

where  $\varepsilon$  is known.



- (i) Given received signals  $(y_1, \dots, y_n)$ , find an expression for the log likelihood,  $\log \text{lik}(\theta \mid y_1, \dots, y_n)$ . [5 marks]
- (ii) Give pseudocode for finding the maximum likelihood estimator  $\hat{\theta}$ . [3 marks]
- (b) The Gaussian Mixture Model with  $m$  components can be written as a two-stage random variable: first generate  $K \in \{1, \dots, m\}$ ,  $\mathbb{P}(K = k) = p_k$ , then generate  $X \sim \text{Normal}(\mu_K, \sigma_K^2)$ . Here  $p_1, \dots, p_m$  and  $\mu_1, \dots, \mu_m$  and  $\sigma_1, \dots, \sigma_m$  are unknown parameters, with  $p_k > 0$  and  $\sigma_k > 0$  for all  $k$ , and  $p_1 + \dots + p_m = 1$ . The number of components  $m$  is known.
- (i) Give formulae for  $\mathbb{P}(X \leq x \mid K = k)$  and for  $\mathbb{P}(X \leq x)$ . [3 marks]
- You should leave your answers in terms of the cumulative distribution function for a Normal distribution,  $\Phi_{\mu, \sigma}(x) = \mathbb{P}(\text{Normal}(\mu, \sigma^2) \leq x)$ .*
- (ii) Calculate the density of  $X$ . [4 marks]
- (iii) Given a dataset  $(x_1, \dots, x_n)$ , explain how to fit the unknown parameters using numerical optimization. [5 marks]