

Q1

$$a) \quad i) \quad x^{(MAP)} = \underset{x}{\operatorname{arg\,max}} p(x|y) = \underset{x}{\operatorname{arg\,max}} \log p(x,y)$$

$$\log p(x,y) = -\frac{1}{2} \left( y - \frac{1}{2} k x^2 \right)^2 - \frac{1}{2\sigma^2} x^2 + c \quad \leftarrow \text{don't depend on } x$$

$$\text{let } u = x^2$$

$$\frac{d}{du} \log p(x,y) = \frac{d}{du} \left( -\frac{1}{2} \left( y - \frac{1}{2} k u \right)^2 - \frac{1}{2\sigma^2} u \right) = 0 \quad \text{NB of } \frac{d}{du} \text{ see note below } y=0 \Rightarrow x_{MAP}=0$$

$$0 = \cancel{2} \cdot \left( y - \frac{1}{2} k u \right) \cdot \left( -\frac{1}{2} k \right) + 1/\sigma^2$$

$$0 = -k y + \frac{1}{2} k^2 u + \frac{1}{\sigma^2}$$

$$u = -\frac{2}{k^2 \sigma^2} + \frac{2y}{k}$$

$$\Rightarrow x_{MAP} = \begin{cases} + \sqrt{\frac{2y}{k} - \frac{2}{k^2 \sigma^2}} & \text{if } y > \frac{1}{k \sigma^2} \\ - & \text{else } x_{MAP} = 0 \end{cases}$$

I was happy if people did not spot that the result only applied if this condition held though

$$ii) \quad \sigma^2 \rightarrow \infty \quad \Rightarrow \quad x_{MAP} \rightarrow \begin{cases} + \sqrt{\frac{2y}{k}} \\ - & \end{cases} \quad \text{c.f. } y = \frac{1}{2} k x^2 + \text{noise}$$

inverts this

This is the maximum likelihood estimate of  $x$  as the prior is uninformative

$$b) i) \quad p(d | y_1, y_2) \propto p(d) p(y_1, y_2 | d) \propto N(d; \mu_{d|y_1, y_2}, \sigma_{d|y_1, y_2}^2)$$

$$\begin{array}{ccc} \uparrow & & \uparrow \\ N(d; 0, 1) & & N\left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}; \begin{bmatrix} d \\ d \end{bmatrix}, \begin{bmatrix} 1 & \alpha \\ \alpha & 1 \end{bmatrix}^{-1}\right) \end{array}$$

$$ii) \quad \log p(d | y_1, y_2) = c - \frac{1}{2} d^2 - \frac{1}{2} \begin{bmatrix} y_1 - d & y_2 - d \end{bmatrix} \begin{pmatrix} 1 & \alpha \\ \alpha & 1 \end{pmatrix}^{-1} \begin{bmatrix} y_1 - d \\ y_2 - d \end{bmatrix}$$

$$= c - \frac{1}{2} d^2 - \frac{1}{2} (y_1 - d)^2 - \frac{1}{2} (y_2 - d)^2 - (y_1 - d)(y_2 - d)\alpha$$

$$= c - \frac{1}{2} d^2 - \frac{1}{2} d^2 + y_1 d - \frac{1}{2} d^2 + y_2 d - \alpha d^2 + (y_1 + y_2)d\alpha$$

$$= c - \frac{1}{2} (3 + 2\alpha)d^2 + (y_1 + y_2)(1 + \alpha)d$$

$$\Rightarrow \sigma_{d|y_1, y_2}^2 = \frac{1}{3 + 2\alpha} \quad \mu_{d|y_1, y_2} = \frac{1}{3 + 2\alpha} \sqrt{(1 + \alpha)(y_1 + y_2)}$$

$$iii) \quad \sigma_{d|y_1, y_2}^2 = \frac{1}{3} \quad \text{when } \alpha = 0$$

$$\text{when } \alpha > 0 \quad \sigma_{d|y_1, y_2}^2 < \frac{1}{3} \quad \text{ie more certain}$$

note

$$\begin{pmatrix} 1 & \alpha \\ \alpha & 1 \end{pmatrix}^{-1} = \frac{1}{1 - \alpha^2} \begin{pmatrix} 1 & -\alpha \\ -\alpha & 1 \end{pmatrix}$$

negative correlation

$$\Rightarrow \alpha = 0 \Rightarrow \text{independent } y_1, y_2$$

$$\alpha > 0 \Rightarrow \text{negative correlation} \Rightarrow \text{forming } y_1 + y_2 \text{ removes some of the noise}$$

$$\Rightarrow \text{can be more certain}$$

$$(\alpha < 0 \Rightarrow \text{positive correlation in } y_1 \text{ \& } y_2 \Rightarrow \sigma_{d|y_1, y_2}^2 > \frac{1}{3} \Rightarrow \text{less certain})$$

Q2

$$a) \quad i) \quad L(w, \sigma^2) = \log \prod_{n=1}^N p(y_n | x_n, w, \sigma^2) \quad (\text{log likelihood})$$

$$ii) \quad \frac{d}{dw} L(w, \sigma^2) = 0 = \frac{d}{dw} \left( -\frac{1}{2\sigma^2} \sum_n (y_n - wx_n)^2 - \frac{N}{2} \log \sigma^2 \right)$$

$$\Rightarrow 0 = \sum_{n=1}^N (y_n - wx_n) x_n$$

$$\Rightarrow w = \frac{\sum_{n=1}^N y_n x_n}{\sum_{n=1}^N x_n^2}$$

$$iii) \quad \text{let } YX_N = \sum_{n=1}^N y_n x_n \quad XX_N = \sum_{n=1}^N x_n^2$$

$$YX_{N+1} = YX_N + y_{N+1} x_{N+1} \quad XX_{N+1} = XX_N + x_{N+1}^2$$

$$w_{N+1} = \frac{YX_{N+1}}{XX_{N+1}}$$

$$b) i) \sigma_{MAP}^2 = \arg \max_{\sigma^2} p(\{y_n\}_{n=1}^N | \{x_n\}_{n=1}^N, w_{MAP}, \sigma^2) p(\sigma^2)$$

$$= \arg \max_{\sigma^2} \left[ \sum_{n=1}^N \log p(y_n | x_n, w_{MAP}, \sigma^2) + \log p(\sigma^2) \right]$$

$$= \arg \max_{\sigma^2} \left( -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - w_{MAP} x_n)^2 - \frac{N}{2} \log \sigma^2 - \alpha \log \sigma^2 - \beta / \sigma^2 \right)$$

$$= \arg \max_{\sigma^2} \left( \underbrace{-\frac{1}{\sigma^2} \left( \beta + \frac{1}{2} \sum_{n=1}^N (y_n - w_{MAP} x_n)^2 \right)}_{\beta'} - \underbrace{\left( \frac{N}{2} + \alpha \right)}_{\alpha'} \log \sigma^2 \right)$$

$$\frac{d}{d\sigma^2} \left( -\frac{1}{\sigma^2} \beta' - \alpha' \log \sigma^2 \right) = 0$$

$$+\frac{1}{\sigma^4} \beta' - \frac{\alpha'}{\sigma^2} = 0$$

$$\Rightarrow \sigma_{MAP}^2 = \frac{\beta'}{\alpha'} = \frac{\beta + \frac{1}{2} \sum_{n=1}^N (y_n - w_{MAP} x_n)^2}{\frac{N}{2} + \alpha}$$

$$ii) \alpha = N'/2 \text{ where } N' = \# \text{ of pseudopoints}$$

$$\beta = \Gamma'/2 \text{ where } \Gamma' = \text{sum of squared errors (noise error) on pseudopoints}$$

Q3

a) E-Step

$$\{q^{(new)}(s_n)\}_{n=1}^N = \arg \max_{\{q(s_n)\}_{n=1}^N} q(\theta_t, \{q(s_n)\}_{n=1}^N)$$

current setting of parameters

$$\Rightarrow q^{(new)}(s_n) = p(s_n | \underline{x}_n, \theta_t) \propto p(s_n | \theta_t) p(\underline{x}_n | s_n, \theta_t)$$

In this case

$$p(s_n = k) = 1/k \quad p(\underline{x}_n | s_n = k) = \prod_{d=1}^D \pi_{kd}^{x_{dn}} (1 - \pi_{kd})^{1-x_{dn}}$$

$$\text{let } u_{nk} = 1/k \cdot \prod_{d=1}^D \pi_{kd}^{x_{dn}} (1 - \pi_{kd})^{1-x_{dn}}$$

$$q^{(new)}(s_n = k) = \frac{u_{nk}}{\sum_k u_{nk}}$$

b) Hard E-Step will perform  $\arg \max_k p(s_n = k | \underline{x}_n) = \arg \max_k \log p(s_n = k | \underline{x}_n)$

$$= \arg \max_k \sum_{d=1}^D (x_{dn} \log \pi_{kd} + (1-x_{dn}) \log (1-\pi_{kd}))$$

$$= \arg \max_k KL(\underline{x}_n \parallel \underline{\pi}_k)$$

i.e. you take each data point  $\underline{x}_n$  & compute how close it is to each  $\underline{\pi}_k$  according to the KL divergence. This is equivalent to the  $s_n = \arg \min_k \|\underline{x}_n - \underline{m}_k\|^2$  assignment step in k-means

$$c) \text{ M-Step } \theta^{(new)} = \arg \max_{\theta} f(\theta, \{q^{(old)}(s_n)\}_{n=1}^n)$$

$$= \arg \max_{\theta} \sum_n \mathbb{E}_{q^{(old)}(s_n)} [\log p(s_n, \underline{x}_n | \theta)]$$

$$\frac{\partial f}{\partial \pi_{l,e}} = \sum_n \mathbb{E}_{q^{(old)}(s_n)} [\log p(\underline{x}_n | s_n, \theta)]$$

$$= \sum_n \sum_{k,d} q^{(old)}(s_n=k) (x_{n,d} \log \pi_{k,d} + (1-x_{n,d}) \log (1-\pi_{k,d}))$$

$$= \sum_n q^{(old)}(s_n=l) \left( \frac{x_{n,e}}{\pi_{l,e}} - \frac{(1-x_{n,e})}{1-\pi_{l,e}} \right) = 0$$

$$= \sum_n q^{(old)}(s_n=l) (x_{n,e} - \pi_{l,e})$$

$$\Rightarrow \pi_{l,e} = \frac{\sum_n q^{(old)}(s_n=l) x_{n,e}}{\sum_n q^{(old)}(s_n=l)}$$

d) For hard assignments  $\pi_{k,d}$  would be updated to be the average of  $x_{n,d}$  for all the points assigned to cluster  $k$ .

$$\pi_{k,d} = \text{mean}(x_{n,d} : s_n=k)$$

Q4

a)  $x_1 \sim N(\mu_1, \sigma_1^2)$        $x_t = \lambda x_{t-1} + \sigma \varepsilon_t$        $\varepsilon_t \sim N(0, 1)$

b)  $z_t = z_{t-1} + x_t$

$$= z_{t-1} + \lambda x_{t-1} + \sigma \varepsilon_t$$

$$= z_{t-1} + \lambda (z_{t-1} - z_{t-2}) + \sigma \varepsilon_t$$

$$= (1+\lambda) z_{t-1} - \lambda z_{t-2} + \sigma \varepsilon_t = \lambda_1 z_{t-1} + \lambda_2 z_{t-2} + \sigma \varepsilon_t$$

$\Rightarrow$  AR(2) process with  $\lambda_1 = 1+\lambda$  &  $\lambda_2 = -\lambda$

c)  $\underline{s}_1 \sim N(\underline{\mu}_1, \underline{\Sigma}_1)$        $\downarrow$   $t=2 \dots T$   $\underline{s}_t | \underline{s}_{t-1} \sim N(\underline{A} \underline{s}_{t-1}, \underline{Q})$        $\overbrace{t=1 \dots T}$   $\underline{y}_t | \underline{s}_t \sim N(\underline{C} \underline{s}_t, \underline{R})$

d) one option:  $\underline{s}_t = \begin{bmatrix} x_t \\ z_t \end{bmatrix}$        $\underline{C} = \underline{I}^{2 \times 2}$        $\underline{R} = \underline{I}^{2 \times 2}$        $\underline{A} = \begin{bmatrix} \lambda & 0 \\ \lambda & 1 \end{bmatrix}$        $\underline{Q} = \sigma^2 \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$    
perfectly correlated state noise

as  $\begin{bmatrix} x_t \\ z_t \end{bmatrix} = \underline{A} \begin{bmatrix} x_{t-1} \\ z_{t-1} \end{bmatrix} + \underline{Q}^{1/2} \begin{pmatrix} \eta_t^{(1)} \\ \eta_t^{(2)} \end{pmatrix} = \begin{bmatrix} \lambda & 0 \\ \lambda & 1 \end{bmatrix} \begin{bmatrix} x_{t-1} \\ z_{t-1} \end{bmatrix} + \frac{\sigma}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \eta_t^{(1)} \\ \eta_t^{(2)} \end{pmatrix} \stackrel{\text{diag}}{=} \begin{bmatrix} \lambda x_{t-1} + \sigma \varepsilon_t \\ \lambda x_{t-1} + z_{t-1} + \sigma \varepsilon_t \end{bmatrix}$

e) Kalman Smoother will return  $p(x_t, z_t | \underline{y}_{1:T}) = N \left( \begin{bmatrix} x_t \\ z_t \end{bmatrix}; \begin{bmatrix} \bar{x}_t \\ \bar{z}_t \end{bmatrix}, \begin{bmatrix} \sigma_{x_t}^2 & \rho_t \\ \rho_t & \sigma_{z_t}^2 \end{bmatrix} \right)$

$\bar{z}_{1:T} = \arg \max_{z_{1:T}} p(z_{1:T} | \underline{y}_{1:T})$       since ① deterministic relation between  $x$  &  $z$

② mode = mean

## Summary of Exam Marks

The examination was taken by 120 candidates in total. The raw marks (from those who had taken IB) had an average of 65.5% and standard deviation 13.0% with the top candidate scoring 92% and bottom candidate scoring 30%

### Q1 Fundamental Inference Concepts

*112 attempts, Ave. raw mark 13.7/20, St.Dev. 2.7, Maximum 20, Minimum 10.*

A popular question. Generally well answered. Many people failed to solve correctly for the MAP estimate in part a (i). Very few candidates realised that the depth sensors in b (iii) are negatively correlated when alpha is positive.

### Q2 Classification and KL divergence

*115 attempts, Ave. raw mark 13.6/20, St.Dev. 4.0, Maximum 20, Minimum 6.*

Generally well answered. A surprising number of candidates could not derive the standard linear regression expression for the weights in part a (ii) which is simple book work, but all other parts were well handled.

### Q3 The EM Algorithm

*84 attempts, Ave. raw mark 11.4/20, St.Dev. 3.0, Maximum 19, Minimum 6.*

This question is on a challenging topic, but was answered reasonably well in general. In part (b) many candidates realised that the hard E-step could be interpreted as minimising a distance, but none identified that this distance is the KL divergence between the binary data vector  $x$  and the cluster prior parameter. Many candidates failed to get to the correct analytic expressions for the E- and M-steps, but generally the attempts got close and used the right method.

### Q4 Auto-regressive Models and Linear Gaussian State Space Models

*49 attempts, Ave. raw mark 11.5/20, Stan. Dev. 3.5, Maximum 19, Minimum 6.*

This was a challenging question and answers were patchy. Many candidates failed to solve part (b) and actually used spurious methods. Consequently most candidates did not realise that the new process on  $z$  is AR(2). Very few candidates constructed a linear Gaussian state space model that correctly captured the correlations between the two variables  $x$  and  $z$ .