# Q1

a) i) local posterior $\quad p(\theta | x_k) = \dfrac{p(\theta) \, p(x_k | \theta)}{p(x_k)}$

$\downarrow$ prior

$\leftarrow$ probability of data given $\theta$ or likelihood of $\theta$

$p(x_k) \leftarrow$ normalising constant

$\propto p(\theta) \, p(x_k | \theta)$

ii) global posterior $\quad p(\theta | x_1, x_2 \cdots x_k) \propto p(\theta) \displaystyle\prod_{k=1}^{K} p(x_k | \theta)$

$\propto p(\theta) \displaystyle\prod_{k=1}^{K} \dfrac{p(\theta | x_k)}{p(\theta)}$

$= \left[ \dfrac{1}{p(\theta)} \right]^{k-1} \displaystyle\prod_{k=1}^{k} p(\theta | x_k)$

iii) Unlikely to happen if we have the correct model & sufficient data.

However it is likely that each hospital is slightly different & so a model like

this would be more appropriate

global parameters
$\downarrow$

local parameters (will each be different in general)

$p(\theta) \displaystyle\prod_{k=1}^{K} p(\theta_k | \theta) \, p(x_k | \theta_k)$

If the data actually come from this model & yet we fit the global model

instead, fitting the local models will perform better when the local parameters

$\theta_k$ are fairly different from one another.

b)

$$p\left(C=1 \mid LFT=1, PCR=0\right) = \frac{p(C=1)\, p\left(LFT=1 \mid C=1\right)\, p\left(PCR=0 \mid C=1\right)}{p(C=1)\, p\left(LFT=1 \mid C=1\right)\, p\left(PCR=0 \mid C=1\right) + p(C=0)\, p\left(PCR=0 \mid C=0\right)\, p\left(LFT=1 \mid C=0\right)}$$

$$= \frac{0.01 \times 0.85 \times 0.05}{0.01 \times 0.85 \times 0.05 + 0.99 \times 0.001 \times 1}$$

$$= \frac{1}{1 + \dfrac{0.99}{0.85} \times \dfrac{0.1}{0.05}} \approx 0.3$$

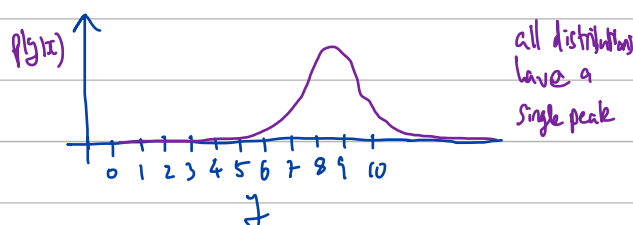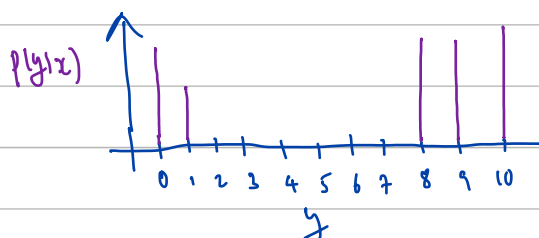i.e. still    30% chance of having covid  even though the PCR

came back negative.

# Q2

## a) i)

$$p(y=k \mid \underline{x}) \propto \exp(\underline{w_k}^T \underline{x})$$   $$p(y \mid \underline{x}) = N(y; \underline{w}^T \underline{x}, \sigma_y^2)$$

- models data as discrete — good as data are discrete
- does not account for ordering of data i.e. that for strong students all high scores are likely, all low scores less likely & vice versa for weak students (linked to the fact that it has weights for each score $\{\underline{w_k}\}_{k=0}^{10}$)

- models data as continuous which is inappropriate
- captures the fact that a strong student is likely to get a high score, but not a low score. (Linked to having a single weight vector $\underline{w}$)



$p(y|x)$ ... can output distributions like this
0 1 2 3 4 5 6 7 8 9 10
$y$

$p(y|x)$ ... all distributions have a single peak
0 1 2 3 4 5 6 7 8 9 10
$y$
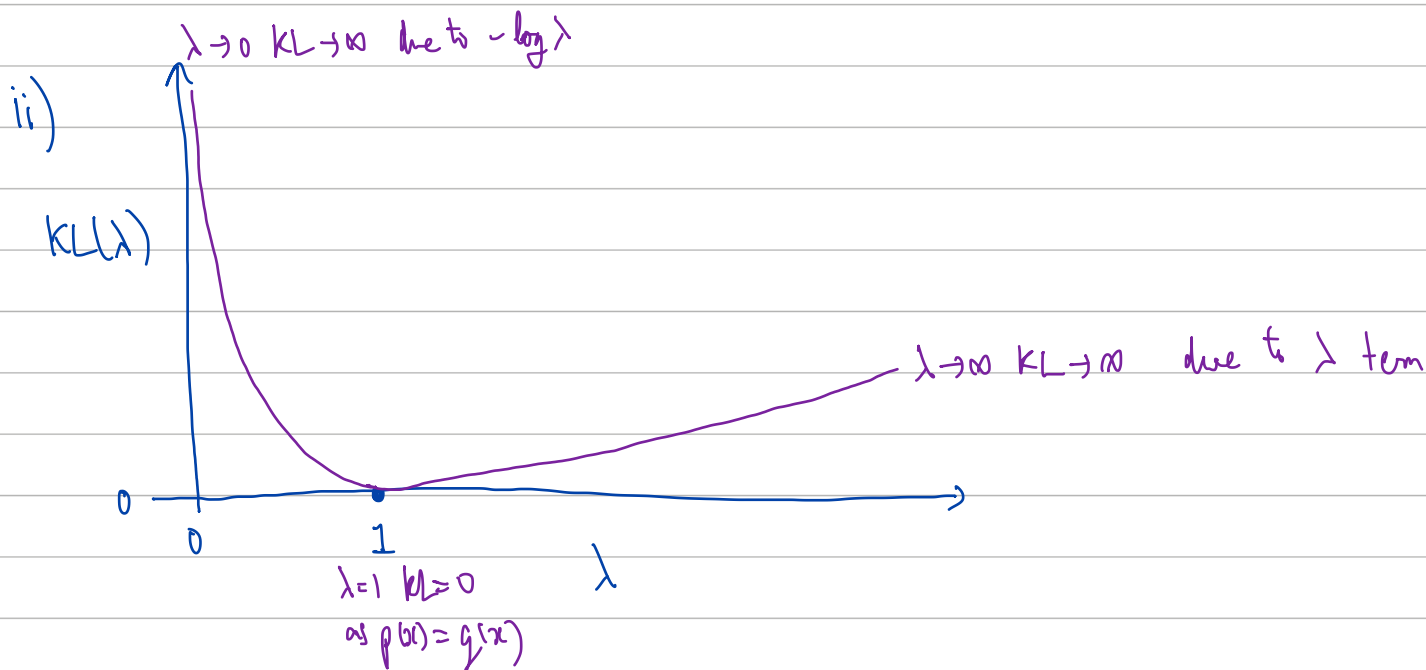
## b) Lots of possible options

one idea    let   $p(z \mid \underline{x}) = N(z; \underline{w}^T \underline{x}, 1)$

$$\begin{cases} 0 & \text{if} & z \leq 0 \\ 1 & \text{if} & 0 < z \leq \mu_1 \\ 2 & \text{if} & \mu_1 < z \leq \mu_2 \\ \quad \vdots \\ K & \text{if} & \mu_{K-1} < z \end{cases}$$

This is called the "ordered probit model" & the general problem is called probit regression.

b) i) $\quad KL\big(q(x) \,\|\, p(x)\big) = \mathbb{E}_{q(x)}\big[\log p(x) - \log q(x)\big]$

$$= \mathbb{E}_{q(x)}\big[-\log \lambda - x/\lambda + x\big]$$

$$\underset{\mathbb{E}_{q}(x) = \lambda}{= \quad -\log \lambda - 1 + \lambda}$$

ii)



$\lambda \to 0 \; KL \to \infty$ due to $-\log \lambda$

KL($\lambda$)

$\lambda \to \infty \; KL \to \infty$ due to $\lambda$ term

0

1

$\lambda = 1 \; KL = 0$

as $p(x) = q(x)$

$\lambda$

NB. do not need to prove properties of KL — can just use them directly

# Q3

a) $p(x_c | N_c, \mu) = N(x_c; N_c \mu, N_c)$  [means & variances add for iid variables]

<span style="color:red">Uniform(1,10)</span> ↓

b) $p(N_c | x_c, \mu) \propto p(N_c) p(x_c | N_c, \mu) \implies p(N_c | x_c, \mu) = \dfrac{N(x_c; \mu N_c, N_c)}{\sum\limits_{N_c = 1}^{10} N(x_c; \mu N_c', N_c')}$

c) i) EM will find maximum likelihood setting of parameters

$$\mu = \arg\max_{\mu} \; \log p(x_{1:c} | \mu)$$

E-Step $\quad \left\{ q^{(new)}(N_c) \right\}_{c=1}^{C} = \arg\max_{\left\{ q(N_c) \right\}_{c=1}^{C}} \mathcal{F}\left( \mu, \left\{ q(N_c) \right\}_{c=1}^{C} \right)$

$$= p(N_c | x_c, \mu)$$

M-Step $\quad \mu^{(new)} = \arg\max_{\mu} \mathcal{F}\left( \mu, \left\{ q(N_c) \right\}_{c=1}^{C} \right)$

ii) $\quad \mu^{(new)} = \arg\max_{\mu} \sum\limits_{c=1}^{C} \sum\limits_{N_c=1}^{10} q(N_c) \log p(x_c | \mu, N_c)$

$$= \arg\max_{\mu} \sum\limits_{c=1}^{C} \mathbb{E}_{q(N_c)} \left[ -\frac{1}{2N_c}(x_c - N_c \mu)^2 - \log N_c - \frac{1}{2}\log 2\pi \right]$$

<span style="color:purple">total weight of all bags</span>

Take derivatives & set to zero:

$$\sum\limits_{c=1}^{C} \mathbb{E}_{q(N_c)} \left[ \frac{1}{2N_c}(x_c - \mu^* N_c) N_c \right] = 0 \implies \mu^* = \frac{\sum\limits_{c=1}^{C} x_c}{\sum\limits_{c=1}^{C} \mathbb{E}_{q(N_c)}(N_c)}$$

<span style="color:red">sum of expected number of oranges in all bags</span>

# Q4

a) $1\ 2222\ 1112\ 111\ |\ 3333\ |\ 222$ ← $y_{1:20}$  (from lectures)

        ↑ 50% 1s & 2s     ↑ all 3s   ↑ 50:50

$1\ 1\ 1\ 1\ 1\ 1\ 1\ 1\ 1\ 111\ |\ 22222\ |\ 111$ ← $x_{1:20}$

b) $p(y_t \mid y_{t-1}) = \sum\limits_{x_t, x_{t-1}} p(y_t \mid x_t)\, p(x_t \mid x_{t-1})\, p(x_{t-1} \mid y_{t-1})$

               if $y_{t-1} = 1$ or $2$ then $x_{t-1} = 1$
               if $y_{t-1} = 3$ then $x_{t-1} = 2$

$$= \begin{bmatrix} 0.5 & 0 \\ 0.5 & 0 \\ 0 & 1 \end{bmatrix} \times \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \times \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.5 & 0 \\ 0.5 & 0 \\ 0 & 1 \end{bmatrix} \times \begin{bmatrix} 0.9 & 0.9 & 0.1 \\ 0.1 & 0.1 & 0.9 \end{bmatrix} = \begin{bmatrix} 0.45 & 0.45 & 0.05 \\ 0.45 & 0.45 & 0.05 \\ 0.1 & 0.1 & 0.9 \end{bmatrix}$$

c) in general

$$p(y_t \mid y_{1:t-1}) = \sum\limits_{x_t, x_{t-1}} p(y_t \mid x_t)\, p(x_t \mid x_{t-1})\, p(x_{t-1} \mid y_{1:t-1})$$

note difference here

However although these two appear different note in this specific case:

if you know $y_{t-1}$ you know $x_{t-1}$

ie. $p(x_{t-1} \mid y_{1:t-1}) = p(x_{t-1} \mid y_{t-1})$

$\therefore p(y_t \mid y_{1:t-1}) = p(y_t \mid y_{t-1}) = \begin{bmatrix} 0.45 & 0.45 & 0.05 \\ 0.45 & 0.45 & 0.05 \\ 0.1 & 0.1 & 0.9 \end{bmatrix}$ & $p(y_1) = \begin{bmatrix} 1/4 \\ 1/4 \\ 1/2 \end{bmatrix}$

d) Not true in general that

$$p(x_{t-1} \mid y_{1:t-1}) = p(x_{t-1} \mid y_{t-1})$$

# Summary of Exam Marks

The examination was taken by 132 candidates in total. The marks had an average of 65.7% and standard deviation 11.5% with the top at 97% and bottom at 35%.

### Q1    Fundamental Inference Concepts

118 attempts, Ave. raw mark 13.3/20, Stan. Dev. 3.0, Maximum 20, Minimum 7.

A popular question. Generally well answered. In part (aii) many solutions did not give the global posterior as a function of the local posteriors, but rather as a product of the prior and the local likelihood functions which is incorrect. Most people failed to substitute numbers into Bayes' rule correctly in part (b) resulting in incorrect numerical answers.

### Q2    Classification and KL divergence

89 attempts, Ave. raw mark 12.5/20, Stan. Dev. 2.7, Maximum 20, Minimum 7.

Generally well answered. In part (a) many failed to spot that the softmax model does not capture the fact that if a student is strong, they are likely to have a high probability for any high score and a low probability of all low scores, and vice versa for weak students. Part (b) was answered more successfully. A significant minority of answers to (bii) sketched KL divergences that went negative, which is impossible by definition.

### Q3    The EM Algorithm

84 attempts, Ave. raw mark 12.9/20, Stan. Dev. 3.4, Maximum 19, Minimum 3.

This question is on a challenging topic, but was well answered in general. Curiously many students did not realise that variances of independent variables add and therefore answered the first part of the question incorrectly. In the last part, many people differentiated the free-energy correctly, but then failed to rearrange the expression for \mu correctly.

### Q4    Discrete Hidden Markov Models

105 attempts, Ave. raw mark 13.5/20, Stan. Dev. 2.811, Maximum 20, Minimum 9.

The first two parts of this question were well done. Part (c) involved computing the transition matrix and initial state distribution of the equivalent bigram model was not successfully completed by about 1/3 of candidates. Most candidates correctly stated for part (d) that it was not usually possible to write an HMM as a bigram model in this way, but they didn't spot the reason that it is possible here is because, for the example considered in the question, the latent state is deterministic given the observed state at the same time point