

Module MLMI1: Introduction to Machine Learning  
Solutions to Example Sheet 2: Clustering, the Expectation  
Maximisation Algorithm, Markov Models,  
and Hidden Markov Models

*Clustering*

1. K-means clustering\*

Consider the K-means algorithm that seeks to minimise the cost function

$$C = \sum_{n=1}^N \sum_{k=1}^K s_{nk} \|x_n - m_k\|^2$$

where  $m_k$  is the mean (centre) of cluster  $k$ ,  $x_n$  is the data point  $n$ ,  $s_{nk} = 1$  indicates that the  $n$ th data point has been assigned to cluster  $k$ , and  $s_{nk} = 0$  indicates that it has not been assigned to that cluster. There are  $N$  data points and  $K$  clusters.

- (a) Given all the assignments  $\{s_{nk}\}$  determine the value of  $m_k$  that minimises the cost  $C$  and give an interpretation in terms of the K-means algorithm.
- (b) You would like to automatically learn the number of clusters  $K$  from data. One possibility is to minimise the cost  $C$  as a function of  $K$ . Explain whether this is a good idea or not, and what the solution to this minimisation is.
- (c) In many real-world applications, data points arrive sequentially and one wants to cluster them as they come in. Devise a sequential variant of the k-means algorithm which takes in one data point at a time and updates the means  $\{m_1, \dots, m_K\}$  sequentially without revisiting previous data points. Describe your sequential algorithm.

$$5 \quad C(\{S\}, \{M\}) = \sum_{n=1}^N \sum_{k=1}^K S_{nk} \|x_n - m_k\|^2 = \sum_n \sum_{k=1}^K S_{nk} (x_{nd} - m_{kd})^2$$

$$a) \frac{dC(\{S\}, \{M\})}{dM_{le}} = -2 \sum_{n=1}^N S_{nl} (x_{ne} - m_{le}) = 0$$

$$\Rightarrow \sum_{n=1}^N S_{nl} x_{ne} = M_{le} \sum_{n=1}^N S_{nl}$$

$$\Rightarrow M_{le} = \frac{\sum_{n=1}^N S_{nl} x_{ne}}{\sum_{n=1}^N S_{nl}}$$

sum of all of the data values assigned to  $k^{th}$  cluster

Or in vector form  $\underline{M}_k = \frac{\sum_{n=1}^N S_{nk} \underline{x}_n}{\sum_{n=1}^N S_{nk}}$

number of datapoints assigned to  $k^{th}$  cluster

b) Bad idea : if  $K=N$  &  $M_n = \underline{x}_n$  cost is zero  
 i.e. each datapoint is a "cluster"

c) Several possible options for this are but a logical soln would be:

$$\text{Since } \underline{M}_k = \frac{\sum_{n=1}^N \underline{x}_n s_{nk}}{\sum_{n=1}^N s_{nk}} = \frac{\sum_{n=1}^{N-1} \underline{x}_n s_{nk} + \underline{x}_N s_{nk}}{\sum_{n=1}^{N-1} s_{nk} + s_{Nk}}$$

$\beta_k^{(N-1)}$

$\alpha_k^{(N-1)}$

New datum arrives

- assign to nearest cluster and set assignment vector  $\{\hat{s}_{nk}\}_{k=1}^K$

$$-\text{update statistics} \quad \beta_k^{(N)} = \beta_k^{(N-1)} + \underline{x}_N s_{nk} \quad [\text{running sum of data vectors associated w/ kth cluster}]$$

$$\alpha_k^{(N)} = \alpha_k^{(N-1)} + s_{nk} \quad [\text{running count of # of data points associated w/ kth cluster}]$$

$$-\text{recompute mean} \quad \underline{M}_k = \frac{\beta_k^{(N)}}{\alpha_k^{(N)}}$$

- repeat

## 2. The KL Divergence

The KL Divergence between two discrete distributions  $p(x = k) = p_k$  and  $p(x = k) = q_k$  is defined as  $\mathcal{KL}(q, p) = \sum_{k=1}^K q_k \log \frac{q_k}{p_k}$ .

- (a) Prove that the KL Divergence is non-negative and that it attains its unique minimum when  $q_k = p_k$ .
- (b) A machine learner has a target distribution  $\mathbf{p} = [p_1, p_2, p_3, p_4, p_5, p_6] = [1, 1, 0, 0, 1, 1]/4$  which they want to fit with approximating distribution  $q$ . The choices for  $q$  available to the machine learner are:  $q_1 = [1, 1, 0, 0, 0, 0]/2$ ,  $q_2 = [0, 1, 1, 0, 0, 0]/2$ ,  $q_3 = [0, 0, 1, 1, 0, 0]/2$ ,  $q_4 = [0, 0, 0, 1, 1, 0]/2$ ,  $q_5 = [0, 0, 0, 0, 1, 1]/2$ , and  $q_6 = [1, 1, 1, 1, 1, 1]/6$ .
  - i. Determine the distribution(s)  $q_i$  that minimises  $\mathcal{KL}(q_i, p)$ . Comment on your result.
  - ii. Determine the distribution(s)  $q_i$  that minimises  $\mathcal{KL}(p, q_i)$ . Comment on your result.

Note that, by convention,  $0 \times \log 0 = 0$  since  $\lim_{\Delta \rightarrow 0} \Delta \log \Delta = 0$ .

6

$$a) \quad \mathcal{KL}(q || p) = \sum_k q_k \log \frac{q_k}{p_k}$$

constraint that probabilities  
 sum to 1

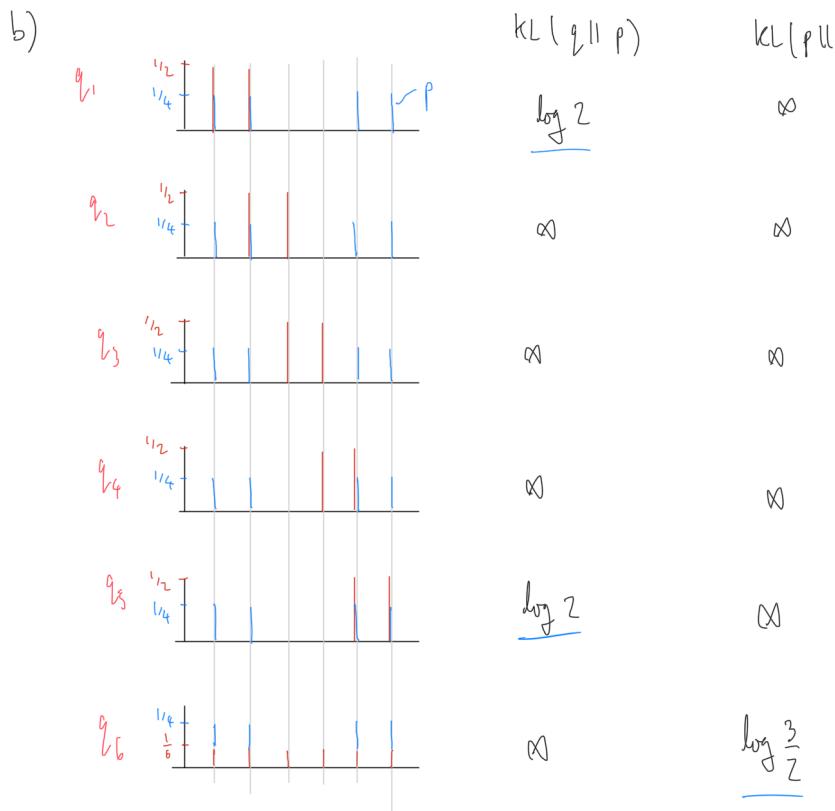
$$\mathcal{O} = \frac{\partial}{\partial q_l} \left[ \sum_k q_k \log \frac{q_k}{p_k} + \lambda \left( \sum_k q_k - 1 \right) \right]$$

↑  
Lagrange multiplier

$$\mathcal{O} = \log \frac{q_l}{p_l} + \frac{q_l}{q_l} + \lambda$$

$$\Rightarrow q_l = p_l \quad @ \text{which} \quad \mathcal{KL}(q^{\text{opt}} || p) = \sum_k p_k \log \frac{p_k}{p_k} = 0$$

$$\frac{\partial^2}{\partial q_l^2} \mathcal{KL}(q || p) = \frac{1}{q_l} > 0 \quad \Rightarrow \text{maximum}$$



$KL(q||p)$  is zero avoiding & tends to fit single modes of  $p$  with  $q$

$KL(p||q)$  is not zero avoiding & tends to cover multiple modes of  $p$  with  $q$

### *The EM Algorithm*

#### 3. Factor Analysis and EM\*

The noisy depth sensor from question 2 in example sheet 1 is used to collect measurements of the distances to a set of  $N$  objects that are unknown distances  $d_n$  metres away. The object depths can be assumed, *a priori*, to be distributed according to independent standard Gaussian distributions  $p(d_n) = \mathcal{N}(d_n; 0, 1)$ . As before, the depth sensor returns  $y_n$  a noisy measurement of the depth, that is also assumed to be Gaussian  $p(y_n|d_n, \sigma_y^2) = \mathcal{N}(y_n; d_n, \sigma_y^2)$ .

The variance of the  $\sigma_y^2$  sensor noise is unknown and must be estimated from the measured depths  $\{y_n\}_{n=1}^N$  using maximum likelihood.

- (a) Derive the steps of the EM algorithm for performing this. Your answer should include the E-Step (you may find your solution to question 2 in example sheet 1 useful here) and the M-Step.
- (b) An alternative approach to maximum-likelihood learning would optimise the log-likelihood  $p(\{y_n\}_{n=1}^N | \sigma_y^2)$  directly. Compute the objective function for this procedure. How do you think this approach will perform compared to EM?

$$p(d_n) = N(d_n; 0, 1)$$

$$p(y_n | d_n, \sigma_y^2) = N(y_n; d_n, \sigma_y^2)$$

in E-Step set  $q_k(d_n) = p(d_n | y_n, \sigma_y^2) = N(d_n; \mu_{d_n|y_n}, \sigma_{d_n|y_n}^2)$

where  $\mu_{d_n|y_n} = \frac{y_n}{1 + \sigma_y^2}$        $\sigma_{d_n|y_n}^2 = \frac{\sigma_y^2}{1 + \sigma_y^2}$       (see Q2 Examples Sheet 1)

in M-Step

$$\sigma_y^2 = \underset{\sigma_y^2}{\text{arg max}} \quad \mathbb{E}_q \left[ \log p(\{d_n\}_{n=1}^N, \{y_n\}_{n=1}^N | \sigma_y^2) \right]$$

$\uparrow$  "average log-joint"

$$G(\sigma_y^2) = \mathbb{E}_q \left[ \sum_n \log p(d_n) + \sum_n \log p(y_n | d_n, \sigma_y^2) \right]$$

$$= C - \frac{1}{2\sigma_y^2} \sum_n \mathbb{E}_{q(d_n)} \left[ (y_n - d_n)^2 \right] - \frac{N}{2} \log 2\pi \sigma_y^2 \quad *$$

$$= C - \frac{1}{2\sigma_y^2} \sum_n (y_n^2 + \mathbb{E}_{q(d_n)} [d_n^2] - 2y_n \mathbb{E}_{q(d_n)} [d_n]) - \frac{N}{2} \log 2\pi \sigma_y^2$$

$$\frac{dG(\sigma_y^2)}{d\sigma_y^2} = \frac{1}{2\sigma_y^4} \sum_n (y_n^2 + \mathbb{E}_{q(d_n)} [d_n^2] - 2y_n \mathbb{E}_{q(d_n)} [d_n]) - \frac{N}{2\sigma_y^2} = 0$$

$$\therefore \sigma_y^2 = \frac{1}{N} \sum_n y_n^2 + \frac{1}{N} \sum_n (\text{Model}_n^2 + \sigma_{\text{data}}^2) - \frac{2}{N} \sum_n y_n / \text{Model}_n$$

Algorithmically this is what you need, although for intuition ~~it~~ yields:

$$\sigma_y^2 = \frac{1}{N} \sum_n \mathbb{E}_{q(d_n)} [(y_n - d_n)^2]$$

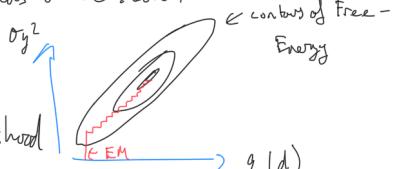
i.e. the noise variance is the average (over data points  $n$  & the posterior) of the squared error between  $y_n$  &  $d_n$ .

b)  $p(y_n | \sigma_y^2) = N(y_n; 0, 1 + \sigma_y^2)$

$$\therefore \log p(\{y_n\}_{n=1}^N | \sigma_y^2) = -\frac{N}{2} \log [2\pi(1 + \sigma_y^2)] - \frac{1}{2(1 + \sigma_y^2)} \sum_n y_n^2$$

Optimising this function for  $\sigma_y^2$  (eg by gradient ascent) would be much faster than EM. EM performs co-ordinate ascent of the free-energy in  $q$  &  $\theta$  which tends to be slow.

It turns out that if EM is analytic, so too is direct optimisation of the log likelihood which is almost always to be preferred.



4. Gaussian Mixture Models and EM\*\* (beyond tripos standard; optional)

A Gaussian Mixture Model for  $D$ -dimensional data  $\{\mathbf{x}_n\}_{n=1}^N$  comprises a categorical distribution over the class membership variables  $p(s_n = k|\theta) = \pi_k$  and a general multivariate Gaussian distribution over the observed data given the class membership variables,  $p(\mathbf{x}_n|s_n = k, \theta) = \mathcal{N}(\mathbf{x}_n; \mathbf{m}_k, \Sigma_k)$ , (i.e.  $\Sigma_k$  is not isotropic). The posterior distribution over the class membership variables has been computed  $p(s_n = k|\mathbf{x}_n, \theta) = q_{n,k}$ .

- (a) Derive the M-Step update of the EM algorithm, i.e. the formulae for updating the parameters  $\{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$  using the posterior probabilities,  $q_{n,k}$ .
- (b) A friend suggests that it might be possible to speed up the EM algorithm by updating the posterior distribution of only a subset of  $K < N$  data points during the E-Step that are selected uniformly at random each time, and then updating the parameters using the same expressions derived in part (a).
  - i. Will this procedure converge to a (local) optimum of the likelihood?
  - ii. Will this partial E-Step update procedure be computationally more efficient?

You may find the following identities useful,

$$\frac{d}{d\alpha} \log \det(\Sigma) = \text{trace} \left( \Sigma^{-1} \frac{d\Sigma}{d\alpha} \right), \quad \frac{d}{d\alpha} \Sigma^{-1} = -\Sigma^{-1} \frac{d\Sigma}{d\alpha} \Sigma^{-1}.$$

$$q. \quad \underline{\mathcal{F}}(\underline{q}(\underline{s}), \theta) = \sum_n q(s_n) \frac{\log p(x_n | s_n, \theta) p(s_n | \theta)}{q(s_n)}$$

$$q(s_n = k) = q_{nk}$$

$$\therefore \underline{\mathcal{F}}(\underline{q}(\underline{s}), \theta) = \sum_n \sum_k q_{nk} \frac{\log p(z_n | s_n = k, \theta) p(s_n = k | \theta)}{q(s_n = k)}$$

where

$$p(x_n | s_n = k, \theta) = N(x_n; \mu_k, \Sigma_k) \quad p(s_n = k | \theta) = \pi_k$$

$$\therefore \underline{\mathcal{F}}(\underline{q}(\underline{s}), \theta) = \sum_n \sum_k q_{nk} \left[ -\frac{1}{2} \log \det(\Sigma_k) - \frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) + \log \pi_k \right] + C$$

$$\begin{aligned} \frac{\partial \underline{\mathcal{F}}}{\partial \mu_{k'd'}} &= \frac{1}{2} \sum_n \sum_k q_{nk} \sum_{ij} (x_{in} - \mu_{ki}) \Sigma_{kij}^{-1} (x_{jn} - \mu_{kj}) = 0 \\ &= \frac{1}{2} \sum_n q_{nk} \left[ \sum_j \Sigma_{k'd'j}^{-1} (x_{jn} - \mu_{kj}) + \sum_i (x_{in} - \mu_{ki}) \Sigma_{k'd'}^{-1} \right] \\ &= \sum_n q_{nk} \sum_j \Sigma_{k'd'j}^{-1} (x_{jn} - \mu_{kj}) = \left[ \sum_{k'}^{-1} \sum_n q_{nk} (x_n - \mu_{k'}) \right]_{d'} \end{aligned}$$

does not depend on parameters

$$\therefore \underline{\mu}_n = \frac{\sum_n q_{nk} \underline{x}_n}{\sum_n q_{nk}}$$

alternatively, you could differentiate wrt  $\Sigma_{kij}^{-1}$  & use  $\det(\Sigma) = \frac{1}{\det(\Sigma^{-1})}$

$$\frac{d\varphi}{d\Sigma_{kij}} = \frac{1}{\det(\Sigma)} \left[ \sum_{n,k} q_{nk} \left[ -\frac{1}{2} \log \det(\Sigma) - \frac{1}{2} (\underline{x}_n - \underline{\mu}_k)^T \Sigma^{-1} (\underline{x}_n - \underline{\mu}_k) \right] \right]$$

$$\frac{d \log(\det \Sigma)}{d\alpha} = \text{trace} \left( \Sigma^{-1} \frac{d\Sigma}{d\alpha} \right) \quad \frac{d \Sigma^{-1}}{d\alpha} = - \Sigma^{-1} \frac{d\Sigma}{d\alpha} \Sigma^{-1}$$

$$\frac{d\varphi}{d\Sigma_{kij}} = \sum_{n,k} q_{nk} \left[ -\frac{1}{2} \sum_{ij} \Sigma_{kij}^{-1} \frac{d\Sigma_{kij}}{d\Sigma_{k'j'}} + \frac{1}{2} \sum_{ijlm} (\underline{x}_{ni} - \underline{\mu}_{ki}) \Sigma_{k'l}^{-1} \frac{d\Sigma_{k'm} \Sigma_{m'l}^{-1}}{d\Sigma_{k'j'}} (\underline{x}_{nj} - \underline{\mu}_{kj}) \right]$$

$$= \sum_{n,k} q_{nk} \left[ -\frac{1}{2} \sum_{ij} \Sigma_{kij}^{-1} \delta_{k'i} \delta_{i'j} \delta_{j'i} + \frac{1}{2} \sum_{ijlm} (\underline{x}_{ni} - \underline{\mu}_{ki}) \Sigma_{k'l}^{-1} \delta_{k'm} \delta_{m'l} \delta_{l'n} \delta_{nj} \Sigma_{kj}^{-1} \right]$$

$$\Rightarrow \sum_n q_{nk} \left( -\Sigma_{k'l}^{-1} + \sum_{i'j'} (\underline{x}_{ni} - \underline{\mu}_{ki}) (\underline{x}_{nj} - \underline{\mu}_{kj})^T \Sigma_{kj}^{-1} \right) = 0$$

$$= \frac{1}{2} \sum_n q_{nk} \left( -\sum_{k'j'i'} \delta_{k'i} \delta_{i'j} \delta_{j'i} + \sum_{ij} (\underline{x}_{ni} - \underline{\mu}_{ki}) \Sigma_{k'i}^{-1} \Sigma_{k'j'i'}^{-1} (\underline{x}_{nj} - \underline{\mu}_{kj}) \right)$$

$$0 = \sum_n q_{nk} \left( -\Sigma_{k'l}^{-1} + \sum_{i'j'} (\underline{x}_{ni} - \underline{\mu}_{ki}) (\underline{x}_{nj} - \underline{\mu}_{kj})^T \Sigma_{kj}^{-1} \right)$$

$$\Rightarrow \Sigma_{k'l} = \sum_n q_{nk} (\underline{x}_{ni} - \underline{\mu}_{ki}) (\underline{x}_{nj} - \underline{\mu}_{kj})^T / \sum_n q_{nk}$$

pre & post multiply by  $\Sigma_{k'l}$

To find  $\Pi$  we need to use Lagrange multipliers :

$$\frac{d}{d\Pi_{k'l}} \left( \varphi(q, \theta) - \lambda \left( \sum_k \Pi_k - 1 \right) \right) = \frac{d}{d\Pi_{k'l}} \left( \sum_k \left( q_{nk} \log \Pi_k - \lambda \Pi_k \right) - 1 \right)$$

$$= \sum_n q_{nk} \left( \frac{1}{\Pi_{k'l}} - \lambda \right) = 0$$

$$\Rightarrow \Pi_{k'l} = \frac{\sum_n q_{nk}}{\sum_n q_{nk}} \quad \text{from constraint that } \sum_k \Pi_k = 1$$

$$= \frac{1}{N} \sum_n q_{nk}$$

- b) i) The tree-energy is guaranteed to improve (or stay the same) if only a fraction of the datapoint posteriors are updated.
  - o. The procedure will converge to a local optimum of the likelihood as for Normal EM
- ii) It's not completely clear that the procedure will help in all cases - for although the cost of the E-Step is reduced (from  $O(N)$  to  $O(K)$ ) convergence will take more steps in general. However, when there are lots of data, you don't need to visit all of them to figure out a good update for the parameters so this means the procedure is likely to be substantially more efficient in these cases.

## Markov Models

### 5. Markov Models: fitting bi-gram models

A data scientist observes part of a long sequence that contains  $K = 3$  characters:  $ABAAABBABCCCCBC$ . She would like to use a bi-gram model to fit the data with parameters  $p(y_1 = k|\theta) = \pi_k^0$  and  $p(y_t = k|y_{t-1} = l, \theta) = T_{k,l}$ .

- Write down the log-likelihood for the model and optimise it with respect to  $\pi^0$  and  $T$  to find the maximum likelihood parameter estimates.
- Is the maximum-likelihood estimate sensible? How might you improve the estimate?

$$a) p(y_1:T | \pi^0, T) = \left( \prod_k \pi_k^0 \mathbb{I}(y_1=k) \right) \left( \prod_{t=2}^T \prod_{k=1}^K \prod_{l=1}^K \mathbb{I}(y_t=k, y_{t-1}=l) \right)$$

where

$$\mathbb{I}(x) = \begin{cases} 1 & \text{if } x \text{ is true} \\ 0 & \text{if } x \text{ is false} \end{cases}$$

called  
the indicator  
function

$$\begin{aligned} \therefore \log p(y_1:T | \pi^0, T) &= \sum_k \mathbb{I}(y_1=k) \log \pi_k^0 \\ &\quad + \sum_{t=2}^T \sum_{k=1}^K \sum_{l=1}^K \mathbb{I}(y_t=k, y_{t-1}=l) \log T_{kl} \\ &= \sum_k \mathbb{I}(y_1=k) \log \pi_k^0 + \sum_{kl} N_{kl} \log T_{kl} = \mathcal{L}(\theta) \end{aligned}$$

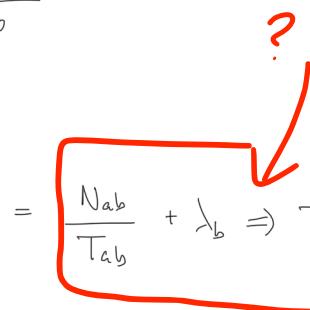
Maximum likelihood:

$$\frac{d}{d\pi_a^0} [\mathcal{L}(\theta) + \lambda \left( \sum_k \pi_k^0 - 1 \right)] = \frac{\mathbb{I}(y_1=a)}{\pi_a^0} + \lambda = 0$$

$$\Rightarrow \pi_a^0 = \mathbb{I}(y_1=a)$$

$$\frac{d}{dT_{ab}} \left[ \mathcal{L}(\theta) + \sum_{l=1}^K \lambda_l \left( \sum_{k=1}^K T_{kl} - 1 \right) \right] = \frac{N_{ab}}{T_{ab}} + \lambda_b \Rightarrow T_{ab} = \frac{N_{ab}}{\sum_a N_{ab}}$$

MLE =  
fraction of  
transitions from  
state b to state a



In the specific case given:

$$\underline{\Pi}^0 = [1, 0, 0]$$

$$\underline{T} = \begin{bmatrix} 2/5 & 2/5 & 0 \\ 3/5 & 1/5 & 1/3 \\ 0 & 2/5 & 2/3 \end{bmatrix}$$

$$\underline{N} = \begin{bmatrix} 2 & 2 & 0 \\ 3 & 1 & 1 \\ 0 & 2 & 2 \end{bmatrix}$$

↓ from  
A      B      C  
↑ to      ↓  
 $\sum_{ab} N_{ab}$

$$\begin{aligned} L(\theta) &= \log \underline{\Pi}_1^{(0)} + 2 \log T_{11} + 2 \log T_{12} + 3 \log T_{21} \\ &\quad + \log T_{22} + \log T_{23} + 2 \log T_{32} + 2 \log T_{33} \\ &= \log \underline{\Pi}_1^{(0)} + 2 [\log T_{11} + \log T_{12} + \log T_{32} + \log T_{33}] \\ &\quad + 3 \log T_{21} + \log T_{22} + \log T_{23} \end{aligned}$$

- b) The MLE is prone to over fitting (as the estimate for  $\underline{\Pi}^0$  is deterministic even though we have only seen one data point from  $\underline{\Pi}^0$ )

To avoid such pathologies in estimating  $\underline{\Pi}^0$  &  $\underline{T}$  we can use priors to get a MAP or Bayesian estimate.

## 6. Markov Models: Gaussian AR(1) models

A data scientist observes a sequence of scalar variables  $y_{1:T} = \{y_t\}_{t=1}^T$  generated from a Gaussian AR(1) process  $y_t = \lambda y_{t-1} + \epsilon_t$  where  $\epsilon_t \sim \mathcal{N}(\mu, \sigma^2)$ .

She knows that the invariant distribution of the process has the following properties

$$\lim_{t \rightarrow \infty} \mathbb{E}(y_t) = \mu_\infty, \quad \lim_{t \rightarrow \infty} \mathbb{E}(y_t^2) = \sigma_\infty^2 + \mu_\infty^2, \quad \lim_{t \rightarrow \infty} \mathbb{E}(y_t y_{t-1}) = \alpha_\infty.$$

- (a) Derive the parameters of the Gaussian AR(1) process  $\{\lambda, \mu, \sigma^2\}$  in terms of the properties of the invariant distribution  $\{\mu_\infty, \sigma_\infty^2, \alpha_\infty\}$ .
- (b) The data scientist reinterprets the original Markov model in terms of a new random variable  $z_t$  such that  $y_t = z_t + \mu_\infty$ . State the form of the distribution  $p(z_{1:T})$  required for this model to be equivalent to the original one.

Q2 a) First the long way!

$$y_t = \lambda y_{t-1} + \varepsilon_t \quad \varepsilon_t \sim N(\mu, \sigma^2)$$

① Taking expectations of both sides wrt to all sources of randomness

$$\langle y_t \rangle = \langle \lambda y_{t-1} \rangle + \langle \varepsilon_t \rangle = \lambda \langle y_{t-1} \rangle + \mu$$

If invariant distribution exists then

$$\mu_\alpha = \lambda \mu_\alpha + \mu \Rightarrow \mu_\alpha = \frac{\mu}{1-\lambda}$$

② Similarly

$$\begin{aligned} \langle y_t^2 \rangle &= \langle (\lambda y_{t-1} + \varepsilon_t)^2 \rangle = \lambda^2 \langle y_{t-1}^2 \rangle + \langle \varepsilon_t^2 \rangle \\ &\quad + 2\lambda \langle \varepsilon_t y_{t-1} \rangle \\ &= \lambda^2 \langle y_{t-1}^2 \rangle + \mu^2 + \sigma^2 + 2\mu \lambda \langle y_{t-1} \rangle \end{aligned}$$

If invariant distribution exists then

$$\sigma_\alpha^2 + \mu_\alpha^2 = \lambda^2 [\sigma_\alpha^2 + \mu_\alpha^2] + \mu^2 + \sigma^2 + 2\mu \lambda \mu_\alpha$$

$$\sigma_\alpha^2 + \mu_\alpha^2 = \frac{\mu^2 + \sigma^2 + 2\mu \lambda \mu_\alpha}{1 - \lambda^2}$$

$$\textcircled{3} \quad \langle y_t y_{t-1} \rangle = \langle (\lambda y_{t-1} + \varepsilon_t) y_{t-1} \rangle = \lambda \langle y_{t-1}^2 \rangle + \mu \mu_\alpha$$

If invariant distribution exists then;

$$\alpha_{\infty} = \lambda [\sigma_{\alpha}^2 + \mu_{\alpha}^2] + \mu_{\alpha} \mu_{\alpha}$$

rearrange:  
 $\lambda = \frac{\alpha_{\infty} - \mu_{\alpha} \mu_{\alpha}}{\sigma_{\alpha}^2 + \mu_{\alpha}^2}$  &  $\mu = \mu_{\alpha}(1-\lambda)$   
 eliminate  $\mu_{\alpha}$

So, a procedure for finding  $\lambda, \mu$  &  $\sigma^2$  is

$$① \quad \mu = \left( \sigma_{\alpha}^2 + \mu_{\alpha}^2 - \alpha_{\infty} \right) \frac{\mu_{\alpha}}{\sigma_{\alpha}^2}$$

$$② \quad \lambda = \frac{\alpha_{\infty} - \mu_{\alpha}^2}{\sigma_{\alpha}^2}$$

$$③ \quad \sigma^2 = (\sigma_{\alpha}^2 + \mu_{\alpha}^2)(1-\lambda^2) - \mu^2 - 2\mu \lambda \mu_{\alpha}$$

$$= \sigma_{\alpha}^2 (1-\lambda^2) + \mu_{\alpha}^2 [(1-\lambda^2) - (1-\lambda)^2 - 2\lambda(1-\lambda)]$$

$$= \sigma_{\alpha}^2 (1-\lambda^2) + \mu_{\alpha}^2 [4 - 2\lambda^2 - 1 + 2\lambda - 2\lambda + 2\lambda^2]$$

$$= \sigma_{\alpha}^2 (1-\lambda^2)$$

$$b) \quad y_t = z_t + \mu_{\alpha} = \lambda y_{t-1} + \varepsilon_t \quad \varepsilon_t \sim N(\mu, \sigma^2)$$

$$= \lambda (z_{t-1} + \mu_{\alpha}) + \varepsilon_t$$

$$\Rightarrow z_t = \lambda (z_{t-1} + \mu_{\alpha}) - \mu_{\alpha} + \varepsilon_t$$

$$= \lambda z_{t-1} + \underbrace{\lambda \mu_{\alpha} - \mu_{\alpha}}_{-\mu} + \mu + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma^2)$$

(see expression for  $\mu_{\alpha}$ )

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

$$\Rightarrow p(z_{1:T}) = \prod_t N(z_t; \lambda z_{t-1}, \sigma^2) \quad \begin{bmatrix} \text{standard} \\ \text{AR(1) model} \end{bmatrix}$$

This leads to a much faster way of doing part (a) & a useful way of checking this result.

$$\sigma_{\alpha}^2 = \langle z_t^2 \rangle \quad \text{[standard AR(1) process]}$$

For standard AR(1):

$$\sigma_{\alpha}^2 = \frac{\sigma^2}{1-\lambda^2} \quad \text{(see lectures)}$$

$$\therefore \sigma^2 = \sigma_{\alpha}^2 (1-\lambda^2) \quad \text{much faster!}$$

## *Hidden Markov Models*

### 7. Discrete Valued Hidden Markov Models\*

- (a) Provide the probabilistic equations that define a Hidden Markov Model (HMM) for observed data that takes discrete values. Indicate what aspects of the model the following terms refer to: *initial state probabilities*, *transition matrix* and *emission matrix*.
- (b) Consider a dataset consisting of the following string of 160 symbols from the alphabet  $\{A, B, C\}$ :

*AABBACABBBACAAAAAAAABBBACAAAAABACAAAAAA  
BBBBACAAAAAAAABACABACAABBACAAABBBBACA  
AABACAAAABACAABACAAABBACAAAABBBBACABBACAA  
AAAABACABACAAABACAABBACAAAABACABBACA*

Carefully analyse the string. Describe an HMM model for the string. Your description should include the number of states in the HMM, the transition matrix including the values of the elements of the matrix, the emission matrix including the values of its elements, and the initial state probabilities. Explain your reasoning.

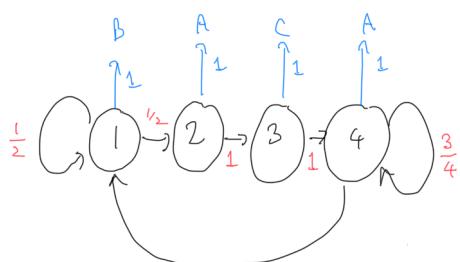
a) bookwork

$$p(x_t = k) = \pi_{tk} \quad \text{initial state probabilities}$$

$$p(x_t = k \mid x_{t-1} = l) = \tau_{kl} \quad \text{transition probabilities}$$

$$p(y_t = m \mid x_t = k) = s_{mk} \quad \text{emission probabilities}$$

b) BACAB common motif & only way 'C' is emitted & only way A follows B



Lots of choices for initial state probability  
eg  $p(1) = 1$        $p(k) \propto n_k$   
 ↑ number of times in that state

[There are a number of other possible solutions]

- only one symbol deterministically emitted from each latent state

- calculate state transition probabilities as follows

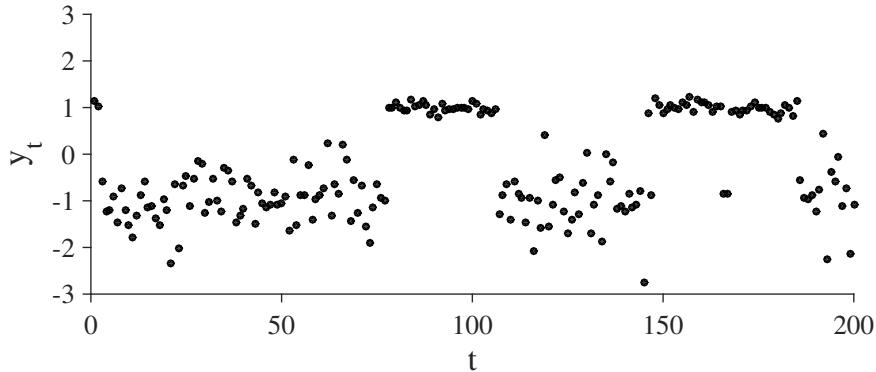
AA BB[BACA] BB[BACA] BB[BACA] AAAAAAAA BB  
 [BACA] AAAA [BACA] AAAAA BBB [BACA] AAAAAA  
 AAAA [BACA] [BACA] & B [BACA] AA BBB [BACA] AA [BACA]  
 AAA [BACA] A [BACA] AA B [BACA] AAA  
 BBB [BACA] B [BACA] AAAA B [BACA] [BACA] AA [BACA]  
 BBB [BACA] AAA [BACA] B [BACA]

from	to		total
(1)	(1)	(2)	45
	23	22	
(4)	(1)	(4)	76
	22	54	

[rough counts - rounded to fractions a human might choose]

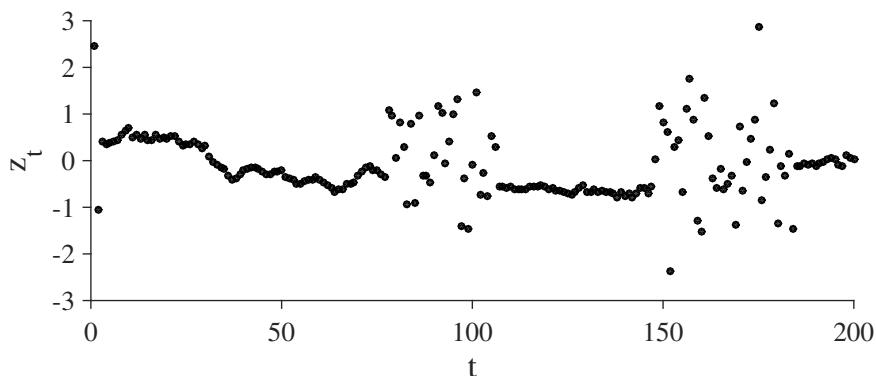
## 8. Probabilistic Modelling using HMMs for continuous valued observations

- (a) A machine learner observes the time-series,  $y_t$ , shown below:



Suggest a suitable Hidden Markov Model (HMM) for this sequence and state the model's probabilistic equations. Indicate plausible numerical values for the parameters where possible.

- (b) The machine learner is provided with a second set of observations  $z_t$  that were measured simultaneously with  $y_t$ , shown below:



Extend the HMM you proposed for part (a) so that it can jointly model the first and second set of observations.

a) binary hidden state  $s \in \{0, 1\}$

$$p(s_t=1) = \frac{1}{2} \quad p(s_t | s_{t-1}) = \begin{bmatrix} 0.99 & 0.01 \\ 0.01 & 0.99 \end{bmatrix}$$

$$p(y_t | s_t=1) = N(y_t; 1, 0.1^2)$$

$$p(y_t | s_t=0) = N(y_t; -1, (1.2)^2)$$

b)

$$p(z_t^{(1)}) = N(z_t^{(1)}; 0, 1)$$

This is known as a  
switching state-space  
model

$$p(z_t^{(1)} | z_{t-1}^{(1)}) = N(z_t^{(1)}; \lambda z_{t-1}^{(1)}, \begin{pmatrix} 1 \\ 2 \end{pmatrix}^T \begin{pmatrix} 1 & -\lambda^2 \end{pmatrix}) \quad \lambda = 0.99$$

$$z_t = s_t z_t^{(1)} + (1-s_t) z_t^{(2)} \quad \text{if } z_t = z_t^{(1)} \text{ if } s_t=1 \& z_t = z_t^{(2)} \text{ if } s_t=0$$

Again rough estimates for parameter values is fine, the general structure is  
the main thing to convey



## 9. Inference in HMMs with Discrete Hidden States<sup>†</sup>

A Hidden Markov Model contains a discrete hidden state variable  $x_t$  that takes one of two values and a discrete observed state  $y_t$  that also takes one of two values. The hidden state has a transition probability,

$$\begin{bmatrix} P(x_t = 1|x_{t-1} = 1) & P(x_t = 1|x_{t-1} = 2) \\ P(x_t = 2|x_{t-1} = 1) & P(x_t = 2|x_{t-1} = 2) \end{bmatrix} = \begin{bmatrix} T_{11} & T_{12} \\ T_{21} & T_{22} \end{bmatrix} = \begin{bmatrix} 2/3 & 1/3 \\ 1/3 & 2/3 \end{bmatrix}.$$

The filtering distribution at time  $t - 1$  is

$$P(x_{t-1}|y_{1:t-1}) = \begin{bmatrix} P(x_{t-1} = 1|y_{1:t-1}) \\ P(x_{t-1} = 2|y_{1:t-1}) \end{bmatrix} = \begin{bmatrix} 1/4 \\ 3/4 \end{bmatrix}.$$

- (a) Compute the predictive distribution for the next hidden state variable,  $P(x_t|y_{1:t-1})$ .
- (b) Explain how your solution to part (a) can be used to compute the filtering distribution  $P(x_t|y_{1:t})$ . What additional piece of information would you require to carry out this computation?

$$\begin{aligned}
&\text{QS a)} \\
&\quad \begin{bmatrix} p(x_t=1|x_{t-1}=1) & p(x_t=1|x_{t-1}=2) \\ p(x_t=2|x_{t-1}=1) & p(x_t=2|x_{t-1}=2) \end{bmatrix} \begin{bmatrix} p(x_{t-1}=1|y_{1:t-1}) \\ p(x_{t-1}=2|y_{1:t-1}) \end{bmatrix} \\
&= \begin{bmatrix} p(x_t=1|x_{t-1}=1)p(x_{t-1}=1|y_{1:t-1}) + p(x_t=1|x_{t-1}=2)p(x_{t-1}=2|y_{1:t-1}) \\ p(x_t=2|x_{t-1}=1)p(x_{t-1}=1|y_{1:t-1}) + p(x_t=2|x_{t-1}=2)p(x_{t-1}=2|y_{1:t-1}) \end{bmatrix} \\
&= \begin{bmatrix} p(x_t=1|y_{1:t-1}) \\ p(x_t=2|y_{1:t-1}) \end{bmatrix}
\end{aligned}$$

$\therefore$  Predictive is

$$\begin{bmatrix} 2/3 & 1/3 \\ 1/3 & 2/3 \end{bmatrix} \begin{bmatrix} 1/4 \\ 3/4 \end{bmatrix} = \begin{bmatrix} 2/12 + 3/12 \\ 1/12 + 6/12 \end{bmatrix} = \frac{1}{12} \begin{bmatrix} 5 \\ 7 \end{bmatrix}$$

b)  $p(x_t|y_{1:t}) \propto p(y_t|x_t) p(x_t|y_{1:t-1})$

$\uparrow$   
need the likelihood to get the posterior

## 10. Forecasting in Linear Gaussian State Space Models\*

A simple linear Gaussian state space model with scalar hidden state variables  $x_t$  has been used to model scalar observations  $y_t$ ,

$$p(x_t|x_{t-1}, \lambda, \sigma^2) = \mathcal{N}(x_t; \lambda x_{t-1}, \sigma^2), \quad p(y_t|x_t, \sigma_y^2) = \mathcal{N}(y_t; x_t, \sigma_y^2).$$

The Kalman filter recursions have been used to process  $T$  observations,  $y_{1:T}$ , in order to return the posterior distribution over the  $T$ th latent state,  $p(x_T|y_{1:T}) = \mathcal{N}(x_T; \mu_T, \sigma_T^2)$ .

- (a) Explain how to transform the posterior distribution over the  $T$ th latent state into a forecast for the observations one time step into the future, i.e. express  $p(y_{T+1}|y_{1:T})$  in terms of  $\mu_T$  and  $\sigma_T^2$ .
- (b) Now provide a forecast for the observations  $\tau$  time steps into the future by expressing  $p(y_{T+\tau}|y_{1:T})$  in terms of  $\mu_T$  and  $\sigma_T^2$ .
- (c) What happens to  $p(y_{T+\tau}|y_{1:T})$  as  $\tau \rightarrow \infty$ ?

a)

$$p(y_{T+1} | y_{1:T}) = N(y_{T+1}; \lambda \mu_T, \lambda^2 \sigma_T^2 + \sigma_y^2 + \sigma_\epsilon^2)$$

Calculated by passing  $N(x_T; \mu_T, \sigma_T^2)$  through  $x_{T+1} = \lambda x_T + \sigma \varepsilon_T$

& then noting  $y_{T+1} = x_{T+1} + \sigma_y \eta_T$  where  $\eta_T, \varepsilon_T \sim N(0, 1)$

$$\begin{aligned} b) \quad x_{T+\tau} &= \lambda x_{T+\tau-1} + \sigma \varepsilon_{T+\tau} \\ &= \lambda (\lambda x_{T+\tau-2} + \sigma \varepsilon_{T+\tau-1}) + \sigma \varepsilon_{T+\tau} \\ &= \lambda (\lambda (\lambda x_{T+\tau-3} + \sigma \varepsilon_{T+\tau-2}) + \sigma \varepsilon_{T+\tau-1}) + \sigma \varepsilon_{T+\tau} \\ &\vdots \\ x_{T+\tau} &= \lambda^\tau x_T + \sigma \sum_{t'=0}^{\tau-1} \lambda^{t'} \varepsilon_{T+\tau-t'} \end{aligned}$$

$\therefore p(y_{T+\tau} | y_{1:T}) = N(y_{T+\tau}; \lambda^\tau \mu_T, \sigma_y^2 + \sigma^2 \sum_{t'=0}^{\tau-1} \lambda^{2t'} + \lambda^{\tau} \sigma_\varepsilon^2)$

Geometric series:  $s_p = \sum_{t'=0}^{\tau-1} \lambda^{2t'} \quad s_{p+1} = \lambda^2 s_p + 1$

c) As  $\tau \rightarrow \infty$  the forecast will tend to the stationary distribution of the chain:

$$p(y_\infty | y_{1:T}) = N(y_\infty; 0, \frac{\sigma^2}{1-\lambda^2} + \sigma_y^2)$$

$$s_\infty = \lambda s_\infty + 1$$