

Paper 3F8: Inference

Solutions to Example Sheet 3: Sequence Modelling and Monte Carlo Methods

Straightforward questions are marked †

*Tripos standard (but not necessarily Tripos length) questions are marked **

Markov Models

1. Markov Models: fitting bi-gram models

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

- (a) Write down the log-likelihood for the model and optimise it with respect to π^0 and T to find the maximum likelihood parameter estimates.
- (b) Is the maximum-likelihood estimate sensible? How might you improve the estimate?

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

where

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

called \uparrow
the indicator
function

$$\therefore \log p(y_{1:T} | \pi^0, T) = \sum_k \mathbb{1}(y_1=k) \log \pi_k^0 + \sum_{t=2}^T \sum_{k=1}^K \sum_{l=1}^K \mathbb{1}(y_t=k, y_{t-1}=l) \log T_{kl}$$

$$= \sum_k \mathbb{1}(y_1=k) \log \pi_k^0 + \sum_{k,l} N_{kl} \log T_{kl} = \mathcal{L}(\theta)$$

number of transitions from
state l to state k
 \downarrow

Maximum likelihood:

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

$$\Rightarrow \pi_a^0 = \mathbb{1}(y_1=a)$$

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

MLE =
fraction of
transitions from b
that end up at a
 \downarrow

In the specific case given:

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

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

\downarrow from
 $N = \begin{matrix} & A & B & C \\ \begin{matrix} A \\ B \\ C \end{matrix} & \begin{bmatrix} 2 & 2 & 0 \\ 3 & 1 & 1 \\ 0 & 2 & 2 \end{bmatrix} \end{matrix}$
 \nearrow to
 $\sum_a N_{ab} \begin{matrix} 5 & 5 & 4 \end{matrix}$

$$\begin{aligned} \mathcal{L}(\theta) &= \log \Pi_1^{(0)} + 2 \log T_{11} + 2 \log T_{12} + 3 \log T_{21} \\ &\quad + \log T_{22} + \log T_{23} + 2 \log T_{32} + 2 \log T_{33} \\ &= \log \Pi_1^{(0)} + 2 [\log T_{11} + \log T_{12} + \log T_{32} + \log T_{33}] \\ &\quad + 3 \log T_{21} + \log T_{22} + \log T_{23} \end{aligned}$$

b) The MLE is prone to over fitting (eg 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.

2. 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

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

- Derive the parameters of the Gaussian AR(1) process $\{\lambda, \mu, \sigma^2\}$ in terms of the properties of the invariant distribution $\{\mu_\infty, \sigma_\infty^2, \alpha_\infty\}$.
- 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

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

② Similarly

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

if invariant distribution exists then

$$\sigma_{\infty}^2 + \mu_{\infty}^2 = \lambda^2 [\sigma_{\infty}^2 + \mu_{\infty}^2] + \mu^2 + \sigma^2 + 2\mu\lambda\mu_{\infty}$$

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

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

If invariant distribution exists then;

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

rearrange:

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

eliminate λ

So, a procedure for finding λ, μ & σ^2 is

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

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

sub (1) into $\lambda = \frac{\alpha_\infty - \mu_\infty \mu}{\sigma_\infty^2 + \mu_\infty^2}$

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

$$= \sigma_\infty^2(1-\lambda^2) + \mu_\infty^2[(1-\lambda^2) - (1-\lambda)^2 - 2\lambda(1-\lambda)]$$

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

$$= \sigma_\infty^2(1-\lambda^2)$$

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

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

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

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

(see expression for μ_∞)

$$z_t = \lambda z_t + \mu_t$$

$$\Rightarrow p(z_{1:t}) = \prod_t N(z_t; \lambda z_{t-1}, \sigma^2) \quad \left[\begin{array}{c} \text{standard} \\ \text{AR(1) model} \end{array} \right]$$

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

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

For standard AR(1):

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

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

Hidden Markov Models

3. 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 $\{A, B, C\}$:

AABBBACABBBACAAAAAAAAAABBBACAAAAABACAAAAAA
BBBBACAAAAAAAAAAAAABACABACAABBACAAABBBBACA
AABACAAAABACAABACAABBBACAAAABBBBACABBACAA
AAAABACABACAABACAABBBACAAAABACABBACA

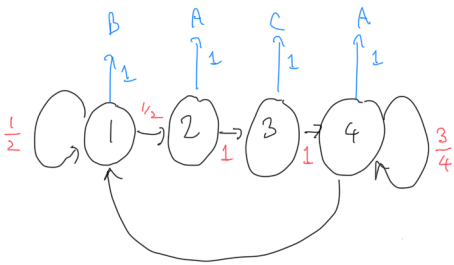
Carefully analyse the string. Describe an HMM model for the string. Your description should include the number of states in the HMM, the transition matrix including the values of the elements of the matrix, the emission matrix including the values of its elements, and the initial state probabilities.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

$$p(x_i = u) = \pi_u$$

initial state probabilities

$$p(x_t = k \mid x_{t-1} = l) = T_{kl} \quad \text{transition probabilities}$$
$$p(y_t = m \mid x_t = k) = S_{mk} \text{ emission probabilities}$$

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



Lots of choices for initial state probability $\frac{1}{4}$

$$\sum_k p(k) = 1 \quad p(k) \propto n_k$$

↑ 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] AAAAAA BB

[BACA] AAAA [BACA] AAAAA BBB [BACA] AAAAAA

AAAA[BACA][BACA]AB[BACA]AABBB[BACA]AA[BACA]

AAA[BACA]A[BACA]AAB[BACA]AAA

BBB(BACA)B(BACA)AAAA(BACA)(BACA)AA(BACA)

ABB [BACA] AAA [BACA] B [BACA]

from to (rough counts - rounded to fractions a human might choose)

①

①

②

total

45

23

22

④

(1)

4

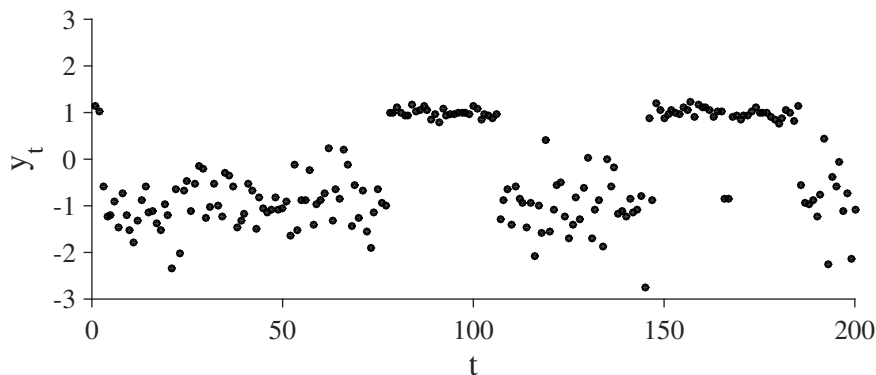
22

54

76

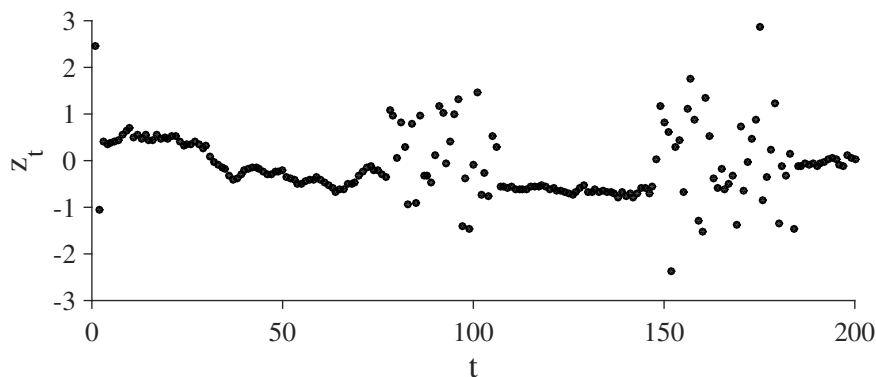
4. Probabilistic Modelling using HMMs for continuous valued observations

(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 $S \in \{0, 1\}$

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

$$p(y_t | s_t=1) = N(y_t; \overset{\text{mean}}{\underset{\downarrow}{1}}, \overset{\text{variance}}{\underset{\nwarrow}{0.1^2}})$$

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

b) $p(z_t^{(1)}) = N(z_t^{(1)}; 0, 1)$ { This is known as a switching state-space model }

$$p(z_t^{(2)} | z_{t-1}^{(2)}) = N(z_t^{(2)}; \lambda z_{t-1}^{(2)}, (\frac{1}{\lambda})^2 (1-\lambda^2)) \quad \lambda = 0.99$$

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

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

5. Inference in HMMs with Discrete Hidden States[†]

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

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

- Compute the predictive distribution for the next hidden state variable, $P(x_t | y_{1:t-1})$.
- 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?

Q5 a)

$$\begin{aligned} & \begin{bmatrix} P(x_t = 1 | x_{t-1} = 1) & P(x_t = 1 | x_{t-1} = 2) \\ P(x_t = 2 | x_{t-1} = 1) & P(x_t = 2 | x_{t-1} = 2) \end{bmatrix} \begin{bmatrix} P(x_{t-1} = 1 | y_{1:t-1}) \\ P(x_{t-1} = 2 | y_{1:t-1}) \end{bmatrix} \\ &= \begin{bmatrix} P(x_t = 1 | x_{t-1} = 1) P(x_{t-1} = 1 | y_{1:t-1}) + P(x_t = 1 | x_{t-1} = 2) P(x_{t-1} = 2 | y_{1:t-1}) \\ P(x_t = 2 | x_{t-1} = 1) P(x_{t-1} = 1 | y_{1:t-1}) + P(x_t = 2 | x_{t-1} = 2) P(x_{t-1} = 2 | y_{1:t-1}) \end{bmatrix} \\ &= \begin{bmatrix} P(x_t = 1 | y_{1:t-1}) \\ P(x_t = 2 | y_{1:t-1}) \end{bmatrix} \end{aligned}$$

∴ Predictive is

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

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

↑↑
need the likelihood to get the posterior

6. 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 μ_T and σ_T^2 .
- (b) Now provide a forecast for the observations τ time steps into the future by expressing $p(y_{T+\tau}|y_{1:T})$ in terms of μ_T and σ_T^2 .
- (c) What happens to $p(y_{T+\tau}|y_{1:T})$ as $\tau \rightarrow \infty$?

a)

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

Calculated by passing $N(x_T ; \mu_T, \sigma_T^2)$ through $x_{T+1} = \lambda x_T + \sigma \varepsilon_T$
 & then noting $y_{T+1} = x_{T+1} + \sigma_y n_T$ where $\varepsilon_T, n_T \sim N(0, 1)$

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

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

geometric series: $S_r = \sum_{t'=0}^{r-1} \lambda^{2t'}$ $S_{r+1} = \lambda^2 S_r + 1$

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

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

$$S_\infty = \lambda^2 S_\infty + 1$$

Selected solutions and hints

1. a) $\log p(y_{1:T}|\theta) = \log \pi_1^0 + 2 \log(T_{11}T_{12}T_{32}T_{33}) + 3 \log T_{21} + \log(T_{22}T_{23})$

$$T = \begin{bmatrix} 2/5 & 2/5 & 0 \\ 3/5 & 1/5 & 1/3 \\ 0 & 2/5 & 2/3 \end{bmatrix}, \quad \pi^0 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

2. a) $\lambda = (\alpha_\infty - \mu_\infty^2)/\sigma_\infty^2, \mu = \mu_\infty(1 + (\mu_\infty^2 - \alpha_\infty)/\sigma_\infty^2), \sigma^2 = \sigma_\infty^2 \left(1 - \left(\frac{\alpha_\infty - \mu_\infty^2}{\sigma_\infty^2}\right)^2\right)$

3. b) pay close attention to repeated patterns and remember that some parts of a HMM can be deterministic

4. b) consider whether the low variance z_t regions are correlated through time and whether a standard HMM could model this

5. a) $p(x_t|y_{1:t-1}) = [5, 7]^\top / 12$

6. b) $p(y_{T+\tau}|y_{1:T}) = \mathcal{N}(y_{T+\tau}; \lambda^\tau \mu_T, \sigma_y^2 + \lambda^{2\tau} \sigma_T^2 + \sigma^2 \sum_{t'=0}^{\tau-1} \lambda^{2t'})$