

## Unit 5 Hidden Markov model for Time-Series

### Lesson 1 Markov Models and Hidden Markov Model's

1. The Markov assumption

2. Markov models for discrete sequences

3. Hidden Markov models

#### The Markov assumption

We want more flexibility than assuming each timestep is iid - but more simpler than letting number of parameters grow with seq length T.

First order Markov assumption:

$z_{t+1}$  is conditionally independent of  $z_1, z_2, \dots, z_{t-1}$  given  $z_t$

ex:

$$p(z_1, z_2, z_3) = p(z_1)p(z_2|z_1)p(z_3|z_1, z_2) \quad \text{under product rule}$$

$$p(z_3|z_2, z_1) = p(z_3|z_2) \quad \text{under Markov assumption}$$

$$p(z_1, \dots, z_T) = p(z_1) \cdot \prod_{t=2}^T p(z_t|z_{t-1})$$

make identical distribution assumption across time

$$\begin{aligned} p(z_2=k|z_1=j) &= A_{jk} \\ p(z_3=k|z_2=j) &= A_{jk} \\ &\vdots \\ p(z_T=k|z_{T-1}=j) &= A_{jk} \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Share same parameters } A \text{ for all steps}$$

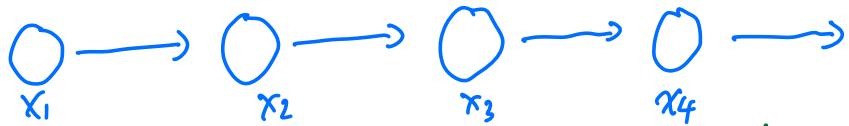
Summary:

(1) 1<sup>st</sup> order Markov:  $p(z_{t+1}|z_t, z_{t-1}, \dots) = p(z_{t+1}|z_t)$

(2) parameter sharing:  $p(z_{t+1}=k|z_t=j) = p(z_2=k|z_1=j)$

achieve simple yet tractable model

1<sup>st</sup> order Markov model



Random Variable: Sequence  $z_1, z_2, \dots, z_T$  with each  $z_t \in \{1, 2, \dots, k\}$

Sample space: All possible sequence of length T using the K symbols.

Parameter:

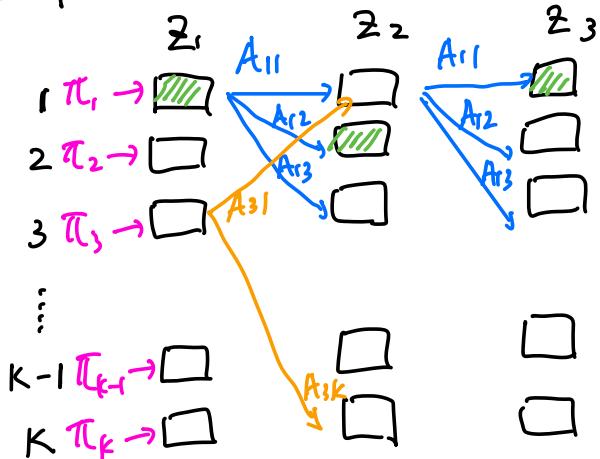
1. initial State Prob:

$$p(z_1 = k) = \pi_k \quad z_1 \sim \text{Cat}(\pi_1, \dots, \pi_K)$$

2. transition prob:

$$p(z_t = k | z_{t-1} = j) = A_{jk}$$

Example:



$$p(1, 2, 1) = \pi_1 A_{12} A_{21}$$

Cat: 单次采样, 3个结果 {0, 1, 2, ..., n}

Bern: 单次采样, 2个结果 {0, 1}

multi: 多次采样, 3个结果

Exercise:

① What is the marginal  $p(z_1)$ ?

② What is the marginal  $p(z_T)$ ?

## Hidden Markov Model

Goal: Tractable model for observed data sequence  $x_1, x_2, \dots, x_T$  with  $x_t \in R$  that has dependence between  $x_t$  and  $x_1, x_2, \dots, x_{t-1}$ , but affordable to do 2 tasks:

(1) Compute data likelihoods:

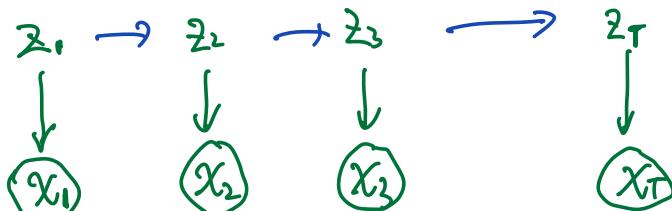
- joint  $p(x_1, x_2, \dots, x_T | \theta)$  Evaluation Forward Backward
- conditional  $p(x_T | x_1, x_2, \dots, x_{T-1}, \theta)$  Decoding Viterbi

(2) Estimate parameters  $\theta$  via penalized Maximum Likelihood.

Learning EM

Idea:

hidden discrete state:



observed data :

$$p(x_{1:T}, z_{1:T}) = \underbrace{p(z_{1:T})}_{\text{First order Markov model}} \cdot \underbrace{p(x_{1:T} | z_{1:T})}_{\text{each timestep iid given } z_{1:T}}$$

$$= \left[ p(z_1) \prod_{t=2}^N p(z_t | z_{t-1}) \right] \cdot \prod_{t=1}^T p(x_t | z_t)$$

$$= \text{Cat}(z_1 | \pi) \prod_{t=2}^T \text{Cat}(z_t | A z_{t-1}) \prod_{t=1}^T \text{Norm PDF}(x_t | \mu_{z_t}, \sigma_{z_t}^2)$$

with parameter  $\theta = \underbrace{\pi, A}_{\text{used for } p(z)} , \underbrace{\mu, \sigma^2}_{\text{used for } p(x|z)}$

Markov

## Lesson 2 EM for HMMs

1. EM for HMM parameter estimation
2. Expected log-likelihood for HMMs
3. M-step for HMMs
4. E-step for HMMs: overview
5. E-step for HMMs: forward-backward algorithm

### Overview of EM for HMMs

Goal: Estimate parameters of an HMM using ML estimation

$$\max_{\pi, A, \mu, \sigma} \log p(x_1, x_2, \dots, x_T | \pi, A, \mu, \sigma)$$

Given: observed sequence  $x_1, x_2, \dots, x_T$

Output:  $\pi$ : initial probability

$A$ : transition probabilities  $K \times K$  matrix, row sum to 1

$\mu$ : means

$\sigma$ : std deviations

Notation:  $\Theta = \{\pi, A, \mu, \sigma\}$

Challenges:

how to compute this likelihood?

• complete likelihood is easy:  $p(X_{1:T}, Z_{1:T} | \Theta)$

• incomplete likelihood is hard:  $p(X_{1:T} | \Theta) = \sum_{Z_{1:T}} p(X_{1:T}, Z_{1:T} | \Theta)$

Complete likelihood:

$$p(X_{1:T}, Z_{1:T}) = p(Z_{1:T}) p(X_{1:T} | Z_{1:T}) \quad *$$

Assume  $T=3$ ,  $Z_1=a$ ,  $Z_2=b$ ,  $Z_3=c$

$$* = p(z_{1:3}) \prod_{t=1}^3 p(x_t | z_t) \\ = p(a) \cdot p(b|a) \cdot p(c|b) \cdot \prod_{t=1}^3 p(x_t | z_t)$$

incomplete likelihood:

$$p(x_{1:3}) = \sum_a \sum_b \sum_c p(x_{1:3}, z_1=a, z_2=b, z_3=c)$$

Big idea:

(1) Let  $q(z|s)$  be an "approx. posterior". Define a valid distribution over the sequence  $z = z_1, z_2, \dots, z_T$

(2) Use the lower bound objective  $L$ :

$$\log p(x|\theta) = E_{q(z|s)} [\log p(x, z|\theta) - \log q(z|s)] \approx L(x, s, \theta)$$

(3) Iteratively optimize lower bound using coordinate ascent.

Init:  $\theta^0$

for iteration  $i=1, 2, \dots$

$$\text{E-step: } \hat{s}_{1:T}^{(i)} \leftarrow \underset{s_{1:T}}{\operatorname{arg\max}} \mathcal{L}(x_{1:T}, s_{1:T}, \theta^{(i-1)})$$

$$\text{M-step: } \theta^{(i)} \leftarrow \underset{\theta}{\operatorname{arg\max}} \mathcal{L}(x_{1:T}, \hat{s}_{1:T}^{(i)}, \theta)$$

Punchline: Can do all key steps in affordable runtime  $O(Tk^2)$  or better. E step, M step, L calculation are all tractable

Define  $q(z)$  and compute expected log likelihood (E-step)

Q: What is  $q(z)$ ? What is parameter  $S$ ?

We can specify proba of adjacent timestamps:

$$q(z_t=k, z_{t+1}=l) = S_{tkl}$$

Parameter:  $S_t : K \times K$  matrix

$$\text{must satisfy } S_{tkl} \geq 0, \sum_k S_{tkl} = 1$$

Besides, for a valid  $q$ , neighbouring pairs should be consistent!

for a single timestamp:

$$p(z_t=k) = \sum_{l=1}^K (z_t=k, z_{t+1}=l) = \sum_l S_{tkl}$$

$$= \sum_{j=1}^K (z_{t-1}=j, z_t=k) = \sum_j S_{jtk}$$

That is ① Valid PMF over  $K \times K$   
② Consistent with Neighbours

Recall:

$$\begin{aligned} p(X_{1:T}, Z_{1:T}) &= P(Z_{1:T}) \cdot P(X_{1:T} | Z_{1:T}) \\ &= P(Z_1) \cdot \prod_{t=2}^T p(z_t | z_{t-1}) \cdot \prod_{t=1}^T p(x_t | z_t) \\ &= p(z_1 | \pi) \prod_{t=2}^T p(z_t | a_{t-1}) \cdot \prod_{t=1}^T N(x_t | \mu_t, \sigma_t^2) \end{aligned}$$

Compute the log-likelihood:

(Know  $z$ )

$$\log p(X_{1:T}, Z_{1:T}) = \log p(Z_{1:T}) + \log p(X_{1:T} | Z_{1:T})$$

Let  $Z$  in One-hot form: (written in  $z'$ )

$$\begin{aligned} \log p(X_{1:T} | Z) &= \log \prod_{t=1}^T \prod_{k=1}^K N(x_t | \mu_k, \sigma_k^2)^{z'_{tk}} \\ &= \sum_{t=1}^T \sum_{k=1}^K z'_{tk} N(x_t | \mu_k, \sigma_k^2) \end{aligned}$$

(linear in  $T$   
linear in  $K$ )

$$\log P(Z_{1:T}) = \log [P(Z_1|\pi) \prod_{t=2}^T P(Z_t|A)] = \log(Z_1|\pi) + \sum_{t=2}^T \log P(Z_t|A)$$

①

$$P(Z_1|\pi) = \frac{\pi}{K} \prod_{k=1}^K \pi_k^{z_{1k}}$$

$$P(Z_t|Z_{t-1}, A) = \frac{1}{K} \sum_{j=1}^K A_{jk} \underbrace{z'_{tk} z'_{t-1,j}}_{=1 \text{ when } \begin{cases} Z_t=k \\ Z_{t-1}=j \end{cases}}$$

$$\Rightarrow \log P(Z_{1:T}|\theta) = \sum_{k=1}^K z_{1k} \log \pi_k + \sum_{t=2}^T \sum_{j=1}^K \sum_{k=1}^K z'_{tk} z'_{t-1,j} \log A_{jk}$$

linear in T  
quadratic in K  $O(TK^2)$

Compute the expectation of log likelihood: (only know  $q(z)$ )

$$E_{q(z|s)} \log P(X_{1:T}, Z_{1:T})$$

2 facts:

$$E_{q(z|s)}[z'_{tk}] = \sum_k S_{tk} \stackrel{\text{Notation}}{=} r_{tk} \quad ?$$

$$E_{q(z|s)}[z'_{tk} z'_{t+1,k}] = S_{tk}$$

Given:

$$\begin{aligned} \log P(X_{1:T}, Z_{1:T}) &= \log(P(X|Z)) + \log(P(Z|\theta)) \\ &= \sum_{k=1}^K z_{1k} \log \pi_k + \sum_{t=2}^T \sum_{j=1}^K \sum_{k=1}^K z'_{tk} z'_{t-1,j} \log A_{jk} \\ &\quad + \sum_{t=1}^T \sum_{k=1}^K z'_{tk} N(\dot{x}_t | \mu_k, \sigma_k^2) \end{aligned}$$

$$\begin{aligned} E(*) &= \sum_{k=1}^K r_{1k} \log \pi_k + \sum_{t=2}^T \sum_{j=1}^K \sum_{k=1}^K S_{tjk} \log A_{jk} + \\ &\quad \sum_{t=1}^T \sum_{k=1}^K r_{tk} N(x_t | \mu_k, \sigma_k^2) \end{aligned}$$

Summary:

1. Can compute  $E[\log p(x, z)]$  easily using  $S$  of  $g(z|s)$
2. We have defined useful notation for marginals at each time  $t$ .  
 $q(z_t=k) = r_{tk} = r_{tk}(s) = \begin{cases} \sum_{\ell=1}^k s_{t\ell k} & \text{for } t=1, 2, \dots, T-1 \\ \sum_{j=1}^k s_{T+1, j, k} & \text{for } t=T \end{cases}$

### M step for HMM

Using simplified expression for expected complete likelihood, we can see  
M-step takes as input:

-  $r_{tk}$  probability of assigning timestep  $t$  to cluster  $k$

$$r_{tk} = \sum_k s_{t\ell k} \quad (\text{deterministic given } S)$$

-  $s_{tjk}$  probability of assigning  $z_t$  to  $j$  and  $z_{t+1}$  to  $k$

$r, s \in t$ -step specific parameters

$\pi, A, \mu, \sigma$ ; global parameters for HMM

Given  $r, s$ , we can see M-step is simplified:

$$E(*) = \sum_{k=1}^K r_{1k} \log \pi_k + \sum_{t=2}^T \sum_{j=1}^K \sum_{k=1}^K s_{tjk} \log A_{jk} + \sum_{t=1}^T \sum_{k=1}^K r_{tk} N(x_t | \mu_k, \sigma_k^2)$$

M step:  $\arg \max_{\theta} L(x, s, \underline{\theta}) = \arg \max_{\theta} E(*)$

$$\pi^* = \arg \max \sum_{k=1}^K r_{1k} \log \pi_k \Rightarrow \pi_k = \frac{r_{1k}}{T} \quad \begin{array}{l} \# \text{times cluster } k \\ \text{to 1st step} \end{array} \Rightarrow r_i \in \Delta^K$$

$$A_j^* \leftarrow \arg \max_{A_j \in \Delta^K} \sum_{t=1}^{T-1} \sum_K S_{tjk} \log A_{ijk} \quad A_{jk}^* = \frac{\sum_t S_{tjk}}{\sum_t \sum_k S_{tjk}} \begin{array}{l} \# \text{ time } j \rightarrow k \\ \# \text{ time } j \rightarrow \text{anything} \\ A_j \text{ sums to 1} \end{array}$$

$$\mu^*, \sigma^* \leftarrow \arg \max_{\mu_k, \sigma_k} \sum_t r_{tk} \log N(x_t | \mu_k, \sigma_k)$$

$\uparrow$   
expected # time state k

Summary:

1.  $S_t \rightarrow$  adjacent time step joint  
 $r_t \rightarrow$  single time step margin

2. E step  $\rightarrow S_t, r_t$



How to do the E-step?

$$\text{Recall: } \log p(x|\theta) \geq L(X, S, \theta) + KL(\underline{q(z|S)} \mid \underline{p(z|x, \theta)})$$

↓ Posterior

Best possible E step update would make  $KL=0$

and thus  $\log p(x|\theta) = L(X, S, \theta)$

This achieve by  $q(z|S) = p(z|x, \theta)$

In other word. we match our learned distribution  $q$  to the hidden-given data posterior  $p(z|x, \theta)$

while we could derive the optimal update by solving

$$S^* = \arg \max_S L(X, S, \theta)$$

S that meet sum  
to one and neighbor  
consistency constraints

We could find the same optimal  $S^*$ , "matching the posterior" will be simpler.

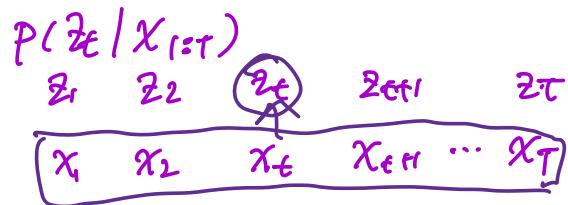
- procedure: Analyse the posterior  $\text{PL} Z(X_{1:T})$  - specifically its moments for margins  $t : \underline{\text{PL} Z_t | X_{1:T}}$   
 pairwise joints  $[t, t+1] : \underline{\text{PL} Z_t, Z_{t+1} | X_{1:T}}$

We'll see both can be computed exactly via dynamiz programmable

### Single Timestep Marginal Posterior

For each timestep  $t$ , we have:

$$\text{PL} Z_t | X_{1:T} = \frac{\text{PL} Z_t, X_{1:T}}{\text{P}(X_{1:T})}$$



$$= \frac{\text{PL} X_{1:t}, Z_t \cdot \text{P}(X_{t+1:T} | X_{1:t}, Z_t)}{\text{PL} X_{1:T}} \quad \text{product rule}$$

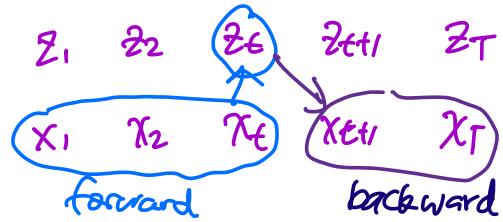
$$= \text{P}(X_{1:t}, Z_t) \cdot \text{P}(X_{t+1:T} | Z_t) \quad \begin{matrix} \text{HMM Conditional Independence} \\ \text{assumption} \end{matrix}$$

$$= \frac{\text{PL} X_{1:t} \cdot \text{PL} Z_t | X_{1:t} \cdot \text{PL} X_{t+1:T} | Z_t}{\text{PL} X_{1:T}}$$

constant

$$\alpha_t^k = \text{PL} Z_t=k | X_{1:t} \rightarrow \text{forward}$$

$$\beta_t^k = \text{PL} X_{t+1:T} | Z_t=k \rightarrow \text{backward}$$



Thus:

$$\text{PL} Z_t | X_{1:T} = \frac{\alpha_t^k \cdot \beta_t^k}{\sum_l \alpha_{tl} \beta_{tl}} \quad \leftarrow \text{normalization}$$

## • Adjacent Timestep Joint Posterior

For each timestep  $t$  in  $1, 2, \dots, T-1$ , we have

$$p(z_t, z_{t+1} | x_{1:T}) = \frac{p(z_t, z_{t+1}, x_{1:T})}{p(x_{1:T})}$$

$$= \frac{1}{p(x_{1:T})} \cdot p(x_{1:t+2} | \cancel{x_{1:t+1}}, \cancel{z_t}, \cancel{z_{t+1}}) \cdot \underline{p(x_{1:t+1}, z_t, z_{t+1})} \quad *$$

$$\begin{aligned} p(x_{1:t}, z_t, z_{t+1}) &= p(x_{t+1} | z_{t+1}, \cancel{z_t}, \cancel{x_{1:t}}) \cdot p(z_{t+1}, z_t, x_{1:t}) \\ &= p(x_{t+1} | z_{t+1}) \cdot p(z_{t+1} | z_t, x_{1:t}) \cdot p(z_t, x_{1:t}) \\ &= p(x_{t+1} | z_{t+1}) \cdot p(z_{t+1} | z_t) \cdot p(z_t | x_{1:t}) \cdot p(x_{1:t}) \end{aligned}$$

$$* \Rightarrow \cancel{\frac{1}{p(x_{1:T})}} \cdot \frac{p(x_{t+1} | z_{t+1})}{p_{t+1}} \cdot \frac{p(x_{t+1} | z_{t+1})}{N(x_{t+1} | \mu_{z_{t+1}}, \sigma_{z_{t+1}}^2)} \cdot \frac{p(z_{t+1} | z_t)}{A_{z_t, z_{t+1}}} \cdot \frac{p(z_t | x_{1:t})}{\partial_t}$$

$$p(z_t=j, z_{t+1}=k | x_{1:T}) = \frac{p_{t+1,k} \cdot L_{t+1,k} \cdot A_{j-k} \cdot \partial_{t,j}}{\sum_m p_{t+1,m} \cdot L_{t+1,m} \cdot A_{e,m} \cdot \partial_{t+1}}$$

$$\text{where } L_{t,k} = p(x_t | z_t=k, \theta) = N(x_t | \mu_k, \sigma_k^2)$$

Punchline:

$$\text{Given } \alpha, \beta, \text{ can set } S_{t,k,l} = p(z_t=k, z_{t+1}=l)$$

$$= \frac{\alpha \cdot \beta}{\text{Normalization}}$$

forward and backward algorithm

Computing forward messages & via Dynamic Programming (Forward Algo)