

SPR Day 18

Hidden Markov Models Parameter Estimation and Likelihood Computation

Reading: Bishop PRML 13.2.1
13.2.2

Outline:

- Overview of EM for HMMs
- Defining $q(z)$ and computing ^{expected} _{log} likelihood
- M step overview
- E step intro
 - FORWARD algorithm
 - BACKWARD algorithm
- Recap of E-step

2

Goal: Estimate parameters of an HMM using maximum likelihood method

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

each $x_t \in \mathcal{R}$

Given: observed sequence $x_1, x_2, x_3, \dots, x_T$

Output: π initial probabilities
 A transition probabilities

μ means
 σ stdeviations

Notation

$$\theta = \{ \pi, A, \mu, \sigma \}$$

All HMM parameters in one symbol θ

Challenges:

- how to compute this likelihood?

Complete is Easy: $p(x_{1:T}, z_{1:T} | \theta)$
likelihood

Incomplete is Hard: $p(x_{1:T} | \theta) = \sum_{z_{1:T} \in \Omega} p(x_{1:T}, z_{1:T} | \theta)$
likelihood (sum rule)

Ω denotes all possible sequences of length T using symbols $\{1, 2, \dots, K\}$

Number of terms in sum is $|\Omega| = K^T$
grows exponentially with T (# timesteps)
Not easy! What else can we do?

Big Idea:

(1) Let $q(z/s)$ be an "approximate" posterior, defining a valid distribution over the sequence $z = z_1, z_2, \dots, z_T$

(2) Use the lower bound objective \mathcal{L}
$$\log p(x/\theta) \geq \int_{q(z/s)} [\log p(x, z/\theta) - \log q(z/s)] = \mathcal{L}(x, s, \theta)$$

(3) Iteratively optimize lower bound using coordinate ascent

Init: θ^0

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

E-step: $s_{1:T}^i \leftarrow \operatorname{argmax}_{s_{1:T}} \mathcal{L}(x_{1:T}, s_{1:T}, \theta^{i-1})$

M-step: $\theta^i \leftarrow \operatorname{argmax}_{\theta} \mathcal{L}(x_{1:T}, s_{1:T}^i, \theta)$

Punchline: Can do all key steps in affordable routine $O(TK^2)$ or better
E-step, M-step, \mathcal{L} calculation are all tractable

Idea: Define a tractable distribution

$q(z_{1:T} | s)$ over sequences z_1, z_2, \dots, z_T
with each $z_t \in \{1, 2, \dots, K\}$
If we want onehot notation, write vector $\text{ONEHOT}(z_t)$

How? Define joint probability at each adjacent pair of timesteps: $(z_1, z_2), (z_2, z_3), \dots, (z_{T-1}, z_T)$
for $t=1, 2, \dots, T-1$: $q(z_t=j, z_{t+1}=k) = S_{tjk}$

Requirements:

$S_t: K \times K$ matrix
non-negative
sums to 1

$\sum_j \sum_k S_{tjk} = 1$ for $t=1, 2, \dots, T-1$

$S_{tjk} \geq 0$ for $t=1, 2, \dots, T-1$
 $j \in \{1, \dots, K\}$
 $k \in \{1, \dots, K\}$

① each pairwise joint is valid PMF over $K \times K$ outcomes

② neighboring pairs have consistent marginals

for all t ,

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

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

Using this distribution, we can compute:

$$E_{q(z|s)} [\text{ONEHOT}(z_t)_k] = \sum_{l=1}^K S_{tkl}$$

for $t=1, 2, \dots, T$

$$E_{q(z|s)} [\text{ONEHOT}(z_{t-1})_j \cdot \text{ONEHOT}(z_t)_k] = S_{tjk}$$

for $t=1, 2, \dots, T-1$

scalar \cdot scalar

If we know z : complete log likelihood for HMM 5

$$\begin{aligned} \log p(z_{1:T}, x_{1:T} | \theta) &= \log p(z_{1:T} | \theta) + \log p(x_{1:T} | z_{1:T}, \theta) \\ &= \log \text{CatPMF}(z_1 | \pi) \\ &\quad + \sum_{t=1}^{T-1} \log \text{CatPMF}(z_{t+1} | \text{ONEHOT}(z_t)^T A) + \sum_{t=1}^T \sum_{k=1}^K \text{ONEHOT}(z_t)_k \log \text{Norm}(x_t | \mu_k, \Sigma_k^2) \end{aligned}$$

$$\begin{aligned} &= \sum_{k=1}^K \text{ONEHOT}(z_1)_k \log \pi_k \\ &\quad + \sum_{t=1}^{T-1} \sum_{j=1}^K \sum_{k=1}^K \text{ONEHOT}(z_t)_j \text{ONEHOT}(z_{t+1})_k \log A_{jk} + \sum_{t=1}^T \sum_{k=1}^K \text{ONEHOT}(z_t)_k \log \text{Norm}(x_t | \mu_k, \Sigma_k^2) \end{aligned}$$

If we only know $q(z)$, replace above with expectation

$$\begin{aligned} \mathbb{E}_{q(z|s)} [\log p(x, z)] &= \sum_{k=1}^K r_{1k}(s) \log \pi_k \\ &\quad + \sum_{t=1}^{T-1} \sum_{j=1}^K \sum_{k=1}^K s_{tjk} \log A_{jk} + \sum_{t=1}^T \sum_{k=1}^K r_{tk}(s) \log \text{Norm}(x_t | \mu_k, \Sigma_k^2) \end{aligned}$$

We have defined useful notation for marginals at each time t

$$q(z_t = k) = r_{tk} = r_{tk}(s) = \begin{cases} \sum_{l=1}^K s_{tkl} & \text{for } t=1, 2, \dots, T-1 \\ \sum_{j=1}^K s_{t-1,jk} & \text{for } t=T \end{cases}$$

M-Step for HMMs

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

- r_{tk} probability of assigning timestep t to cluster k (deterministic given S)
 $r_{tk} = \sum_z s_{ztk}$
- s_{tijk} probability of assigning z_t to j and z_{t+1} to k

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

$$\pi^* \leftarrow \operatorname{argmax}_{\pi \in \Delta^K} \sum_k r_{1k} \log \pi_k \quad \boxed{\pi_k^* = \frac{r_{1k}}{1}}$$

$$A_{ij}^* \leftarrow \operatorname{argmax}_{A_j \in \Delta^K} \sum_t \sum_k s_{tijk} \log A_{jk} \quad \boxed{A_{jk}^* = \frac{\sum_t s_{tijk}}{\sum_t \sum_k s_{tijk}}}$$

$$\mu_k^*, \sigma_k^* \leftarrow \operatorname{argmax}_{\mu_k, \sigma_k} \sum_z r_{zk} \log \operatorname{Norm}(x_t / \mu_k, \sigma_k)$$

μ_k^* same as S
 σ_k^* GMM M-step using $r = r(S)$

Note: likely want to use Penalty or MAP for σ and π , maybe also A

How to do the E-step?

Recall: $\log p(x|\theta) \geq \alpha(x, s, \theta) + \text{KL}(q(z|s) | p(z|x, \theta))$

lower bound objective KL term ≥ 0

Best possible E step update would make $\text{KL} = 0$
and thus $\log p(x|\theta) = \alpha(x, s, \theta)$ [bound is tight]

This is achieved by finding s such that
 $q(z|s) = p(z|x, \theta)$

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

While we could, ^{instead} derive the optimal update by solving

$$s^* = \operatorname{argmax}_s \alpha(x, s, \theta)$$

s that meet
sum to one
and
neighbor consistency
constraints

We would find the same optimal s^* , "matching the posterior" will be simpler.

Procedure: Analyse the posterior $p(z|x, \theta)$, specifically its moments for marginals $(t): p(z_t | x_{1:T})$ and pairwise joints $(t, t+1): p(z_t, z_{t+1} | x_{1:T})$

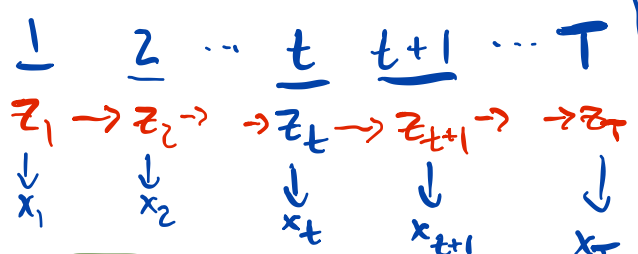
We'll see both can be computed exactly via dynamic programming

Single Timestep Marginal Posterior 8

For each timestep t , we have:

$$P(z_t | x_{1:T}) = \frac{P(x_{1:T}, z_t)}{P(x_{1:T})} \quad \text{by Bayes rule}$$

↑ which of K states at time t
 ↑ given all data from start to end of sequence



x_{t+1} is conditionally indep of x_t given z_t

$$= \frac{P(x_{1:t}, z_t) P(x_{t+1:T} | x_{1:t}, z_t)}{P(x_{1:T})} \quad \text{product rule}$$

$$= \frac{P(x_{1:t}, z_t) P(x_{t+1:T} | z_t)}{P(x_{1:T})} \quad \text{via HMM conditional independence assumption}$$

$$= \frac{\overset{\text{const}}{P(x_{1:t})} P(z_t | x_{1:t}) P(x_{t+1:T} | z_t)}{\overset{\text{const}}{P(x_{1:T})}}$$

Path forward: Suppose we can compute

$$\alpha_{tk} = P(z_t = k | x_{1:t}) \quad \text{for } t=1,2,\dots,T \text{ and } k=1,2,\dots,K$$

$$\beta_{tk} = P(x_{t+1:T} | z_t = k) \quad \text{for } t=1,2,\dots,T \text{ and } k=1,2,\dots,K$$

Then, we can compute single timestep posterior marginal as:

$$P(z_t = k | x_{1:T}) = \frac{\alpha_{tk} \beta_{tk}}{\sum_l \alpha_{tl} \beta_{tl}}$$

Adjacent Timestep Joint Posterior

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

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

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

product rule & x_{t+2} indep of $x_{1:t+1}$ given z_{t+1}

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

product rule

$$= \boxed{\frac{1}{p(x_{1:T})}} \frac{\alpha_t p(z_t | x_{1:t})}{\boxed{\frac{1}{p(x_{1:t})}}} \underbrace{p(z_{t+1} | z_t)}_{= A_{z_t, z_{t+1}}} \underbrace{p(x_{t+1} | z_{t+1})}_{\text{NormPDF}(x_{t+1} | \mu_{z_{t+1}}, \Sigma_{z_{t+1}})} p(x_{t+2:T} | z_{t+1})$$

conditional independence assumptions

- Const wrt z
- easy given θ
- using α, β defined on prev page

Thus, we can compute the adjacent timestep posterior if we have precomputed $\{\alpha_t, \beta_t\}_{t=1}^T$ as:

$$p(z_t=j, z_{t+1}=k | x_{1:T}) = \frac{\alpha_{tj} A_{jk} L_{t+1,k} \beta_{t+1,k}}{\sum_{l=1}^K \sum_{m=1}^K \alpha_{tl} A_{lm} L_{t+1,m} \beta_{t+1,m}} \quad \text{for } t = 1, 2, \dots, T-1$$

where $L_{tk} = p(x_t | z_t=k, \theta) = \text{NormPDF}(x_t | \mu_k, \Sigma_k)$ when emission model is Gaussian

Computing Forward Messages α via Dynamic Programming (FORWARD alg.)

Definition: $\alpha_{tk} = p(z_t = k | x_{1:t})$ for $t=1, 2, \dots, T$

Procedure: Dynamic Program, with base case $t=1$ and recurrence relation that computes α_{t+1} from α_t
 $\alpha_{t+1, :} = f(\alpha_{t, :})$

$$\alpha_{t+1, k} = p(z_{t+1} = k | x_{1:t+1}) = \sum_j p(z_t = j, z_{t+1} = k | x_{1:t+1})$$

multiply by $\frac{p(x_{t+1} | x_{1:t})}{p(x_{t+1} | x_{1:t})}$

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

product rule

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

cond. indep. cond. indep.

Recursive update:

$$\alpha_{t+1, k} = \frac{\sum_{j=1}^K \alpha_{tj} A_{jk} L_{t+1, k}}{\sum_l \sum_j \alpha_{tj} A_{jl} L_{t+1, l}}$$

alpha-update (α_t, A, L_{t+1})
 "Forward" update from α_t to α_{t+1}
 Input: α_t : length K vector, sums to 1
 A : $K \times K$ transition proba matrix
 L_{t+1} : length K vector "likelihood"
 $L_{t+1, k} = p(x_{t+1} | z_t = k)$

FORWARD algorithm

Base Case $t=1$: $\alpha_{1k} = \frac{\pi_k L_{1k}}{\sum_l \pi_l L_{1l}}$

for $t=2, 3, 4, \dots, T$:

$$\alpha_{tk} = \text{alpha_update}(\alpha_{t-1}, A, L_t)$$

return $\{\alpha_t\}_{t=1}^T$

Runtime
 linear in T
 quadratic in K

Using Forward Messages to compute incomplete
log likelihood.

Recall $\log p(x_{1:T} | \theta)$ is useful to know
difficult to compute
by naively summing out
 $z_{1:T}$

Studying the recurrence relation in alpha update, we focus
on the denominator and see

$$p(x_{t+1} | x_{1:t}) = \sum_{j=1}^K \sum_{k=1}^K \alpha_{tj} A_{jk} L_{t+1,k}$$
$$= \text{mat_mult}(\alpha_t, A) \cdot L_{t+1}$$

dot product
or inner product

Insight: Compute

$$\log p(x_{1:T}) = \underbrace{\log p(x_1)}_{\text{use denominator in base case (t=1) of FORWARD alg.}} + \sum_{t=2}^T \underbrace{\log p(x_t | x_{1:t-1})}_{\text{use denominator in step t of FORWARD alg.}}$$

Computing backward messages β_t via dynamic programming (BACKWARD alg.)

Definition: $\beta_{tk} = P(x_{t+1:T} | z_t = k)$, for $t = 1, 2, \dots, T$

special case: $\beta_{Tk} = 1 \forall k$
 because T is final timestep. x_{T+1} does not exist

Procedure: Dynamic Programming. Need to fill in table $\beta \in K \times K$

BACKWARD
 Init: $\beta_T = 1$
 for $T-1, T-2, \dots, 1$:
 $\beta_{tk} = \text{beta_update}(\beta_{t+1}, L_{t+1}, A)$ $\beta_{t-1} = f(\beta_t)$ for $t = T, T-1, \dots, 2, 1$
 Start with $\beta_T = 1$
 Work back in time using recurrence relation

Relation

$$\begin{aligned} \beta_{t-1,k} = P(x_{t:T} | z_{t-1} = k) &= \sum_{l=1}^K P(x_{t:T}, z_t = l | z_{t-1} = k) \quad \text{sum rule} \\ &= \sum_{l=1}^K P(x_t, x_{t+1:T} | z_t = l, z_{t-1} = k) P(z_t = l | z_{t-1} = k) \\ &= \sum_{l=1}^K P(x_{t+1:T} | z_t = l) P(x_t | z_t = l) P(z_t = l | z_{t-1} = k) \end{aligned}$$

x_{t+1} indep of prev x or $z_{1:t-1}$ given z_t x_t indep of all else given z_t

beta_update
 Input: β_t K -len vector
 L_t K -len vector of likelihoods
 A $K \times K$ trans, proba

$$\beta_{t-1,k} = \sum_{l=1}^K \beta_{tl} L_{tl} A_{kl}$$

✓
 A_{kl}

Recap: E Step

13

Goal: Update S parameters of $q(z|s)$ to optimal values

Procedure:

Input: π, A, μ, Σ (HMM params)

Step 1: Calc likelihood $L = \{L_{tk}\}$
 $L_{tk} = \text{NormPDF}(x_t | \mu_k, \Sigma_k^2)$

Step 2: Calc forward messages $\alpha = \{\alpha_{tk}\}$
 $\log p(x|\theta), \alpha = \text{FORWARD}(\pi, A, L)$ $O(TK^2)$

Can compute $\log p(x|\theta)$ in FORWARD easily

Calc backward messages $\beta = \{\beta_{tk}\}$
 $\beta = \text{BACKWARD}(A, L)$ $O(TK^2)$

Step 3: Update s_t to match adjacent timestep joint posterior for $t=1, 2, \dots, T-1$:

$$s_{tjk} = p(z_t=j, z_{t+1}=k | x_{1:T}) \quad \text{for } j, k \text{ in } 1 \dots K$$
$$\propto \alpha_{tj} A_{jk} L_{t+1,k} \beta_{t+1,k}$$

Step 4: Update $r = \{r_{tk}\}$ to match single timestep posterior for $t=1, 2, \dots, T$

$$r_{tk} \propto \alpha_{tk} \beta_{tk}$$

return $\log p(x_{1:T}|\theta), s, r$

Remember: $r = r(s)$
We can deterministically calculate r from s .
Just useful to have direct expression.