COMS30035, Machine learning: Sequential Data 3: EM for HMMs

Edwin Simpson edwin.simpson@bristol.ac.uk

Department of Computer Science, SCEEM University of Bristol

December 12, 2023

Agenda

- Markov Models
- Hidden Markov Models
- EM for HMMs
- Linear Dynamical Systems

Hidden Markov Models (HMMs)

We want to use maximum likelihood to estimate the HMM parameters:

- 1. A transition matrix
- 2. π initial state probabilities
- 3. ϕ parameters of the emission distributions
- We examine the unsupervised case where the sequence of states Z is not observed.

$$\ln p(\boldsymbol{X}|\boldsymbol{A}, \boldsymbol{\pi}, \boldsymbol{\phi}) = \\ \ln \sum_{\boldsymbol{Z}} \left\{ p(\boldsymbol{z}_1|\boldsymbol{\pi}) \prod_{n=2}^{N} p(\boldsymbol{z}_n|\boldsymbol{z}_{n-1}, \boldsymbol{A}) \prod_{n=1}^{N} p(\boldsymbol{x}_n|\boldsymbol{\phi}, \boldsymbol{z}_n) \right\}$$

Likelihood for an HMM

- As with GMMs, there is no closed-form solution to the MLE, so we turn to EM
- Unlike GMM, the likelihood doesn't factorise over the data points:

1.
$$\ln p(\mathbf{X}|\mathbf{A}, \pi, \phi) = \ln \sum_{\mathbf{Z}} \left\{ p(\mathbf{z}_1|\pi) \prod_{n=2}^{N} p(\mathbf{z}_n|\mathbf{z}_{n-1}, \mathbf{A}) \prod_{n=1}^{N} p(\mathbf{x}_n|\phi, \mathbf{z}_n) \right\}$$

- 2. The distribution of z_n depends on z_{n-1} , which also depends on z_{n-2} ...
- 3. Can't just sum over the values of z_n independently for each data point.
- 4. So we have to sum over all K^N possible sequences **Z**!

Expectation Maximisation (EM)

- Goal: maximise the expected log likelihood
- First, we define $Q(\theta|\theta^{old}) = \sum_{\boldsymbol{Z}} p(\boldsymbol{Z}|\boldsymbol{X}, \theta^{old}) \ln p(\boldsymbol{X}, \boldsymbol{Z}|\theta)$.
- 1. Initialise the parameters with a random guess: $\theta^{old} = \{ \mathbf{A}, \pi, \phi \}$.
- 2. **E-step**: use θ^{old} to compute expectations over **Z** required to compute $Q(\theta|\theta^{old})$.
- 3. **M-step**: choose the values of $\theta = \{A, \pi, \phi\}$ that maximise $Q(\theta|\theta^{old})$.
- 4. Set $\theta^{old} = \theta$.
- 5. Repeat steps 2-4 until convergence.

E step

- We need to compute expectations of the latent states and pairs of latent states:
- Responsibilities: $\gamma(z_{nk}) = p\left(z_n = k | \boldsymbol{X}, \boldsymbol{\theta}^{(old)}\right)$

State pairs:
$$\xi(z_{n-1,j}, z_{nk}) = p\left(z_{n-1} = j, z_n = k | \boldsymbol{X}, \boldsymbol{\theta}^{(old)}\right)$$

 To compute these efficiently, we need the *forward-backward* algorithm (coming up in a few slides...)

M step

- $\pi_k = \gamma(z_{1k})$ • $A_{ik} = \sum_{n=2}^{N} \xi(z_{n-1,i}, z_{nk}) / \sum_{n=2}^{N} \gamma(z_{n-1,i})$
- φ_k: parameters of posterior emission distributions, with observations weighted by responsibilities, γ(z_{nk})
 - If we have Gaussian emissions, the equations are the same as for GMM.
 - Discrete observations with value i:

$$\phi_{ki} = \rho(x_n = i | z_n = k) = \frac{\sum_{n=1}^{N} \gamma(z_{nk}) [x_n = i]}{\sum_{n=1}^{N} \gamma(z_{nk})}$$
(1)

Forward-backward Algorithm

(

A specific example of the *sum-product algorithm* used in the E-step
 Forward pass computes for each time-step *n* and state value *k*:

$$\begin{aligned}
\alpha(z_{nk}) &= p(\mathbf{x}_{1}, ..., \mathbf{x}_{n}, z_{n} = k | \pi, \mathbf{A}, \phi) \\
&= p(\mathbf{x}_{n} | z_{n} = k, \phi_{k}) \sum_{l=1}^{K} A_{lk} \alpha(z_{n-1,l})
\end{aligned} (2)$$

Backward pass computes:

$$\beta(z_{nk}) = p(\mathbf{x}_{n+1}, ..., \mathbf{x}_N | z_n = k, \mathbf{A}, \phi)$$

= $\sum_{l=1}^{K} A_{kl} p(\mathbf{x}_{n+1} | z_{n+1} = l, \phi_l) \beta(z_{n+1,l})$ (3)

Forward-backward Algorithm

• Use the computed α and β terms to compute our expectations over z:

$$\tilde{\xi}(z_{n-1,l}, z_{nk}) = p(\mathbf{x}_1, ..., \mathbf{x}_{n-1}, z_{n-1} = l | \mathbf{A}, \pi, \phi)$$
before

$$p(z_n = k | z_{n-1} = l, \mathbf{A}) p(\mathbf{x}_n | z_n = k, \phi)$$
current

$$p(\mathbf{x}_{n+1}, ..., x_N | z_n = k, \mathbf{A}, \phi)$$
after

$$= \alpha(z_{n-1,l}) A_{lk} p(\mathbf{x}_n | z_n = k, \phi) \beta(z_{nk})$$
(4)

•
$$\xi(z_{n-1,l}, z_{nk}) = \tilde{\xi}(z_{n-1,l}, z_{nk}) / \sum_{l=1}^{K} \sum_{k=1}^{K} \tilde{\xi}(z_{n-1,l}, z_{nk})$$

• $\gamma(z_{nk}) = \sum_{l=1}^{K} \xi(z_{n-1,l}, z_{nk})$

Putting It All Together...

- 1. Initialise the parameters with a random guess: $\theta^{old} = \{A, \pi, \phi\}$.
- **2. E-step** using θ^{old} :
 - 2.1 Run forward pass to compute $\alpha(z_{nk})$
 - 2.2 Run backward pass to compute $\beta(z_{nk})$
 - **2.3** Use $\alpha(z_{n-1,l})$ and $\beta(z_{nk})$ to compute $\xi(z_{n-1,l}, z_{nk})$ and $\gamma(z_{nk})$.
- 3. M-step using $\xi(z_{n-1,l}, z_{nk})$ and $\gamma(z_{nk})$, update $\theta = \{\pi, \mathbf{A}, \phi\}$.
- 4. Set $\theta^{old} = \theta$.
- 5. Repeat steps 2-4 until convergence.

By summing inside each forward and backward computations, we now have an algorithm that is linear $(\mathcal{O}(N))$ rather than exponential $(\mathcal{O}(K^N))$ in the sequence length \mathfrak{D} .

Viterbi Algorithm

- Given our estimated model parameters θ = {π, A, φ}, how can we predict a sequence of hidden states Z?
- Most probable labels (given by the values of \(\gamma(z_{nk})\)) are not the same as the most probable sequence!
- ► We apply a max-sum algorithm called viterbi to "decode" the sequence with O(N) computational cost.

Viterbi Algorithm

- Forward pass: compute the probability of the most likely sequence that leads to each possible state at time n.
- Backward pass: starting with the most likely final state and recursing backwards, choose the previous state n - 1 that makes the chosen state at n most likely.

Viterbi Algorithm

► Forward pass: 1. $\omega(z_{1k}) = \ln \pi_k + \ln p(\mathbf{x}_1 | z_1 = k)$ 2. For n = 2 to N compute for each state value k: 2.1 $\omega(z_{nk}) = \max_l \{ \omega(z_{n-1,l}) + \ln p(z_n = k | z_{n-1} = l) \} + \ln p(\mathbf{x}_n | z_n = k)$. 2.2 $\psi(z_{nk}) = \operatorname*{argmax}_l \{ \omega(z_{n-1,l}) + \ln p(z_n = k | z_{n-1} = l) \} + \ln p(\mathbf{x}_n | z_n = k)$. 2.3 Passes messages from the start of the sequence to the end.

Backward pass:

- 1. Most likely final state: $\hat{z}_N = \operatorname{argmax} \omega(z_{Nk})$.
- 2. For n = N 1 to 1: $\hat{z}_n = \psi(z_{n+1,\hat{z}_{n+1}})$.
- There are multiple paths leading to each possible state at each step n. We keep only the path with the highest probability, so we don't have to compute the likelihood of every complete path from 1 to N.

Summary

By computing sums and maximums at each timestep we can perform inference over an exponential number of sequences. We use the...

- Forward-backward algorithm, an instance of the more general sum-product algorithm to marginalise over sequences of hidden states.
- Viterbi algorithm, an instance of the more general max-sum algorithm to find the most likely sequence of hidden states.

Please do the quiz for this lecture on Blackboard. Next up: linear dynamical systems for modelling continuous states.