

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

Edwin Simpson

`edwin.simpson@bristol.ac.uk`

Department of Computer Science, SCEEM
University of Bristol

November 5, 2020

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. \mathbf{A} - transition matrix
 2. π - initial state probabilities
 3. ϕ - parameters of the emission distributions
- ▶ We assume the *unsupervised* case where the sequence of states \mathbf{Z} is not observed.
- ▶ $\ln p(\mathbf{X}|\mathbf{A}, \pi, \phi) = \ln \sum_{\mathbf{Z}} \left\{ p(\mathbf{Z}|\mathbf{A}, \pi) \prod_{n=1}^N p(\mathbf{x}_n|\phi, \mathbf{z}_n) \right\}$

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 \mathbf{z}_n depends on \mathbf{z}_{n-1} , which also depends on $\mathbf{z}_{n-2} \dots$
 3. Can't just sum over the values of \mathbf{z}_n independently for each data point.
 4. So we have to sum over all K^N possible sequences \mathbf{Z} !

Expectation Maximisation (EM)

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

E step

- ▶ Responsibilities: $\gamma(z_{nk}) = \sum_{\mathbf{Z}} p(\mathbf{Z}) z_{nk}$
- ▶ State pairs: $\xi(z_{n-1,j}, z_{nk}) = \sum_{\mathbf{Z}} p(\mathbf{Z}) z_{n-1,j} z_{nk}$
- ▶ To compute these efficiently, we need the *forward-backward* algorithm.

M step

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

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

Forward-backward Algorithm

- ▶ A specific example of the *sum-product algorithm*
- ▶ We want to compute:
 - ▶ Responsibilities: $\gamma(z_{nk}) = p(z_{nk} | \mathbf{A}, \pi, \phi, \mathbf{x}) = \sum_{l=1}^K \xi(z_{n-1,l}, z_{nk})$
 - ▶ State pairs: $\xi(z_{n-1,l}, z_{nk}) = p(z_{n-1,l}, z_{nk} | \mathbf{A}, \pi, \phi, \mathbf{x})$
- ▶ Forward pass computes:

$$\mathbb{E}_{z_1, \dots, z_{n-2}} [p(z_{n-1,l} = 1, z_{nk} = 1 | \mathbf{z}_1, \dots, \mathbf{z}_{n-2}, \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{A}, \pi, \phi)] \quad (2)$$

- ▶ Backward pass computes:

$$\mathbb{E}_{z_1, \dots, z_{n-2}, z_{n+1}, \dots, z_N} [p(z_{n-1,l} = 1, z_{nk} = 1 | \mathbf{z}_1, \dots, \mathbf{z}_{n-2}, \mathbf{z}_{n+1}, \dots, \mathbf{z}_N, \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{A}, \pi, \phi)] \quad (3)$$

Forward-backward Algorithm

- ▶ Let's take the final output of the backward pass and break it down:
 $\mathbb{E}_{z_1, \dots, z_{n-2}, z_{n+1}, \dots, z_N} [p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{z}_1, \dots, \mathbf{z}_{n-2}, \mathbf{z}_{n+1}, \dots, \mathbf{z}_N, \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{A}, \pi, \phi)]$
- ▶ We can write it as separate terms for time-steps *before* n , the *current* time-step n , and time-steps *after* n .

$$\begin{aligned} \tilde{\xi}(z_{n-1,l}, z_{nk}) &= p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, z_{n-1,l} = 1 | \mathbf{A}, \pi, \phi) && \text{before} \\ &\quad p(z_{nk} = 1 | z_{n-1,l} = 1, \mathbf{A}) p(\mathbf{x}_n | z_{nk} = 1, \phi) && \text{current} \\ &\quad p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_{nk} = 1, \mathbf{A}, \phi) && \text{after} \\ &= \alpha(z_{n-1,l}) A_{kl} p(\mathbf{x}_n | z_{nk} = 1, \phi) \beta(z_{nk}) && (4) \end{aligned}$$

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

Forward Pass

Compute and save the α terms for all states at all time-steps¹.

$$\begin{aligned}\alpha(z_{nk}) &= p(\mathbf{x}_1, \dots, \mathbf{x}_n, z_{nk} = 1 | \pi, \mathbf{A}, \phi) \\ &= p(\mathbf{x}_n | z_{nk} = 1, \phi_k) \sum_{l=1}^K A_{lk} \alpha(z_{n-1, l})\end{aligned}\tag{5}$$

1. Initialise $\alpha(z_{1k}) = \pi_k p(\mathbf{x}_1 | z_{1k} = 1, \phi_k)$ for all states k .
2. Compute $\alpha(z_{nk})$ for each n from 2 to N .
3. To avoid $\alpha(z_{nk})$ values becoming increasingly small, normalise over $\sum_{k=1}^K \alpha(z_{nk})$ at each iteration.

¹See Bishop (2006) Section 13.2.2 for full derivation of the algorithm.

Backward Pass

Compute and save the β terms for all states and all time-steps².

$$\begin{aligned}\beta(z_{nk}) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | z_{nk} = 1, \mathbf{A}, \phi) \\ &= \sum_{l=1}^K A_{kl} p(\mathbf{x}_{n+1} | z_{n+1,l} = 1, \phi_l) \beta(z_{n+1,l})\end{aligned}\tag{6}$$

1. Initialise $\beta(z_{Nk}) = 1$ for all states k .
2. Compute $\beta(z_{nk})$ for each n from $N - 1$ to 1.
3. To avoid $\beta(z_{nk})$ values becoming too small after many iterations, normalise over $\sum_{k=1}^K \alpha(z_{n+1,k})$ at each iteration..

²See Bishop (2006) Section 13.2.2 for full derivation of the algorithm.

Putting It All Together...

1. Initialise the parameters with a random guess: $\theta^{old} = \{\mathbf{A}, \pi, \phi\}$.
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 😊.

Viterbi Algorithm

- ▶ Given our estimated model parameters $\theta = \{\pi, \mathbf{A}, \phi\}$, how can we predict a sequence of hidden states \mathbf{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 $\mathcal{O}(N)$ 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. $\omega(\mathbf{z}_1) = \ln \pi + \ln p(\mathbf{x}_1|\mathbf{z}_1)$
2. For $n = 2$ to N compute:
 - 2.1 $\omega(\mathbf{z}_n) = \max_{\mathbf{z}_{n-1}} \{\omega(\mathbf{z}_{n-1}) + \ln p(\mathbf{z}_n|\mathbf{z}_{n-1})\} + \ln p(\mathbf{x}_n|\mathbf{z}_n).$
 - 2.2 $\psi(\mathbf{z}_n) = \operatorname{argmax}_{\mathbf{z}_{n-1}} \{\omega(\mathbf{z}_{n-1}) + \ln p(\mathbf{z}_n|\mathbf{z}_{n-1})\} + \ln p(\mathbf{x}_n|\mathbf{z}_n).$
 - 2.3 Passes messages from the start of the sequence to the end.

► Backward pass:

1. Most likely final state: $\hat{\mathbf{z}}_N = \operatorname{argmax}_{k \in \{1, \dots, K\}} \omega(\mathbf{z}_N)_k.$
 2. For $n = N - 1$ to 1 : $\hat{\mathbf{z}}_n = \psi(\mathbf{z}_{n+1})_{\hat{\mathbf{z}}_{n+1}}.$
- 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.

Now do the quiz!

Please do the quiz for this lecture on Blackboard.

Next up: linear dynamical systems for modelling continuous states.