Introduction to Hidden Markov Models

Peter Bell

Automatic Speech Recognition— ASR Lecture 4 26 January 2023

HMMs

- Introduction to HMMs models
- HMMs for ASR
- Likelihood computation with the forward algorithm

Fundamental Equation of Statistical Speech Recognition

If X is the sequence of acoustic feature vectors (observations) and W denotes a word sequence, the most likely word sequence W^* is given by

$$\mathsf{W}^* = \arg \max_{\mathsf{W}} \mathsf{P}(\mathsf{W} \mid \mathsf{X})$$

Fundamental Equation of Statistical Speech Recognition

If X is the sequence of acoustic feature vectors (observations) and W denotes a word sequence, the most likely word sequence W^* is given by

$$\mathsf{W}^* = \arg \max_{\mathsf{W}} \mathsf{P}(\mathsf{W} \mid \mathsf{X})$$

Applying Bayes' Theorem:

$$P(W \mid X) = \frac{p(X \mid W)P(W)}{p(X)}$$

$$\propto p(X \mid W)P(W)$$

$$W^* = \arg \max_{W} \underbrace{p(X \mid W)}_{\text{Acoustic}} \underbrace{P(W)}_{\text{Language}}$$

$$model \qquad model$$



Hierarchical modelling of speech



- A statistical model for time series data with a set of **discrete** states $\{1, \ldots, J\}$ (we index them by *j* or *k*)
- At each time step *t*:
 - the model is in a fixed state q_t .
 - the model generates an observation, x_t, according to a probability distribution that is specific to the state
- We don't actually observe which state the model is in at each time step hence **"hidden"**.
- Observations can be either continous or discrete (usually the former)

HMM probabilities



- Imagine we know the state at a given time step t, $q_t = k$
- Then the probability of being in a new state, *j* at the next time step, is dependent only on *q*_t. This is the **Markov** assumption.
- Alternatively: q_{t+1} is conditionally independent of q_1, \ldots, q_{t-1} , given q_t .

HMM assumptions



Observation x_t is *conditionally independent* of other observations, given the state that generated it, q_t

The parameters of the model, λ , are given by:

- Transition probabilities $a_{kj} = P(q_{t+1} = j | q_t = k)$
- Observation probabilities $b_j(x) = P(x|q = j)$

• The HMM topology determines the set of allowed transitions between states

- The HMM topology determines the set of allowed transitions between states
- In principle any topology is possible

- The HMM topology determines the set of allowed transitions between states
- In principle any topology is possible



- The HMM topology determines the set of allowed transitions between states
- In principle any topology is possible



Not all transition probabilities are shown

Example topologies



Traditional speech recognition: Speaker recognition: left-to-right HMM with 3 \sim 5 states ergodic HMM

We generally model words or phones with a left-to-right topology with self loops.



Traditional HMMs for ASR tend to model each phone with three distinct states (this also enforces a minimum phone duration of three frames of observations)



Traditional HMMs for ASR tend to model each phone with three distinct states (this also enforces a minimum phone duration of three frames of observations)



HMMs for ASR

Traditional HMMs for ASR tend to model each phone with three distinct states (this also enforces a minimum phone duration of three frames of observations)



The phone model topologies can be concatenated to form a HMM for the whole word

HMMs for ASR

Traditional HMMs for ASR tend to model each phone with three distinct states (this also enforces a minimum phone duration of three frames of observations)



This model naturally generates an alignment between states and observations (and hence words/phones).

Suppose we have a sequence of observations of length T, $X = (x_1, \ldots, x_T)$, and Q is a known state sequence, (q_1, \ldots, q_T) . Then we can use the HMM to compute the joint likelihood of Xand Q:

$$P(X, Q; \lambda) = P(q_1)P(x_1|q_1)P(q_2|q_1)P(x_2|q_2)\dots$$
(1)
= $P(q_1)P(x_1|q_1)\prod_{t=2}^{T}P(q_t|q_{t-1})P(x_t|q_t)$ (2)

 $P(q_1)$ denotes the initial occupancy probability of each state

Working with HMMs requires the solution of three problems:

Likelihood Determine the overall likelihood of an observation sequence X = (x₁,...,x_t,...,x_T) being generated by a known HMM topology, *M*.

Working with HMMs requires the solution of three problems:

- Likelihood Determine the overall likelihood of an observation sequence X = (x₁,...,x_t,...,x_T) being generated by a known HMM topology, *M*.
- Oecoding and alignment Given an observation sequence and an HMM, determine the most probable hidden state sequence

Working with HMMs requires the solution of three problems:

- Likelihood Determine the overall likelihood of an observation sequence X = (x₁,...,x_t,...,x_T) being generated by a known HMM topology, *M*.
- Oecoding and alignment Given an observation sequence and an HMM, determine the most probable hidden state sequence
- Training Given an observation sequence and an HMM, find the state occupation probabilities

- Likelihood Determine the overall likelihood of an observation sequence X = (x₁,...,x_t,...,x_T) being generated by a known HMM topology, *M*.
 - \rightarrow the forward algorithm
- NB. We do **not** know the state sequence!

By talking about HMM topologies in the context of speech recognition, \mathcal{M} , we can mean:

- A restricted left-to-right topology based on a known word/sentence, leading to a "trellis-like" structure over time
- A much less restricted topology based on a grammar or language model – or something in between
- Some algorithms are not (generally) suitable for unrestricted topologies

Example: trellis for a 3-state left-to-right phone HMM



• Goal: determine p(X | M)

- Goal: determine p(X|M)
- Sum over all possible state sequences Q = (q₁,...,q_T) that could result in the observation sequence X

$$\begin{split} p(\mathsf{X}|\mathcal{M}) &= \sum_{Q \in \mathcal{Q}} P(\mathsf{X}, Q|\mathcal{M}) \\ &= \sum_{Q \in \mathcal{Q}} P(q_1) P(\mathsf{x}_1|q_1) \prod_{t=2}^T P(q_t|q_{t-1}) P(\mathsf{x}_t|q_t) \end{split}$$

- Goal: determine p(X | M)
- Sum over all possible state sequences Q = (q₁,..., q_T) that could result in the observation sequence X

$$\begin{split} p(\mathsf{X}|\mathcal{M}) &= \sum_{Q \in \mathcal{Q}} P(\mathsf{X}, Q|\mathcal{M}) \\ &= \sum_{Q \in \mathcal{Q}} P(q_1) P(\mathsf{x}_1|q_1) \prod_{t=2}^T P(q_t|q_{t-1}) P(\mathsf{x}_t|q_t) \end{split}$$

• How many paths Q do we have to calculate?

$$\sim \underbrace{N \times N \times \cdots N}_{T \text{ times}} = N^T \qquad N: \text{ number of HMM states} \\ T: \text{ length of observation}$$

e.g. $N^{T} \approx 10^{10}$ for $N \!=\! 3, T \!=\! 20$

- Goal: determine p(X | M)
- Sum over all possible state sequences Q = (q₁,..., q_T) that could result in the observation sequence X

$$\begin{split} p(\mathsf{X}|\mathcal{M}) &= \sum_{Q \in \mathcal{Q}} P(\mathsf{X}, Q|\mathcal{M}) \\ &= \sum_{Q \in \mathcal{Q}} P(q_1) P(\mathsf{x}_1|q_1) \prod_{t=2}^T P(q_t|q_{t-1}) P(\mathsf{x}_t|q_t) \end{split}$$

• How many paths Q do we have to calculate?

$$\sim \underbrace{N \times N \times \cdots N}_{T \text{ times}} = N^T \qquad N: \text{ number of HMM states} \\ T: \text{ length of observation}$$

e.g. $N^{T} pprox 10^{10}$ for $N\!=\!3, T\!=\!20$

• Computation complexity of multiplication: $O(2TN^{T})$

Likelihood: The Forward algorithm

The Forward algorithm:

- Rather than enumerating each sequence, compute the probabilities recursively (exploiting the Markov assumption)
- Reduces the computational complexity to $O(TN^2)$
- State time trellis for an arbitrary HMM topology



The forward probability

Define the *Forward probability*, $\alpha_j(t)$: the probability of observing the observation sequence $x_1 \dots x_t$ and being in state j at time t:

$$\alpha_j(t) = p(\mathsf{x}_1, \ldots, \mathsf{x}_t, q_t = j | \mathcal{M})$$

We can recursively compute this probability

Initial and final state probabilities

It what follows it is convenient to define:

• an additional single initial state $S_I = 0$, with transition probabilities

$$a_{0j} = P(q_1 = j)$$

denoting the probability of starting in state j

- a single final state, *S_E*, with transition probabilities *a_{jE}* denoting the probability of the model terminating in state *j*.
- S_I and S_E are both non-emitting

Likelihood: The Forward recursion

Initialisation

$$\begin{array}{ll} \alpha_j(0) = 1 & j = 0 \\ \alpha_j(0) = 0 & j \neq 0 \end{array}$$

Recursion

$$\alpha_j(t) = \sum_{i=0}^J \alpha_i(t-1)a_{ij}b_j(\mathbf{x}_t) \qquad 1 \le j \le J, \ 1 \le t \le T$$

Termination

$$p(\mathsf{X}|\mathcal{M}) = \alpha_E = \sum_{i=1}^J \alpha_i(T) \mathbf{a}_{iE}$$

 s_I : initial state, s_E : final state

Likelihood: Forward Recursion

$$\alpha_j(t) = p(\mathsf{x}_1, \ldots, \mathsf{x}_t, q_t = j | \mathcal{M}) = \sum_{i=1}^J \alpha_i(t-1) a_{ij} b_j(\mathsf{x}_t)$$



More HMM algorithms

- Finding the most likely path with the Viterbi algorithm
- Parameter estimation:
 - the Forward-Backward algorithm
 - the Expectation-Maximisation algorithm

- Gales and Young (2007). "The Application of Hidden Markov Models in Speech Recognition", *Foundations and Trends in Signal Processing*, 1 (3), 195–304: section 2.2.
- Jurafsky and Martin (2008). Speech and Language Processing (2nd ed.): sections 6.1-6.5; 9.2; 9.4. (Errata at http://www.cs.colorado.edu/~martin/SLP/Errata/SLP2-PIEV-Errata.html)
- Rabiner and Juang (1989). "An introduction to hidden Markov models", *IEEE ASSP Magazine*, 3 (1), 4–16.
- Renals and Hain (2010). "Speech Recognition", *Computational Linguistics and Natural Language Processing Handbook*, Clark, Fox and Lappin (eds.), Blackwells.