

School of Computing and Information Systems
The University of Melbourne
COMP30027 Machine Learning (Semester 1, 2021)

Sample solutions: Week 10

1. Hidden Markov Models (HMMs) are best used when the observables are a **univariate time series**: we are just observing a single variable, which changes over time due to some factor that can be estimated from previous observations.

- (a) Recall the two main assumptions (Markov, output independence) that are built into an HMM.

Markov assumption: the likelihood of transitioning into a given state depends only on the current state, and not the previous state(s) (or output(s)), i.e. $P(q_t | q_1 \dots q_{t-1}) \approx P(q_t | q_{t-1})$

Output independence assumption: the likelihood of a state producing a certain observation (as output) does not depend on the preceding (or following) state(s) (or output(s)), i.e. $P(o_t | q_1 \dots q_t, o_1 \dots o_{t-1}) \approx P(o_t | q_t)$

- (b) Could we construct the HMM in such a way to relax these assumptions? What would the model look like, and what is the major downside?

Well, we could have pairs of states in the conditions for our transition probability matrix A , and pairs of states in the conditions for our output probability matrix B , but this will vastly increase the number of parameters in the model.

- (c) Could we build an HMM for a **multivariate time series**, where we have a number of observed variables for a given (hidden) state?

We could represent the outputs as a tuple, again at the cost of vastly increasing the number of parameters. However, sometimes coupling the outputs like this is unnecessary; it might be possible to just generate independent HMMs for each output.

2. **Natural language processing** is one common application for HMMs: we have a single observation (a “word”) that varies over time (a “sentence” or “document”), where each observation is associated with some property (like “part of speech”).

Consider the following HMM: $\Pi[J, N, V] = [0.3, 0.4, 0.3]$

| A | J (adj) | N (noun) | V (verb) |
|---|---------|----------|----------|
| J | 0.4 | 0.5 | 0.1 |
| N | 0.1 | 0.4 | 0.5 |
| V | 0.4 | 0.5 | 0.1 |

| B | <i>brown</i> | <i>leaves</i> | <i>turn</i> |
|---|--------------|---------------|-------------|
| J | 0.8 | 0.1 | 0.1 |
| N | 0.3 | 0.4 | 0.3 |
| V | 0.1 | 0.3 | 0.6 |

- (a) How might we go about obtaining the values in the matrices Π , A , and B given above, in a **supervised** context?

Each element a_{ij} of A is the count of how many times the state sequence i, j was observed in the labelled data, out of the number of times the state i was observed. In this case, this is a pair of part-of-speech tags (adj, verb, or noun).

Each element b_{ik} of B is the count of how many times the observation k was observed labelled with state i in the training data, out of the number of times the state i was observed. In this case, this is a word being labelled as the equivalent part-of-speech.

Each element π_i is the count of how many times state i was the start of an observation sequence, out of the number of observation sequences. In this case, this is where some part-of-speech starts a sentence.

(b) Use the **forward** algorithm to find the probability of the “sentence” brown leaves turn.

| α | | 1:brown | 2:leaves | 3:turn |
|----------|---|---|----------|--------|
| J: | J | $\pi[J]B[J, \text{brown}]$ $0.3 \times 0.8 = 0.24$ | | |
| N: | N | $\pi[N]B[N, \text{brown}]$ $0.4 \times 0.3 = 0.12$ | | |
| V: | V | $\pi[V]B[V, \text{brown}]$ $0.3 \times 0.1 = 0.03$ | | |

| α | 1:brown | 2:leaves | 3:turn |
|----------|---------|--|--|
| J: | 0.24 | $J \rightarrow J : \alpha_1(J) = 0.24$ $N \rightarrow J : \alpha_1(N) = 0.12$ $V \rightarrow J : \alpha_1(V) = 0.03$ $\sum_i (\alpha_1(i) a_{iJ}) b_{Jl} =$ | $\alpha_1(J) \times A[J, J] = 0.24 \times 0.4 = 0.096$ $\alpha_1(N) \times A[N, J] = 0.12 \times 0.1 = 0.012$ $\alpha_1(V) \times A[V, J] = 0.03 \times 0.4 = 0.012$ $(0.096 + 0.012 + 0.012) \times 0.1 = 0.012$ |
| N: | 0.12 | $J \rightarrow N : \alpha_1(J) = 0.24$ $N \rightarrow N : \alpha_1(N) = 0.12$ $V \rightarrow N : \alpha_1(V) = 0.03$ $\sum_i (\alpha_1(i) a_{iN}) b_{Nl} =$ | $\alpha_1(J) \times A[J, N] = 0.24 \times 0.5 = 0.12$ $\alpha_1(N) \times A[N, N] = 0.12 \times 0.4 = 0.048$ $\alpha_1(V) \times A[V, N] = 0.03 \times 0.5 = 0.015$ $(0.12 + 0.048 + 0.015) \times 0.4 = 0.0732$ |
| V: | 0.03 | $J \rightarrow V : \alpha_1(J) = 0.24$ $N \rightarrow V : \alpha_1(N) = 0.12$ $V \rightarrow V : \alpha_1(V) = 0.03$ $\sum_i (\alpha_1(i) a_{iV}) b_{Vl} =$ | $\alpha_1(J) \times A[J, V] = 0.24 \times 0.1 = 0.024$ $\alpha_1(N) \times A[N, V] = 0.12 \times 0.5 = 0.06$ $\alpha_1(V) \times A[V, V] = 0.03 \times 0.1 = 0.003$ $(0.024 + 0.06 + 0.003) \times 0.3 = 0.0261$ |

| α | 1:brown | 2:leaves | 3:turn | |
|----------|---------|----------|---|--|
| J: | 0.24 | 0.012 | $J \rightarrow J : \alpha_2(J) = 0.012$ $N \rightarrow J : \alpha_2(N) = 0.0732$ $V \rightarrow J : \alpha_2(V) = 0.0261$ $\sum_i (\alpha_2(i) a_{iJ}) b_{Jt} =$ | $\alpha_2(J) \times A[J, J] = 0.012 \times 0.4 = 0.0048$ $\alpha_2(N) \times A[N, J] = 0.0732 \times 0.1 = 0.00732$ $\alpha_2(V) \times A[V, J] = 0.0261 \times 0.4 = 0.01044$ $(0.0048+0.00732+0.01044) \times 0.1=0.002256$ |
| N: | 0.12 | 0.0732 | $J \rightarrow N : \alpha_2(J) = 0.012$ $N \rightarrow N : \alpha_2(N) = 0.0732$ $V \rightarrow N : \alpha_2(V) = 0.0261$ $\sum_i (\alpha_2(i) a_{iN}) b_{Nl} =$ | $\alpha_2(J) \times A[J, N] = 0.012 \times 0.5 = 0.006$ $\alpha_2(N) \times A[N, N] = 0.0732 \times 0.4 = 0.02928$ $\alpha_2(V) \times A[V, N] = 0.0261 \times 0.5 = 0.01305$ $(0.006+0.02928+0.01305) \times 0.3=0.014499$ |
| V: | 0.03 | 0.0261 | $J \rightarrow V : \alpha_2(J) = 0.012$ $N \rightarrow V : \alpha_2(N) = 0.0732$ $V \rightarrow V : \alpha_2(V) = 0.0261$ $\sum_i (\alpha_2(i) a_{iV}) b_{Vl} =$ | $\alpha_2(J) \times A[J, V] = 0.012 \times 0.1 = 0.0012$ $\alpha_2(N) \times A[N, V] = 0.0732 \times 0.5 = 0.0366$ $\alpha_2(V) \times A[V, V] = 0.0261 \times 0.1 = 0.00261$ $(0.0012+0.0366+0.00261) \times 0.6=0.024246$ |

The overall probability can be obtained by summing the values in the final column:

$$0.002256 + 0.014499 + 0.024246 = 0.041001$$

(c) Use the **Viterbi** algorithm to find the most likely state sequence for the sentence **brown leaves turn**.

| α | | 1: <i>brown</i> | 2: <i>leaves</i> | 3: <i>turn</i> |
|----------|---|---|------------------|----------------|
| J: | J | $\pi[J]B[J, \text{brown}]$ $0.3 \times 0.8 = 0.24$ | | |
| N: | N | $\pi[N]B[N, \text{brown}]$ $0.4 \times 0.3 = 0.12$ | | |
| V: | V | $\pi[V]B[V, \text{brown}]$ $0.3 \times 0.1 = 0.03$ | | |

| α | 1: <i>brown</i> | 2: <i>leaves</i> | 3: <i>turn</i> |
|----------|-----------------|---|----------------|
| J: | 0.24 | $J \rightarrow J: \alpha_1(J) = 0.24$ $\alpha_1(J) A[J,J]B[J, \text{leaves}] = 0.24 \times 0.4 \times 0.1 = \mathbf{0.0096}$ $N \rightarrow J: \alpha_1(N) = 0.12$ $\alpha_1(N) A[N,J]B[J, \text{leaves}] = 0.12 \times 0.1 \times 0.1 = 0.0012$ $V \rightarrow J: \alpha_1(V) = 0.03$ $\alpha_1(V) A[V,J]B[J, \text{leaves}] = 0.03 \times 0.4 \times 0.1 = 0.0012$ $\psi_2(J) \rightarrow J$ | |
| N: | 0.12 | $J \rightarrow N: \alpha_1(J) = 0.24$ $\alpha_1(J) A[J,N]B[N, \text{leaves}] = 0.24 \times 0.5 \times 0.4 = \mathbf{0.048}$ $N \rightarrow N: \alpha_1(N) = 0.12$ $\alpha_1(N) A[N,N]B[N, \text{leaves}] = 0.12 \times 0.4 \times 0.4 = 0.0192$ $V \rightarrow N: \alpha_1(V) = 0.03$ $\alpha_1(V) A[V,N]B[N, \text{leaves}] = 0.03 \times 0.5 \times 0.4 = 0.006$ $\psi_2(N) \rightarrow J$ | |
| V: | 0.03 | $J \rightarrow V: \alpha_1(J) = 0.24$ $\alpha_1(J) A[J,V]B[V, \text{leaves}] = 0.24 \times 0.1 \times 0.3 = 0.0072$ $N \rightarrow V: \alpha_1(N) = 0.12$ $\alpha_1(N) A[N,V]B[V, \text{leaves}] = 0.12 \times 0.5 \times 0.3 = \mathbf{0.018}$ $V \rightarrow V: \alpha_1(V) = 0.03$ $\alpha_1(V) A[V,V]B[V, \text{leaves}] = 0.03 \times 0.1 \times 0.3 = 0.0009$ $\psi_2(V) \rightarrow N$ | |

| α | 1: <i>brown</i> | 2: <i>leaves</i> | 3: <i>turn</i> |
|----------|-----------------|-----------------------------|---|
| J: | 0.24 | 0.0096 $J \rightarrow J$ | $J \rightarrow J: \alpha_2(J) = 0.0096$ $\alpha_2(J) A[J,J]B[J, \text{turn}] = 0.0096 \times 0.4 \times 0.1 = 0.000384$ $N \rightarrow J: \alpha_2(N) = 0.048$ $\alpha_2(N) A[N,J]B[J, \text{turn}] = 0.048 \times 0.1 \times 0.1 = 0.00048$ $V \rightarrow J: \alpha_2(V) = 0.018$ $\alpha_2(V) A[V,J]B[J, \text{turn}] = 0.018 \times 0.4 \times 0.1 = \mathbf{0.00072}$ $\psi_3(J) \rightarrow V$ |
| N: | 0.12 | 0.048 $J \rightarrow N$ | $J \rightarrow N: \alpha_2(J) = 0.0096$ $\alpha_2(J) A[J,N]B[N, \text{turn}] = 0.0096 \times 0.5 \times 0.3 = 0.00144$ $N \rightarrow N: \alpha_2(N) = 0.048$ $\alpha_2(N) A[N,N]B[N, \text{turn}] = 0.048 \times 0.4 \times 0.3 = \mathbf{0.00576}$ $V \rightarrow N: \alpha_2(V) = 0.018$ $\alpha_2(V) A[V,N]B[N, \text{turn}] = 0.018 \times 0.5 \times 0.3 = 0.0027$ $\psi_3(N) \rightarrow N$ |
| V: | 0.03 | 0.018 $N \rightarrow V$ | $J \rightarrow V: \alpha_2(J) = 0.0096$ $\alpha_2(J) A[J,V]B[V, \text{turn}] = 0.0096 \times 0.1 \times 0.6 = 0.000576$ $N \rightarrow V: \alpha_2(N) = 0.048$ $\alpha_2(N) A[N,V]B[V, \text{turn}] = 0.048 \times 0.5 \times 0.6 = \mathbf{0.0144}$ $V \rightarrow V: \alpha_2(V) = 0.018$ $\alpha_2(V) A[V,V]B[V, \text{turn}] = 0.018 \times 0.1 \times 0.6 = 0.00108$ $\psi_3(V) \rightarrow N$ |

The most likely tag sequence can be read right-to-left, based upon the maximum probability we've observed: in this case, 0.0144 when *turn* is a V; this value is derived from the $N \rightarrow V$ transition, so we can infer that *leaves* is an N; that in turn comes from the $J \rightarrow N$ transition, so *brown* is a J.