6.S064 INTRODUCTION TO MACHINE LEARNING Phase 2: Hidden Markov Models (lectures 15 and 16)

Motivation

In many practical problems, we would like to model pairs of sequences. Consider, for instance, the task of part-of-speech (POS) tagging. Given a sentence, we would like to compute the corresponding tag sequence:

Input: "Faith is a fine invention"

Output: "Faith/N is/V a/D fine/A invention/N"

More generally, a sequence labeling problem involves mapping a sequence of observations x_1, x_2, \ldots, x_n into a sequence of tags y_1, y_2, \ldots, y_n . In the example above, every word x is tagged by a single label y. One possible approach for solving this problem would be to label each word independently. For instance, a classifier could predict a part-of-speech tag based on the word, its suffix, its position in the sentence, etc. In other words, we could construct a feature vector on the basis of the observed "context" for the tag, and use the feature vector in a linear classifier. However, tags in a sequence are dependent on each other and this classifier would make each tagging decision independently of other tags. We would like our model to directly incorporate these dependencies. For instance, in our example sentence, the word "fine" can be either noun (N), verb (V) or adjective (A). The label V is not suitable since a tag sequence "D V" is very unlikely. Today, we will look at a model – a Hidden Markov Model – that allows us to capture some of these correlations.

Generative Tagging Model

Assume a finite set of words Σ and a finite set of tags \mathcal{T} . Define S to be the set of all sequence tag pairs $(x_1, \ldots, x_n, y_1, \ldots, y_n)$, $x_i \in \Sigma$ and $y_i \in \mathcal{T}$ for $i = 1 \ldots n$. S here contains sequences of different lengths as well, i.e., n varies as well. A generative tagging model is a probability distribution p over pairs of sequences:

- $p(x_1,\ldots,x_n,y_1,\ldots,y_n) \ge 0 \ \forall (x_1,\ldots,x_n,y_1,\ldots,y_n) \in S$
- $\sum_{(x_1,\dots,x_n,y_1,\dots,y_n)\in S} p(x_1,\dots,x_n,y_1,\dots,y_n) = 1$

If we have such a distribution, then we can use it to predict the most likely sequence of tags y_1, \ldots, y_n for any observed sequence of words x_1, \ldots, x_n , as follows

$$f(x_1,\ldots,x_n) = \operatorname*{argmax}_{y_1,\ldots,y_n} p(x_1,\ldots,x_n,y_1,\ldots,y_n)$$
(1)

where we view f as a mapping from word sequences to tags.

Three key questions:

- How to specify $p(x_1, \ldots, x_n, y_1, \ldots, y_n)$ with a few number of parameters (degrees of freedom)
- How to estimate the parameters in this model based on observed sequences of words (and tags).
- How to predict, i.e., how to find the most likely sequence of tags for any observed sequence of words: evaluate $\operatorname{argmax}_{y_1,\ldots,y_n} p(x_1,\ldots,x_n,y_1,\ldots,y_n)$

1 Model definition

Let X_1, \ldots, X_n and Y_1, \ldots, Y_n be sequences of random variables of length n. We wish to specify a joint probability distribution

$$P(X_1 = x_1, \dots, X_n = x_n, Y_1 = y_1, \dots, Y_n = y_n)$$
(2)

where $x_i \in \Sigma$, $y_i \in \mathcal{T}$. For brevity, we will write it as $p(x_1, \ldots, x_n, y_1, \ldots, y_n)$, i.e., treat it as a function of values of the random variables without explicating the variables themselves. We will define one additional random variable Y_{n+1} , which always takes the value STOP. Since our model is over variable length sequences, we will use the end symbol to model when to stop. In other words, if we observe x_1, \ldots, x_n , then clearly the symbol after y_1, \ldots, y_n , i.e., y_{n+1} , had to be STOP (otherwise we would have continued generating more symbols).

Now, let's start by rewriting the distribution a bit according to general rules that apply to any distribution. The goal is to put the distribution in a form where we can easily explicate our assumptions. First,

$$p(x_1, \dots, x_n, y_1, \dots, y_{n+1}) = p(y_1, \dots, y_{n+1})p(x_1, \dots, x_n | y_1, \dots, y_{n+1})$$
 (chain rule)

Then we will use the chain rule repeatedly along the sequence of tags

$$p(y_1, \dots, y_{n+1}) = p(y_1)p(y_2|y_1)p(y_3|y_1, y_2)\dots p(y_{n+1}|y_1, \dots, y_n)$$
 (chain rule)

So far, we have made no assumptions about the distribution at all. Since we don't expect tags to have very long dependences along the sequence, we will simply say that the next tag only depends on the current tag. In other words, we will "drop" the dependence on tags further back

$$p(y_1, \dots, y_{n+1}) \approx p(y_1)p(y_2|y_1)p(y_3|y_2)\dots p(y_{n+1}|y_n) \quad \text{(independence assumption)}$$
$$= \prod_{i=1}^{n+1} p(y_i|y_{i-1})$$

Put another way, we assume that the tags form a Markov sequence (future tags are independent of the past tags given the current one). Let's now make additional assumptions about the observations as well

$$p(x_1, \dots, x_n | y_1, \dots, y_{n+1}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}, y_1, \dots, y_{n+1})$$
(chain rule)
$$\approx \prod_{i=1}^n p(x_i | y_i)$$
(independence assumption)

In other words, we say that the identity of each word only depends on the corresponding tags. This is a drastic assumption but still (often) leads to a reasonable tagging model, and simplifies our calculations. A more formal statement here is that the random variable X_i is conditionally independent of all the other variables in the model given Y_i (see more on conditional independence in the Bayesian networks lectures).

Now, we have a much simpler tagging model

$$p(x_1, \dots, x_n, y_1, \dots, y_{n+1}) = p(y_1) \prod_{i=2}^{n+1} p(y_i | y_{i-1}) \prod_{i=1}^n p(x_i | y_i)$$
(3)

For notational convenience, we also assume a special fixed START symbol $y_0 = *$ so that $p(y_1)$ becomes $p(y_1|y_0)$. As a result, we can write

$$p(x_1, \dots, x_n, y_1, \dots, y_{n+1}) = \prod_{i=1}^{n+1} p(y_i | y_{i-1}) \prod_{i=1}^n p(x_i | y_i)$$
(4)

Let's understand this model a bit more carefully by looking at how the pairs of sequences could be generated from the model. Here's the recipe

- 1. Set $y_0 = *$ (we always start from the START symbol) and let i = 1.
- 2. Generate tag y_i from the conditional distribution $p(y_i|y_{i-1})$ where y_{i-1} already has a value (e.g., $y_0 = *$ when i = 1)
- 3. If $y_i = \text{STOP}$, we halt the process and return $y_1, \ldots, y_i, x_1, \ldots, x_{i-1}$. Otherwise we generate x_i from the output/emission distribution $p(x_i|y_i)$
- 4. Set i = i + 1, and return to step 2.

HMM formal definition

The model we have defined is a Hidden Markov Model or HMM for short. An HMM is defined by a tuple $\langle N, \Sigma, \theta \rangle$, where

- N is the number of states $1, \ldots, N$ (assume the last state N is the final state, i.e. what we called "STOP" earlier).
- Σ is the alphabet of output symbols. For example, $\Sigma = \{$ "the", "dog" $\}$.
- $\theta = \langle a, b, \pi \rangle$ consists of three sets of parameters
 - Parameter $a_{i,j} = p(y_{next} = j | y = i)$ for i = 1, ..., N 1 and j = 1, ..., N is the probability of transitioning from state *i* to state *j*: $\sum_{k=1}^{N} a_{i,k} = 1$
 - Parameter $b_j(o) = p(x = o|y = j)$ for j = 1...N 1 and $o \in \Sigma$ is the probability of emitting symbol o from state $j: \sum_{o \in \Sigma} b_j(o) = 1$.
 - Parameter $\pi_i = p(y_1 = i)$ for $i = 1 \dots N$ specifies probability of starting at state $i: \sum_{i=1}^{N} \pi_i = 1$.

Note that θ is a vector of $N + N \cdot (N - 1) + (N - 1) \cdot |\Sigma|$ parameters.

Example:

- N = 3. States are $\{1, 2, 3\}$
- Alphabet $\Sigma = \{the, dog\}$
- Distribution over initial states: $\pi_1 = 1, \pi_2 = \pi_3 = 0.$
- Parameters $a_{i,j}$ are

	j = 1	j = 2	j = 3
i = 1	0.5	0.5	0
i = 2	0	0.8	0.2

• Parameters $b_i(o)$ are

	o = the	o = dog
i = 1	0.9	0.1
i = 2	0.1	0.9

An HMM specifies a probability for each possible (x, y) pair, where $x = (x_1, \ldots, x_n)$ is a sequence of symbols drawn from Σ and $y = (y_1, \ldots, y_n)$ is a sequence of states drawn from the integers $1, \ldots, (N-1)$.



Figure 1: Transition graph for the example.

$$p(x, y|\theta) = \pi_{y_1} \cdot$$
 (prob. of choosing y_1 as an initial step)
 $a_{y_n,N} \cdot$ (prob. of transitioning to the final step)

$$\prod_{i=2}^{n} a_{y_{i-1},y_i} \cdot$$
 (transition probability)

$$\prod_{i=1}^{n} b_{y_i}(x_i)$$
 (emission probability)

Consider the example: "the/1, dog/2, the/1". The probability of such sequence is:

$$\pi_1 b_1(the) a_{1,2} b_2(dog) a_{2,1} b_1(the) a_{2,3} = 0 \tag{5}$$

Parameter estimation

We will first look at the fully observed case (complete data case), where our training data contains both xs and ys. We will do maximum likelihood estimation. Consider the examples: $\Sigma = \{e, f, g, h\}, N = 3$. Observation: (e/1, g/2), (e/1, h/2), (f/1, h/2), (f/1, g/2). To find the MLE, we will simply look at the counts of events, i.e., the number of transitions between tags, the number of times we saw an output symbol together with an specific tag (state). After normalizing the counts to yield valid probability estimates, we get

$$a_{i,j} = \frac{count(i,j)}{count(i)} \tag{6}$$

$$a_{1,2} = \frac{count(1,2)}{count(1)} = \frac{4}{4} = 1, \ a_{2,2} = \frac{count(2,2)}{count(2)} = \frac{0}{4} = 0, \dots$$
(7)

where count(i, j) is the number of times we have (i, j) as two successive states and count(i) is the number of times state *i* appears in the sequence. Similarly,

$$b_i(o) = \frac{count(i \to o)}{count(i)} \tag{8}$$

$$b_1(e)\frac{count(1 \to e)}{count(1)} = \frac{2}{4} = 0.5, \dots$$
 (9)

where $count(i \rightarrow o)$ is the number of times we see that state *i* emits symbol *o*. count(i) was already defined above.

2 Decoding with HMM

Suppose now that we have the HMM parameters θ (see above) and the problem is to infer the underlying tags y_1, \ldots, y_n corresponding to an observed sequence of words x_1, \ldots, x_n . In other words, we wish to evaluate

$$\operatorname*{argmax}_{y_1, \dots, y_n} p(x_1, \dots, x_n, y_1, \dots, y_{n+1})$$
(10)

where

$$p(x_1, \dots, x_n, y_1, \dots, y_{n+1}) = \prod_{i=1}^{n+1} a_{y_{i-1}, y_i} \prod_{i=1}^n b_{y_i}(x_i)$$
(11)

and $y_0 = *, y_{n+1} = \text{STOP}$. Note that by defining a fixed starting state, we have again folded the initial state distribution π into the transition probabilities $\pi_i = a_{*,i}$, $i \in \{1, \ldots, N\}$ (where N = STOP).

One possible solution for finding the most likely sequence of tags is to do brute force enumeration. Consider the example: $\Sigma = \{\text{the, dog}\}, x = \text{"the the dog"}$. The possible state sequences include:

$$1\ 1\ 1\ 2\ \text{STOP}$$
 (12)

$$1\ 1\ 2\ 2\ \text{STOP}$$
 (13)

$$1 \ 2 \ 2 \ 2 \ \text{STOP}$$
 (14)

But there are $|\mathcal{T}|^n$ possible sequences in total! Solving the tagging problem by enumerating the tag sequences will be prohibitively expensive.

Viterbi algorithm

The HMM has a simple dependence structure (recall, tags form a Markov sequence, observations only depend on the underlying tag). We can exploit this structure in a dynamic programming algorithm.

Input: $x = x_1, \ldots, x_n$ and model parameters θ .

Output: $\operatorname{argmax}_{y_1, \dots, y_{n+1}} p(x_1, \dots, x_n, y_1, \dots, y_{n+1}).$

Now, let's look at a truncated version of the joint probability, focusing on the first k tags for any $k \in \{1, ..., n\}$. In other words, we define

$$r(y_1, \dots, y_k) = \prod_{i=1}^k a_{y_{i-1}, y_i} \prod_{i=1}^k b_{y_i}(x_i)$$
(16)

where y_k does not equal STOP. Note that our notation $r(y_1, \ldots, y_k)$ suppresses any dependence on the observation sequence. This is because we view x_1, \ldots, x_n as given (fixed). Note that, according to our definition,

$$p(x_1, \dots, x_n, y_1, \dots, y_{n+1}) = r(y_1, \dots, y_n) \cdot a_{y_n, y_{n+1}}$$
$$= r(y_1, \dots, y_n) \cdot a_{y_n, \text{STOP}}$$

Let S(k, v) be the set of tag sequences y_1, \ldots, y_k such that $y_k = v$. In other words, S(k, v) is a set of all sequences of length k whose last tag is v. The dynamic programming algorithm will calculate

$$\pi(k,v) = \max_{(y_1,\dots,y_k) \in S(k,v)} r(y_1,\dots,y_k)$$
(17)

recursively in the forward direction. In other words, $\pi(k, v)$ can be thought as solving the maximization problem partially, over all the tags y_1, \ldots, y_{k-1} with the constraint that we use tag v for y_k . If we have $\pi(k, v)$, then $\max_v \pi(k, v)$ evaluates $\max_{y_1,\ldots,y_k} r(y_1,\ldots,y_k)$. We leave v in the definition of $\pi(k, v)$ so that we can extend the maximization one step further as we unravel the model in the forward direction. More formally,

• Base case:

 $\pi(0,*) = 1$ (starting state, no observations) $\pi(0,v) = 0$, if $v \neq *$ (an actual state has observations)

This definition reflects the assumption that $y_0 = *$.

• Moving forward recursively: for any $k \in \{1, ..., n\}$

$$\pi(k,v) = \max_{u \in \mathcal{T}} \{ \pi(k-1,u) \cdot a_{u,v} \cdot b_v(x_k) \}$$

$$(18)$$

In other words, when extending $\pi(k-1, u)$, $u \in \mathcal{T}$, to $\pi(k, v)$, $v \in \mathcal{T}$, we must transition from $y_{k-1} = u$ to $y_k = v$ (part $a_{u,v}$) and generate the corresponding observation x_k (part $b_v(x_k)$). Then we maximize over the previous tag $y_{k-1} = u$ so that $\pi(k, v)$ only depends on the value of y_k .

Finally, we must transition from y_n to STOP so that

$$\max_{y_1,\dots,y_n} p(x_1,\dots,x_n,y_1,\dots,y_n,y_{n+1} = \text{STOP}) = \max_{v \in \mathcal{T}} \{\pi(n,v) \cdot a_{v,\text{STOP}}\}$$
(19)

The whole calculation can be done in time $O(n|\mathcal{T}|^2)$, linear in length, quadratic in the number of tags.

Now, having values $\pi(k, v)$, how do we reconstruct the most likely sequence of tags which we denote as $\hat{y}_1, \ldots, \hat{y}_n$? We can do this via *back-tracking*. In other words, at the last step, $\pi(n, v)$ represents maximizations of all y_1, \ldots, y_n such that $y_n = v$. What is the best value for this last tag v, i.e., what is \hat{y}_n ? It is

$$\hat{y}_n = \operatorname*{argmax}_{v} \left\{ \pi(n, v) a_{v, \text{STOP}} \right\}$$
(20)

Now we can fix \hat{y}_n and work backwards. What is the best value \hat{y}_{n-1} such that we end up with tag \hat{y}_n in position n? It is simply

$$\hat{y}_{n-1} = \operatorname*{argmax}_{u} \left\{ \pi(n-1, u) a_{u, \hat{y}_n} \right\}$$
(21)

and so on.

3 Hidden Variable Problem

When we no longer have the tags, we must resort to other ways of estimating the HMMs. It is not trivial to construct a model that agrees with the observations except in very simple scenarios. Here is one:

- We have an HMM with $N = 3, \Sigma = \{a, b, c\}$
- We see the following output sequences in training data: (a, b), (a, c), (a, b).

How would you choose the parameter values for $\pi_i, a_{i,j}$ and $b_i(o)$? A reasonable choice is:

$$\pi_1 = 1.0, \pi_2 = \pi_3 = 0 \tag{22}$$

$$b_1(a) = 1.0, b_1(b) = b_1(c) = 0$$
(23)

$$b_2(a) = 0, b_2(b) = b_2(c) = 0.5$$
 (24)

$$a_{1,2} = 1.0, a_{1,1} = a_{1,3} = 0 \tag{25}$$

$$a_{2,3} = 1.0, a_{2,1} = a_{2,2} = 0 \tag{26}$$

Expectation-Maximization (EM) for HMM

Suppose now that we have multiple observed sequences of outputs (no observed tags). We will denote these sequences with superscripts, i.e., x^1, x^2, \ldots, x^m . In the context of each sequence, we must evaluate a posterior probability over possible tag sequences. For estimation, we only need expected counts that are used in the re-estimation step (M-step). To this end, let $count(x^i, y, p \to q)$ be the numbers of times a transition from state p to state q occurs in a tag sequence y corresponding to observation x^i . We will only show here the derivations for transition probabilities; the equations for emission and initial state parameters are obtained analogously.

E-step: calculate expected counts, added across sequences

$$\overline{count}(u \to v) = \sum_{i=1}^{m} \sum_{y} p(y|x^{i}, \theta^{t-1}) count(x^{i}, y, u \to v)$$
(27)

M-step: re-estimate transition probabilities based on the expected counts

$$a_{u,v} = \frac{\overline{count}(u \to v)}{\sum_{k=1}^{N} \overline{count}(u \to k)}$$
(28)

where the denominator ensures that $\sum_{k=1}^{N} a_{u,k} = 1$.

The main problem in running the EM algorithm is calculating the sum over the possible tag sequences in the E-step.

$$\sum_{y} p(y|x^{i}, \theta^{t-1}) count(x^{i}, y, u \to v)$$
(29)

The sum is over an exponential number of possible hidden state sequences y. Next we will discuss a dynamic programming algorithm – forward-backward algorithm. The algorithm is analogous to the Viterbi algorithm for maximizing over the hidden states.

The Forward-Backward Algorithm for HMMs

Suppose we could efficiently calculate marginal posterior probabilities

$$p(y_j = p, y_{j+1} = q | x, \theta) = \sum_{y: y_j = p, y_{j+1} = q} p(y | x, \theta)$$
(30)

for any $p \in 1 \dots (N-1)$, $q \in 1 \dots N$, $j \in 1 \dots n$. These are the posterior probabilities that the state in position j was p and we transitioned into q at the next step. The probability is conditioned on the observed sequence x and the current setting of the model parameters θ . Now, under this assumption, we could rewrite the difficulty computation in Eq. 29 as:

$$\sum_{y} p(y|x^{i}, \theta^{t-1}) count(x^{i}, y, p \to q) = \sum_{j=1}^{n} p(y_{j} = p, y_{j+1} = q|x^{i}, \theta^{t-1})$$
(31)

The key remaining question is how to calculate these posterior marginals effectively. In other words, our goal is to evaluate $p(y_j = p, y_{j+1} = q | x^i, \theta^{t-1})$.

Now, consider a single observation sequence x_1, \ldots, x_n . We will make use of the following forward probabilities:

$$\alpha_p(j) = p(x_1, \dots, x_{j-1}, y_j = p|\theta) \tag{32}$$

for all $j \in 1...n$, for all $p \in 1...N - 1$. $\alpha_p(j)$ is the probability of emitting the symbols x_1, \ldots, x_{j-1} and ending in state p in position j without (yet) emitting the corresponding output symbol. These are analogous to the $\pi(k, v)$ probabilities in the Viterbi algorithm with the exception that $\pi(k, v)$ included generating the corresponding observation in position k. Note that, unlike before, we are summing over all the possible sequences of states that could give rise to the observations x_1, \ldots, x_{j-1} . In the Viterbi algorithm, we maximized over the tag sequences.

Similarly to the forward probabilities, we can define the backward probabilities:

$$\beta_p(j) = p(x_j, \dots, x_n | y_j = p, \theta)$$
(33)

for all $j \in 1...n$, for all $p \in 1...N - 1$. $\beta_p(j)$ is the probability of emitting symbols x_j, \ldots, x_n and transitioning into the final (STOP) state, given that we begun in state p in position j. Again, this definition involves summing over all the tag sequences that could have generated the observations from x_j onwards, provided that the tag at j is p.

Why are these two definitions useful? Suppose we had been able to evaluate α and β probabilities effectively. Then the marginal probability we were after could be calculated as:

$$p(y_j = p, y_{j+1} = q | x, \theta) = \frac{1}{Z} \alpha_p(j) a_{p,q} b_p(o_j) \beta_q(j+1)$$
$$Z = p(x_1, \dots, x_n | \theta) = \sum_p \alpha_p(j) \beta_p(j) \text{ for any } j = 1 \dots n$$

This is just the sum over all possible tag sequences that include the transition $y_j = p$ and $y_{j+1} = q$ and generates the observations, divided by the sum over all tag sequences that generate the observations. As a result, we obtain the relative probability of the transition, relative to all the alternatives given the observations, i.e., the posterior probability. Note that $\alpha_p(j)$ involves all the summations over tags y_1, \ldots, y_{j-1} , and $\beta_q(j+1)$ involves all the summations over the tags y_{j+2}, \ldots, y_n .

Let's finally discuss how we can calculate α and β .

As Fig. 2 shows, for every state sequence y_1, y_2, \ldots, y_n there is

p=1 • (1,1) • (2,1)



Figure 2: A path associated with state sequence.

- a path through graph that has the sequence of states $s, \langle 1, y_1 \rangle, \ldots, \langle n, y_n \rangle, f$.
- The path associated with state sequence y_1, \ldots, y_n has weight equal to $p(x, y|\theta)$.
- $\alpha_p(j)$ is the sum of weights at all paths from s to the state (j, p).
- $\beta_p(j)$ is the sum of weights at all paths from state $\langle j, p \rangle$ to the final state f.

Given an input sequence x_1, \ldots, x_n , for any $p \in 1 \ldots N, j \in 1 \ldots n$, the forward and backward probability can be calculated recursively.

Forward probability:

$$\alpha_p(j) = p(x_1, \dots, x_{j-1}, y_j = p|\theta) \tag{34}$$

• Base case:

$$\alpha_p(1) = \pi_p \quad \forall p \in 1 \dots N - 1 \tag{35}$$

• Recursive case

$$\alpha_p(j+1) = \sum_q \alpha_q(j) a_{q,p} b_q(x_j) \quad \forall p \in 1 \dots N - 1, j = 1 \dots n - 1$$
(36)

Backward probability:

$$\beta_p(j) = p(x_j, \dots, x_n | y_j = p, \theta)$$
(37)

• Base case:

$$\beta_p(n) = a_{p,N} b_p(x_n) \quad \forall p \in 1 \dots N - 1 \tag{38}$$

• Recursive case

$$\beta_p(j) = \sum_q a_{p,q} b_p(x_j) \beta_q(j+1) \quad \forall p \in 1 \dots N - 1, j = 1 \dots n - 1$$
(39)