

# NLP: Hidden Markov Models

Dan Garrette  
dhg@cs.utexas.edu

October 16, 2013

## 1 Tagging

- Named entities
- Parts of speech

## 2 Parts of Speech

Tagsets

- Google Universal Tagset, 12: Noun, Verb, Adjective, Adverb, Pronoun, Determiner, Adposition (prepositions and postpositions), Numerals, Conjunctions, Particles, Punctuation, Other
- Penn Treebank, 45. Includes 4 categories of noun, 4 categories of pronoun, and 6 categories of verb.
- Brown Corpus, 87

Deciding what categories

- Auxilliary verbs? “I *had* gone”
- Numbers as adjectives? “I saw 4 cars”
- Count vs. mass nouns? “some apple” vs. “two apples”, “some snow” vs. “two snows”

Uses

- Parsing: determiner and noun should connect first, then to verb
- Speech synthesis: OBJECT (noun) vs. obJECT (verb), CONtent (noun) vs. conTENT (adj)

Rule-based

- “if it ends in ‘-tion’ ” → Noun
- “if it ends in ‘-ize’ ” → Verb
- “if it starts with ‘re-’ ” → Verb
- “if it starts with a capital letter” → Proper Noun
- “if it’s preceded by ‘the’ ” → Noun
- “if it’s followed by ‘s’ ” → Noun
- “if it’s preceded by an adjective” → Noun
- Or just memorize a big list of words and tags?
  - Ambiguity?
  - Picking the most common tag for a word → 90% accuracy (thought state-of-the-art is 97%)

Open vs. Closed

- Open class tags: new words are frequently added
  - noun, verb, adjective, adverb
  - “Twitter”, to “tweet”, “twitterish”
- Closed class tags: new words are rarely added
  - determiner, preposition, pronouns, ...
  - Don’t want to completely close off new words
  - Maybe *alongside* (preposition) wasn’t seen in training, only test
  - New domains: “u” as a pronoun

Complexities:

- Ambiguity:
  - “buy a book” (noun) vs. “book a flight” (verb)
  - “talk over the deal” (particle) vs. “talk over the phone” (prep)
- Phrasal verbs: “turn down”, “rule out”
  - “he went on for days”
  - “he went on a boat”

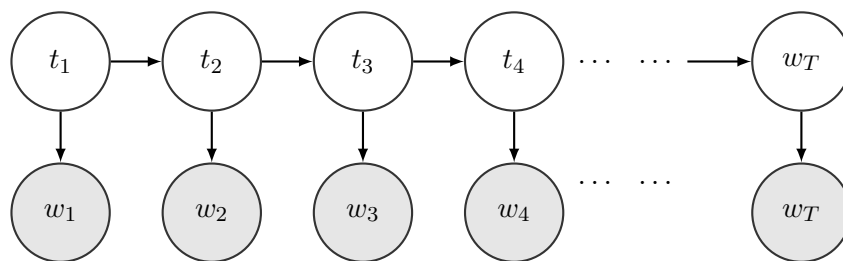


Figure 1: Second-order HMM showing conditional dependencies

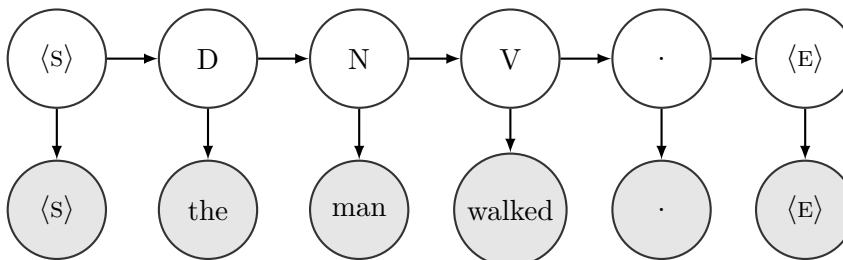


Figure 2: Second-order HMM for the sentence “the man walks .”

### 3 Hidden Markov Models

Similar to N-Gram models

- Model the text as a *sequence*
  - Bad assumption, but less sparse
- For ngrams, we modeled the probability of each word conditioned on the previous n-1 words.
- Here, we model each *tag* conditioned on previous *tags*
- Still uses Markov assumption: only look back a few tags
  - Again, bad assumption, but less sparse
- Note that we have to condition words on tags because otherwise

\\ Day 2

Derivation

- We want the most likely tag sequence for a sequence of words:

$$p(\langle S \rangle t_1 t_2 \dots t_n \langle E \rangle \mid \langle S \rangle w_1 w_2 \dots w_n \langle E \rangle)$$

Remember that order matters!

- For simplicity, we'll write this as  $t_1^n = \langle t_1, t_2, \dots, t_n \rangle$

- So we want

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} p(\langle S \rangle t_1 t_2 \dots t_n \langle E \rangle \mid \langle S \rangle w_1 w_2 \dots w_n \langle E \rangle)$$

but it's hard to estimate

- Like for ngrams, use Bayes rule:

$$\begin{aligned} \hat{t}_1^n &= \operatorname{argmax}_{t_1^n} \frac{p(\langle S \rangle w_1 w_2 \dots w_n \langle E \rangle \mid \langle S \rangle t_1 t_2 \dots t_n \langle E \rangle) \cdot p(\langle S \rangle t_1 t_2 \dots t_n \langle E \rangle)}{p(\langle S \rangle w_1 w_2 \dots w_n \langle E \rangle)} \\ &= \operatorname{argmax}_{t_1^n} p(\langle S \rangle w_1 w_2 \dots w_n \langle E \rangle \mid \langle S \rangle t_1 t_2 \dots t_n \langle E \rangle) \cdot p(\langle S \rangle t_1 t_2 \dots t_n \langle E \rangle) \end{aligned}$$

- Two major independence assumptions:

- Like ngrams, assume probability of a sequence is dependent only on recent past:

$$\begin{aligned} p(\langle S \rangle t_1 t_2 \dots t_n \langle E \rangle) &\approx p(t_1 \mid \langle S \rangle) \cdot p(t_2 \mid t_1) \cdot p(t_2 \mid t_2) \cdot \dots \cdot p(t_n \mid t_{n-1}) \cdot p(\langle E \rangle \mid t_n) \\ &= \prod_{i=1}^n p(t_i \mid t_{i-1}) \end{aligned}$$

- Also assume word is only dependent on its tag:

$$\begin{aligned} p(\langle S \rangle w_1 w_2 \dots w_n \langle E \rangle \mid \langle S \rangle t_1 t_2 \dots t_n \langle E \rangle) &\approx p(w_1 \mid t_1) \cdot p(w_2 \mid t_2) \cdot \dots \cdot p(w_n \mid t_n) \\ &= \prod_{i=1}^n p(w_i \mid t_i) \end{aligned}$$

- Together:

$$\begin{aligned} \hat{t}_1^n &= \operatorname{argmax}_{t_1^n} p(\langle S \rangle t_1 t_2 \dots t_n \langle E \rangle \mid \langle S \rangle w_1 w_2 \dots w_n \langle E \rangle) \\ &= \operatorname{argmax}_{t_1^n} p(\langle S \rangle w_1 w_2 \dots w_n \langle E \rangle \mid \langle S \rangle t_1 t_2 \dots t_n \langle E \rangle) \cdot p(\langle S \rangle t_1 t_2 \dots t_n \langle E \rangle) \\ &\approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n p(w_i \mid t_i) \cdot \prod_{i=1}^n p(t_i \mid t_{i-1}) \\ &= \operatorname{argmax}_{t_1^n} \prod_{i=1}^n p(w_i \mid t_i) \cdot p(t_i \mid t_{i-1}) \end{aligned}$$

## 4 Estimating Parameters: MLE

Two probability distributions to estimate:

- Transitions: probability of a tag, given previous tag,  $p(t_i \mid t_{i-1})$
- Emissions: probability of a word, given its tag,  $p(w_i \mid t_i)$

MLE

- MLE estimation is just like before (naïve Bayes, ngrams, ...): normalized counts

- Transitions:  $p(t_i | t_{i-1}) = \frac{C(t_{i-1} t_i)}{\sum_x C(t_{i-1} x)} = \frac{C(t_{i-1} t_i)}{C(t_{i-1})}$
- Emissions:  $p(w_i | t_i) = \frac{C(t_i, w_i)}{\sum_x C(t_i, x)} = \frac{C(t_i, w_i)}{C(t_i)}$

Example dataset (punctuation excluded for simplicity!):

```
<S>|<S> the|D man|N walks|V the|D dog|N <E>|<E>
<S>|<S> the|D dog|N runs|V <E>|<E>
<S>|<S> the|D dog|N walks|V <E>|<E>
<S>|<S> the|D man|N walks|V <E>|<E>
<S>|<S> a|D man|N saw|V the|D dog|N <E>|<E>
<S>|<S> the|D cat|N walks|V <E>|<E>
```

Some probabilities:

- $p(t_i = V | t_{i-1} = N) = \frac{C(N V)}{\sum_x C(N x)} = \frac{6}{8} = 0.75$
- $p(t_i = D | t_{i-1} = N) = \frac{C(N D)}{\sum_x C(N x)} = \frac{0}{8} = 0.0$
- $p(w_i = \text{dog} | t_i = N) = \frac{C(N, \text{dog})}{\sum_x C(N, x)} = \frac{4}{8} = 0.50$
- $p(w_i = \text{the} | t_i = N) = \frac{C(N, \text{the})}{\sum_x C(N, x)} = \frac{0}{8} = 0.0$

## 5 Add- $\lambda$ Smoothing

Again, just like before.

- Transitions:  $p(t_i | t_{i-1}) = \frac{C(t_{i-1} t_i) + \lambda}{\sum_x (C(t_{i-1} x) + \lambda)} = \frac{C(t_{i-1} t_i) + \lambda}{(\sum_x C(t_{i-1} x) + \lambda \cdot |T|)} = \frac{C(t_{i-1} t_i)}{C(t_{i-1})}$
- Emissions:  $p(w_i | t_i) = \frac{C(t_i, w_i) + \lambda}{\sum_x (C(t_i, x) + \lambda)} = \frac{C(t_i, w_i) + \lambda}{(\sum_x C(t_i, x) + \lambda \cdot |V|)} = \frac{C(t_i, w_i)}{C(t_i)}$

Some probabilities (with  $\lambda = 2$ ):

- The Vocabulary  $V$  is the set of all word types that might be associated with a tag  $t_i$ :  
 $V = \{a, cat, dog, man, runs, saw, the, walks\}, \quad |V| = 8$
- The Tagset  $T$  is the set of all possible tags that could follow a tag  $t_{i-1}$ :  
 $T = \{D, N, V, \langle E \rangle\}, \quad |T| = 4$
- $p(t_i = V | t_{i-1} = N) = \frac{C(N V) + 2}{\sum_{x \in T} (C(N x) + 2)} = \frac{C(N V) + 2}{(\sum_{x \in T} C(N x) + (2 \cdot 4))} = \frac{6+2}{8+8} = \frac{8}{16} = 0.50$
- $p(t_i = D | t_{i-1} = N) = \frac{C(N D) + 2}{\sum_{x \in T} (C(N x) + 2)} = \frac{C(N D) + 2}{(\sum_{x \in T} C(N x) + (2 \cdot 4))} = \frac{0+2}{8+8} = \frac{2}{16} = 0.125$

- $p(w_i = \text{dog} \mid t_i = N) = \frac{C(N, \text{dog})}{\sum_{x \in V} (C(N, x) + 2)} = \frac{C(N, \text{dog})}{(\sum_{x \in V} C(N, x)) + (2 \cdot 8)} = \frac{4+2}{8+16} = 0.25$
- $p(w_i = \text{the} \mid t_i = N) = \frac{C(N, \text{the})}{\sum_{x \in V} (C(N, x) + 2)} = \frac{C(N, \text{the})}{(\sum_{x \in V} C(N, x)) + (2 \cdot 8)} = \frac{0+2}{8+16} = 0.08$

Differences:

- Two distributions, two kinds of smoothing
- Can do add- $\lambda$  for both, but don't need to use same  $\lambda$
- Can use totally different smoothing for each.

## 6 One-Count Smoothing

From Chen and Goodman (1996)

Idea: Smooth conditional distribution by interpolating the MLE with a unigram distribution where the degree of smoothing is proportional to the “openness” of a tag.

$$p(t_i \mid t_{i-1}) = \frac{C(t_{i-1}, t_i) + \alpha(t_{i-1}) \cdot p(t_i)}{C(t_{i-1}) + \alpha(t_{i-1})} \quad \text{where } \alpha(t_{i-1}) = |t_i : C(t_{i-1}, t_i) = 1|$$

$$p(w_i \mid t_i) = \frac{C(t_i, w_i) + \beta(t_i) \cdot p(w_i)}{C(t_i) + \beta(t_i)} \quad \text{where } \beta(t_i) = |w_i : C(t_i, w_i) = 1|$$

- For emissions, if a tag  $t_i$  appears one time with a large number of word types, then we assume that it is very “open”: there is a large number of words that it has only been seen with once, so the likelihood that it *could have* been seen with some other word, even though it didn't *happen* to be seen with it, is high.
- For transitions, the reasoning is the same. Instead of a tag being associated with a wide variety of words once, it's a wide variety of subsequent tags once.
- Note that it is still necessary to smooth the prior distributions  $p(t_{i-1})$  and  $p(w_i)$  since if the transition or emission count is zero, we need to ensure that we back off to some non-zero value; add-1 smoothing should be sufficient. Additionally, It is necessary to add some small amount to  $\alpha(t_{i-1})$  and  $\beta(t_i)$  in case there are no singletons

## 7 Three Tasks for HMM

1. Likelihood: Given a tagged sequence, determine its likelihood
2. Decoding: Given an untagged sequence, determine the best tag sequence for it
3. Learning: Given an untagged sequence, and a set of tags, learn the HMM parameters.

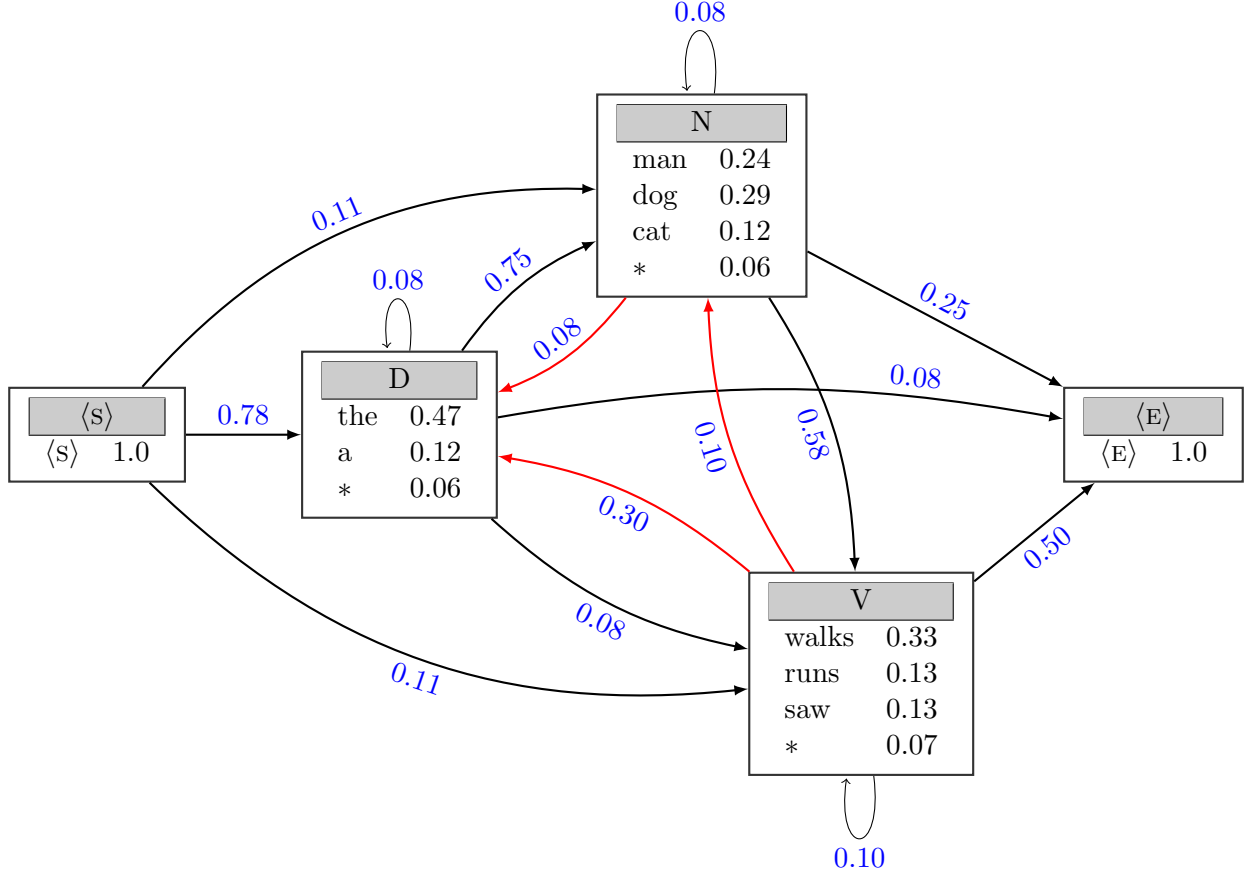


Figure 3: Finite state machine of the data with add-1 smoothing. Missing arrows are assumed to be zero probabilities. With smoothing, there is an arrow in *both directions* between every pair of tags.

## 8 Likelihood of a tagged sentence

We can compute the likelihood of a particular sequence of tags for a sentence:

- $p(t_1 \dots t_n \mid w_1 \dots w_n) = \prod_{i=1}^n p(w_i \mid t_i) \cdot p(t_i \mid t_{i-1})$

Example: “the|D dog|N walks|V”

$$\begin{aligned}
 p(t_1 \dots t_n \mid w_1 \dots w_n) &\approx \prod_{i=1}^n p(w_i \mid t_i) \cdot p(t_i \mid t_{i-1}) \\
 &= p(D \mid \langle S \rangle) \cdot p(\text{the} \mid D) \cdot p(N \mid D) \cdot p(\text{dog} \mid N) \cdot p(V \mid N) \cdot p(\text{walks} \mid V) \cdot p(\langle E \rangle \mid V)
 \end{aligned}$$

- This says that the probability of the tag sequence is the probability of seeing  $D$  as the first tag,  $\text{the}$  given that the chosen first tag was  $D$ ,  $N$  given that the previous tag was  $D$ ,  $\text{dog}$  given that the chosen second tag was  $N$ , etc...
- Transition ( $p(t_i \mid t_{i-1})$ ) and emission ( $p(w_i \mid t_i)$ ) values are calculated as above

## 9 Decoding

Tagging a sentence

- $\hat{t}_i^n = \operatorname{argmax}_{t_1^n} p(t_1 \dots t_n \mid w_1 \dots w_n) \approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n p(w_i \mid t_i) \cdot p(t_i \mid t_{i-1})$
- For “the dog walks”, we have to check all combinations:  
 $\langle S \rangle D D D \langle E \rangle$   
 $\langle S \rangle N D D \langle E \rangle$   
 $\langle S \rangle V D D \langle E \rangle$   
 $\langle S \rangle . D D \langle E \rangle$   
 $\langle S \rangle N N D \langle E \rangle$   
 $\langle S \rangle N V D \langle E \rangle$   
...  
• But this is insane.

\\ Day 3

Choosing a tag

- Because of our independence assumptions, the probability of each tag is dependent only on its “Markov blanket”: the previous tag, emitted word, and next tag.

$$\begin{aligned} p(t_i \mid \langle S \rangle w_1 w_2 \dots w_n \langle E \rangle, \langle S \rangle t_1 t_2 \dots \mathbf{t}_{i-1} \mathbf{t}_{i+1} \dots t_n \langle E \rangle) &= p(t_i \mid t_{i-1}, w_i, t_{i+1}) \\ &= p(t_i \mid t_{i-1}) \cdot p(w_i \mid t_i) \cdot p(t_{i+1} \mid t_i) \end{aligned}$$

Tagging a sentence

- We want to make global decisions about tags.
- Due to independence, global decisions are made in terms of local decisions
- Tagging “forward” through the sentence allows us to predict tags based on previous decisions.
- But we can’t calculate picking a tag changes the probabilities of previous tags.

The Viterbi Algorithm

- We need to check all possible tag combinations *efficiently*: dynamic programming
- We do this by making two passes through the sentence: one forward through the sentence and one backward
- The forward pass



- Compute, for each potential tag value  $k$  of  $t_i$  in the sentence, the highest probability sequence of tags for the previous  $i-1$  words and ending with  $t_i = k$ , given *all possible ways of tagging the  $i-1$  previous words*.

$$v_i(k) = \max_{t_1, \dots, t_{i-1}} p(w_1, \dots, w_i, t_1, \dots, t_{i-1}, t_i = k)$$

- With our independence assumptions:

$$v_i(k) = p(w_i | k) \cdot \max_{k' \in T} v_{i-1}(k') \cdot p(k | k')$$

- $v_i(k)$  is the probability of the most likely subsequence of tags that accounts for the first  $i$  words and ends with tag  $t_i = k$ .
- So  $k$  is the potential tag for the current token  $i$
- $k'$  is a potential tag for the previous token  $i-1$ .
- In terms of the **Markov blanket**, this covers the connections to the previous tag and the emitted word

- The backward pass

- Starting at the end of the sentence with  $\hat{t}_N = \langle E \rangle$ , use each chosen “next tag” along with the viterbi scores (probability of assigning a tag  $k$  to  $t_i$  given the best possible sequence of tags), to figure out the most likely tag  $k$  for  $t_i$ .
- So:

$$\hat{t}_i = \operatorname{argmax}_{k \in T} v_i(k) \cdot P(\hat{t}_{i+1} | t_i = k)$$

- So the chosen tag is one that has a maximally plausible preceding sequence of tags and best connects to the just-selected “next” tag.
- In terms of the **Markov blanket**, this adds in the connection to the next tag (since it has now been decided) to the forward probability which is made up of connections to the previous tag and the emitted word

**Example:** Viterbi algorithm for tagging the sentence “the dog walks” is given in the figures below

## 10 Tag Dictionary

- Mapping from words to potential tags
- Get from train, assume it works on test
- Unseen words usually assumed to be any possible tag. Could, instead, assume only open-class tags.
- Prune low-probability tags.

## 11 Learning: Forward-Backward

1. Given some unlabeled text, and a set of labels, learn the best set of parameters for an HMM tagger.
2. Usually needs a tag dictionary to constrain the search space.

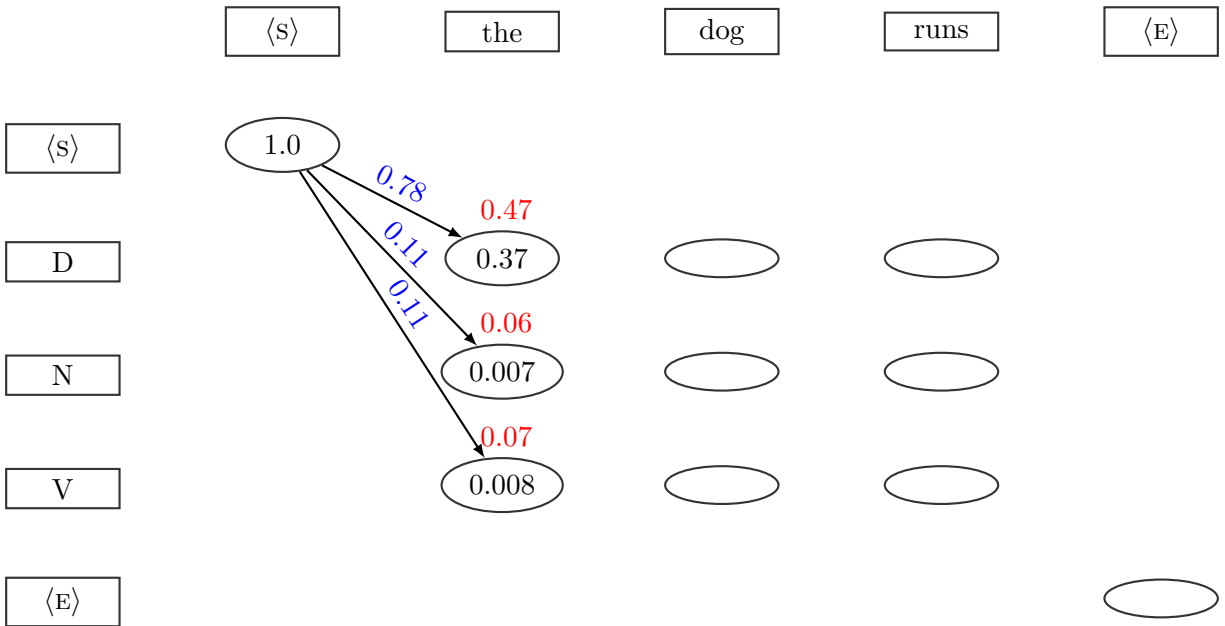


Figure 4: Viterbi Algorithm: Calculating  $v_1(D)$ ,  $v_1(N)$ , and  $v_1(V)$ . Since  $\langle s \rangle$  is only choice for the previous tag, each score  $v_1(k)$  is calculated as  $p(\text{the} \mid k) \cdot v_0(\langle s \rangle) \cdot p(k \mid \langle s \rangle)$ . So, for example,

$$v_1(D) = p(\text{the} \mid D) \cdot v_0(\langle s \rangle) \cdot p(D \mid \langle s \rangle) = 0.47 \cdot 1.0 \cdot 0.78 = 0.37$$



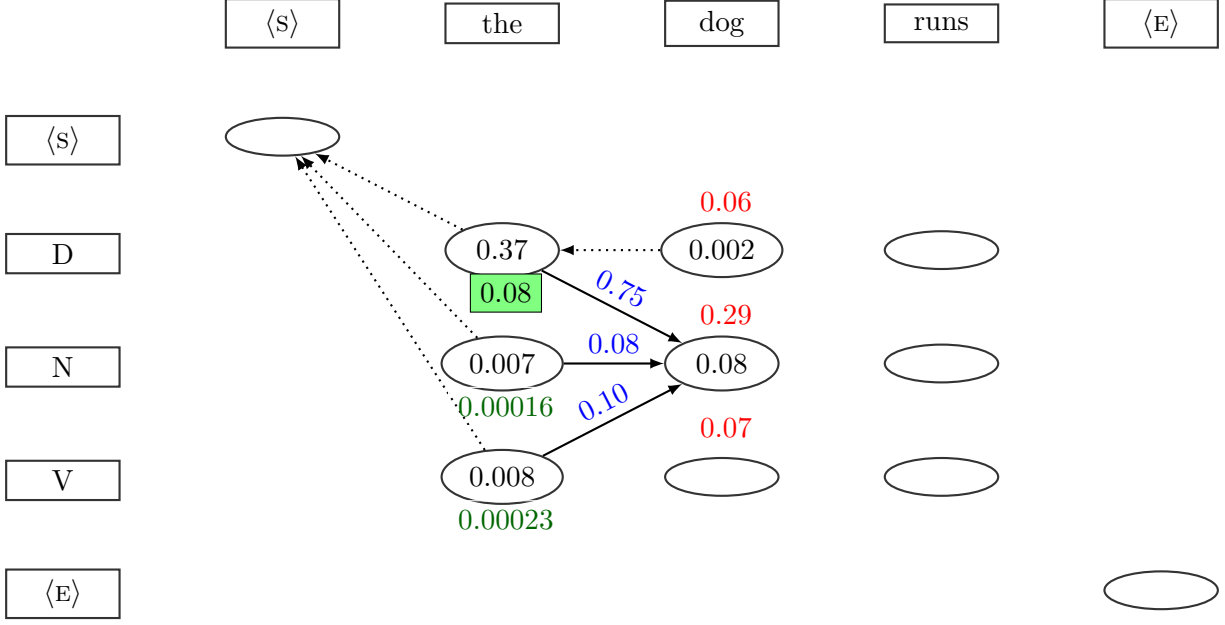


Figure 6: Viterbi Algorithm: Calculate  $v_2(N)$  and find the best path to it, which turns out to be the path ending with  $t_1 = D$ .

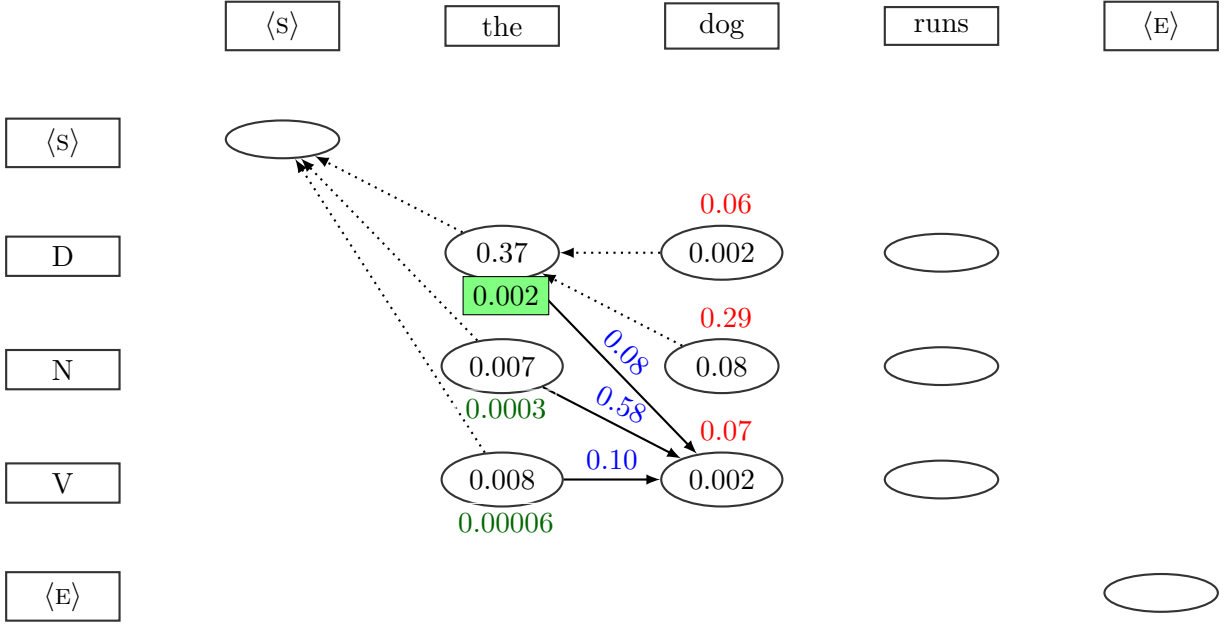


Figure 7: Viterbi Algorithm: Calculate  $v_2(V)$  and find the best path to it, which turns out to be the path ending with  $t_1 = D$ .

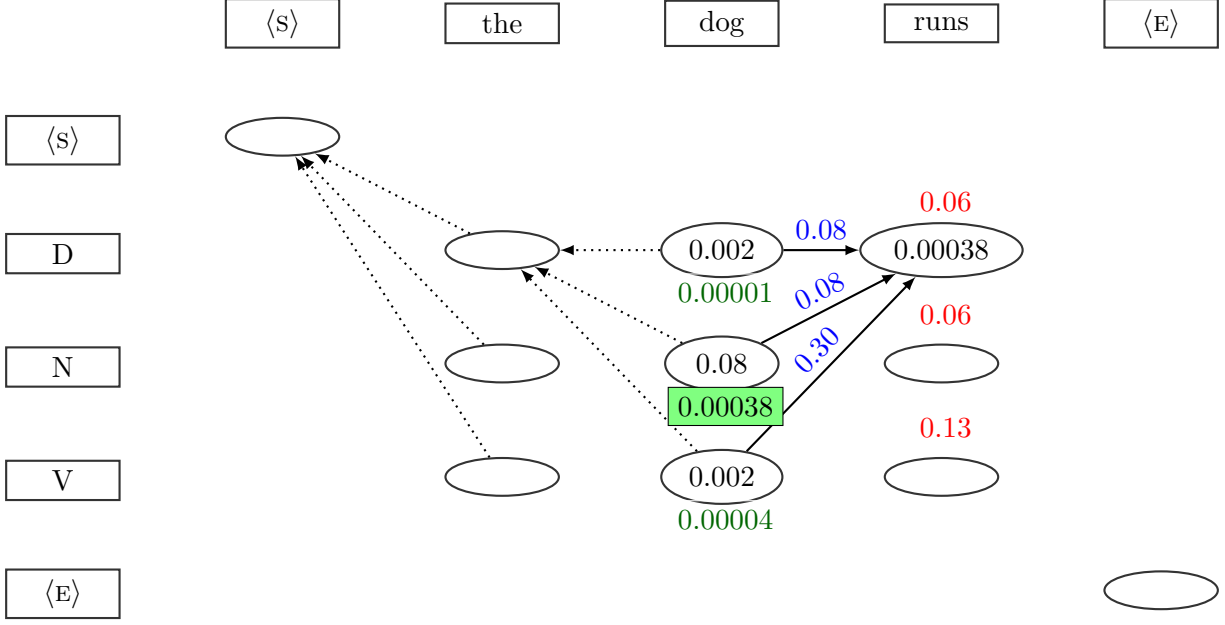


Figure 8: Viterbi Algorithm: Calculate  $v_3(D)$  and find the best path to it, which turns out to be the path ending with  $t_2 = N$

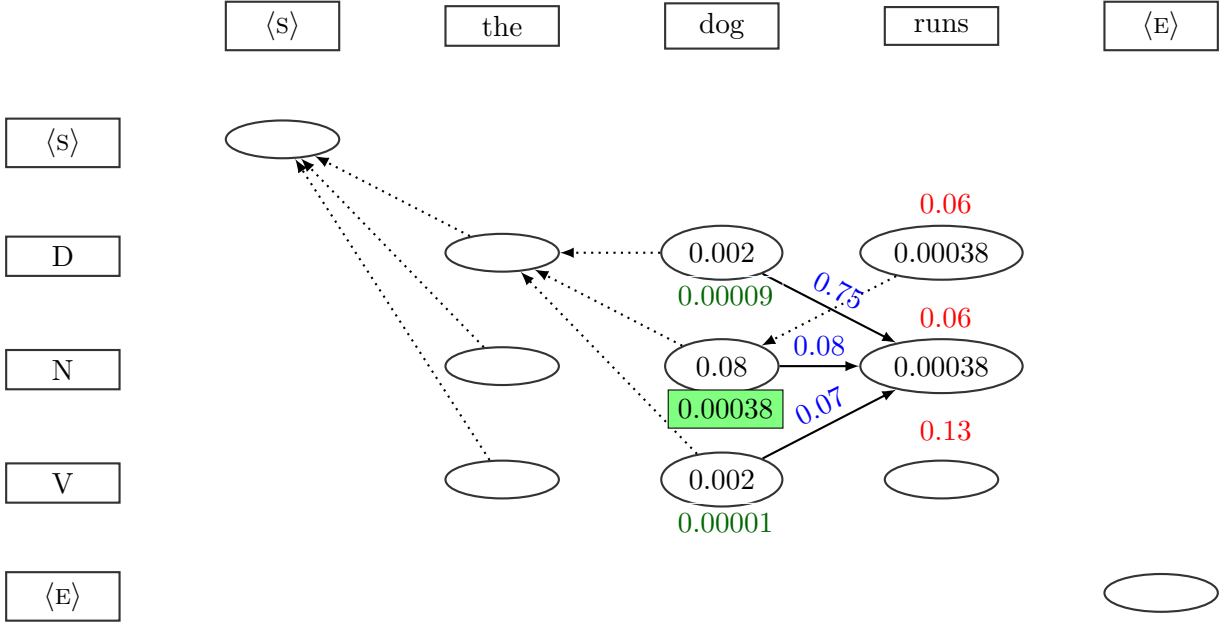


Figure 9: Viterbi Algorithm: Calculate  $v_3(N)$  and find the best path to it, which turns out to be the path ending with  $t_2 = N$



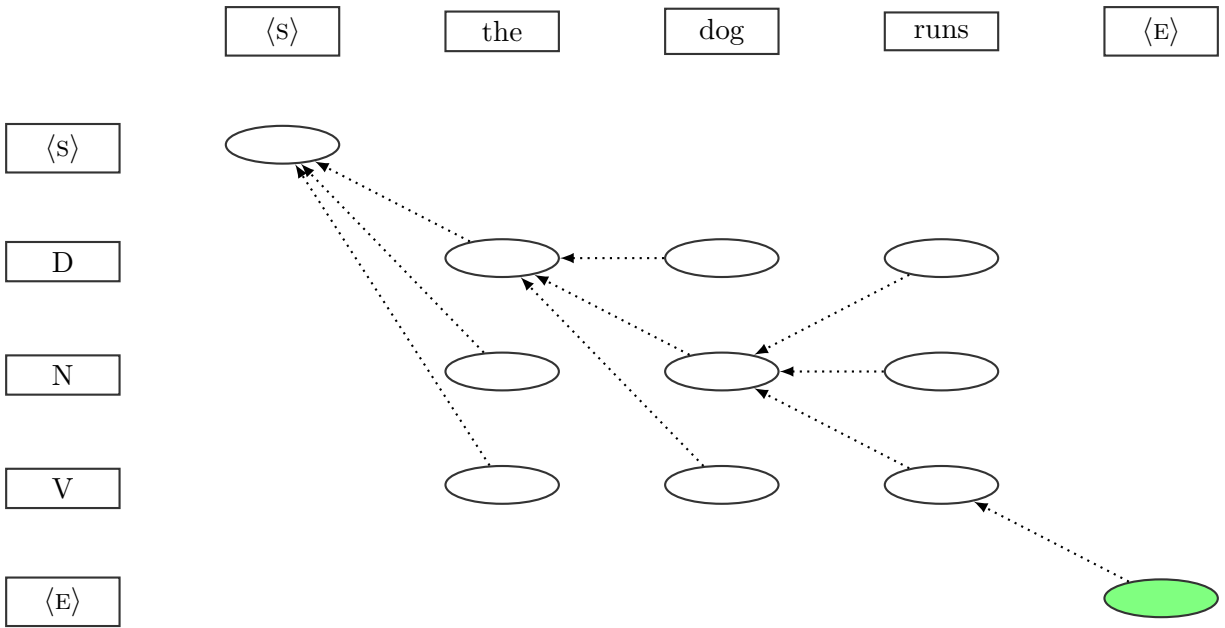


Figure 12: Viterbi Algorithm: The tag for  $\langle E \rangle$  is always  $\langle E \rangle$ , so we can always start by selecting that tag.

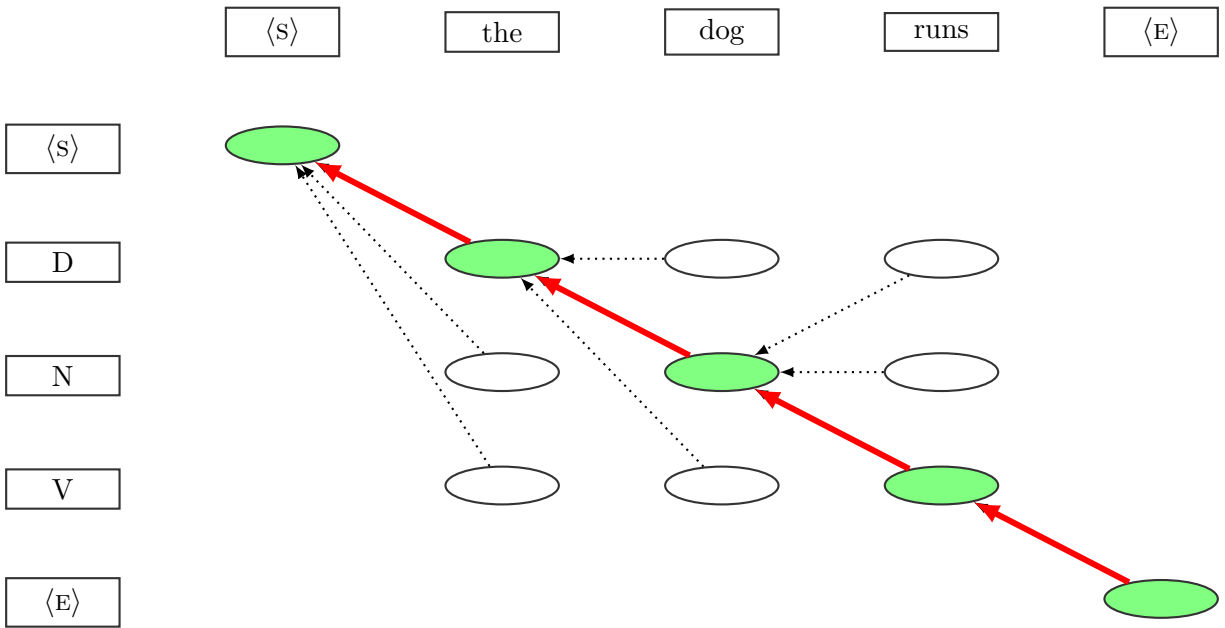


Figure 13: Viterbi Algorithm: For each tag  $t_{i+1}$ , we examine the backpointer that indicates the way to the best path to  $t_{i+1}$ , and choose  $t_i$  accordingly.