

Statistical methods in NLP

Part-of-speech tagging



**UNIVERSITY OF
GOTHENBURG**

Richard Johansson

February 23, 2016

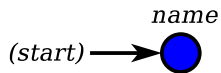
hidden Markov models



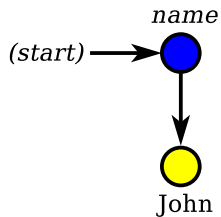
$$P(t_n | t_{n-1}) \quad P(w_n | t_n)$$

- ▶ a model where we have an unknown underlying sequence is called a **hidden Markov** model (HMM)

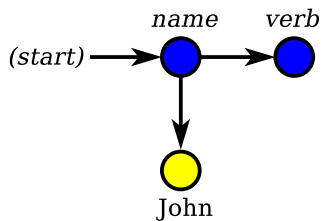
generative story in hidden Markov models



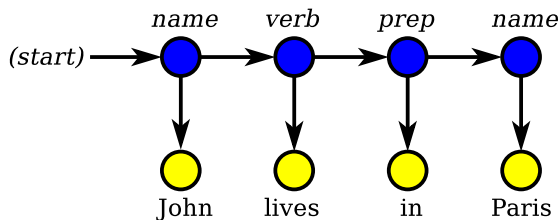
generative story in hidden Markov models



generative story in hidden Markov models



generative story in hidden Markov models



smoothing...

- ▶ **smoothing** may be useful, in particular if the corpus is small
- ▶ for instance, Laplace smoothing for transition probabilities:

$$P(t_n|t_{n-1}) = \frac{\text{count}(t_{n-1}, t_n) + \lambda}{\text{count}(t_{n-1}) + \lambda \cdot T}$$

where T is the number of distinct tags

- ▶ and for emission probabilities:

$$P(w|t) = \frac{\text{count}(w, t) + \lambda}{\text{count}(t) + \lambda \cdot V}$$

where V is the number of distinct words

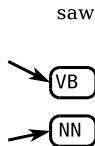
- ▶ usually there is some “special treatment” for the emission probability $P(w_n|t_n)$ if w_n is unseen in the training corpus
 - ▶ taking for instance punctuation, capitalization, numbers, suffixes into account

tagging

- ▶ how do we use our probability model to tag?
- ▶ conceptually: enumerate all possible tag sequences; use the probabilities to find the best one
- ▶ however in long sentences, the number of possible tag sequences is very large
- ▶ the **Viterbi algorithm** finds the most probable underlying tag sequence
 - ▶ Viterbi runs in **linear time** with respect to the length of the sentence

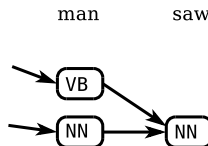
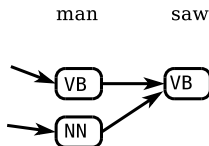
the Viterbi algorithm

- ▶ for each possible tag t_i of a word w_i , we compute **the best tag sequence leading to t_i**
 - ▶ for instance: for the word *saw*, we find the best sequence ending with *saw* as a verb, and the best ending with *saw* as a noun



the trick

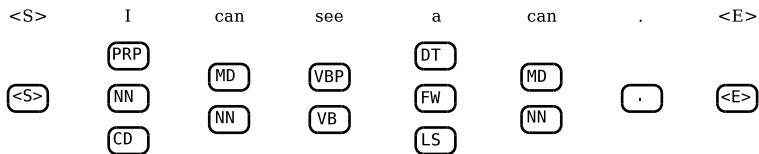
- ▶ to compute the best path ending with *saw* as a verb, **consider the best paths for the previous word** and the **transition probabilities**
- ▶ assume the previous word is e.g. *man*, which can be a noun or a verb
- ▶ select the highest of
 - ▶ the LP of the best path ending in *man* as a verb + the LP of the transition verb \rightarrow verb
 - ▶ the LP of the best path ending in *man* as a noun + the LP of the transition noun \rightarrow verb



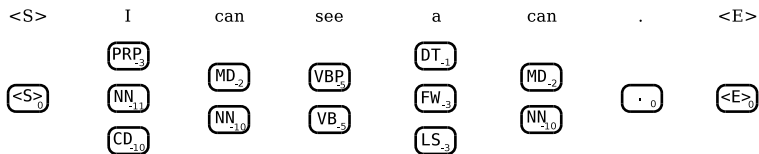
Viterbi example

<S> I can see a can . <E>

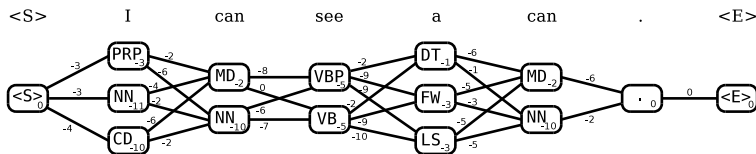
Viterbi example



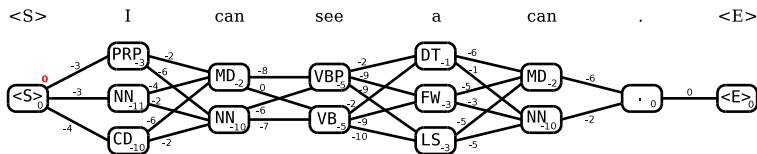
Viterbi example



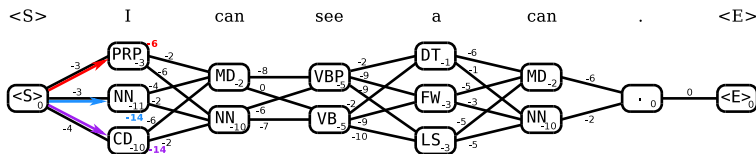
Viterbi example



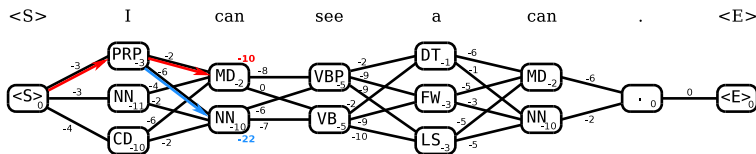
Viterbi example



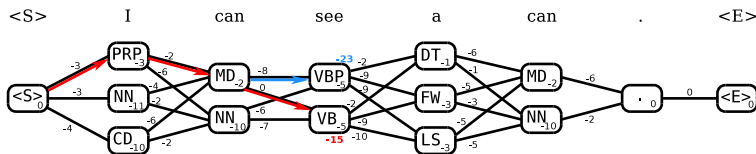
Viterbi example



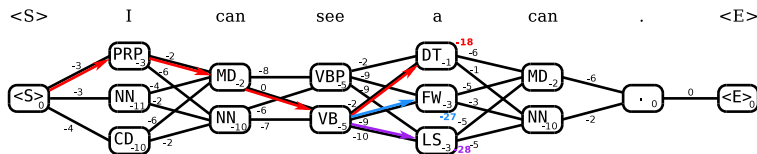
Viterbi example



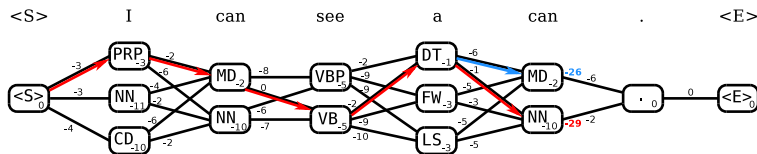
Viterbi example



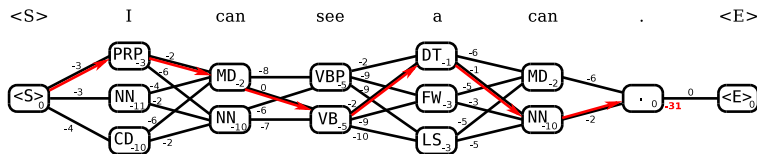
Viterbi example



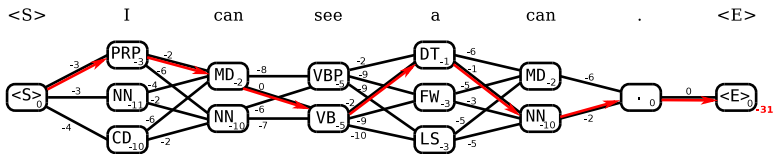
Viterbi example



Viterbi example

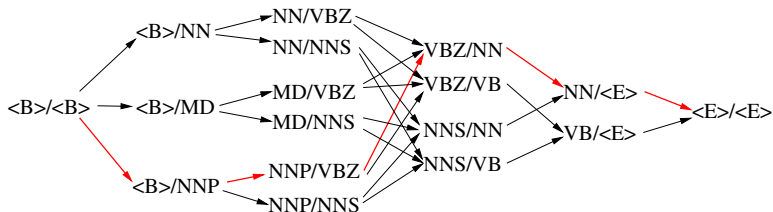
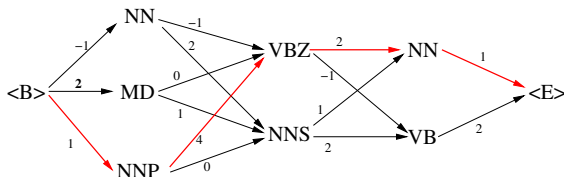


Viterbi example



Search spaces...

- ▶ example: *Will plays golf*



some hints: estimation (2)

- ▶ it can be useful to use a “double dictionary” pattern

```
word_tag_counter[word][tag] += 1
```

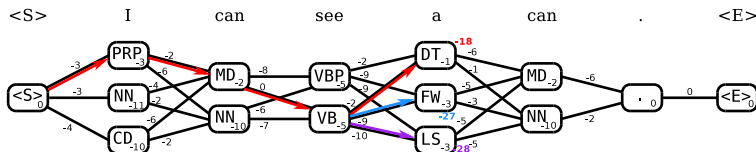
- ▶ to get rid of key checks that clutter the code, you can use a `defaultdict(Counter)`
 - ▶ (recall: `Counter` is a frequency table)

some hints: estimation (4)

- ▶ the code will be a bit simpler if you make sure that there are transition probabilities for all possible transitions
 - ▶ even those you haven't seen, so use smoothing!
- ▶ also, it's probably best if you use special tags for the start and end of the sentence

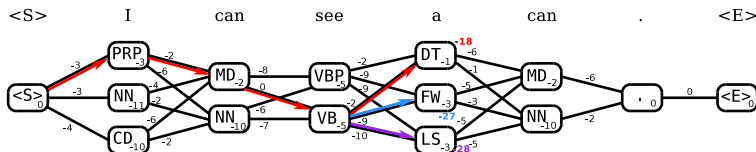
some hints: Viterbi (1)

- ▶ what do you think we should use to represent the **links** that we have drawn?



some hints: Viterbi (1)

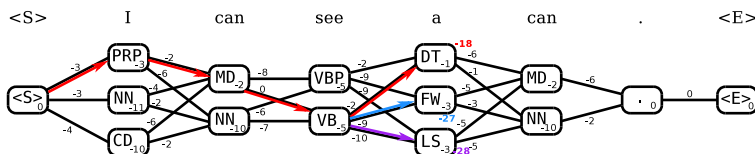
- ▶ what do you think we should use to represent the **links** that we have drawn?



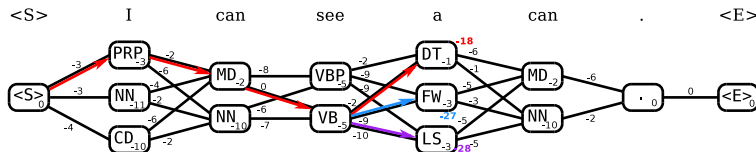
- ▶ we need something that remembers
 - ▶ the (log) **probability** of the path ending in that link
 - ▶ the **tag**
 - ▶ the **previous link**
- ▶ simplest solution is probably to use a **tuple** of these three

some hints: Viterbi (2)

- ▶ at each step of the algorithm, we keep a list of links
- ▶ there's one link for each possible tag at that step
- ▶ the link represents the best path leading up to that tag

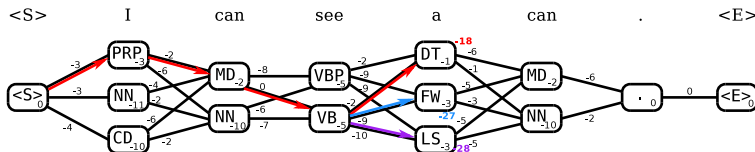


some hints: Viterbi (3)



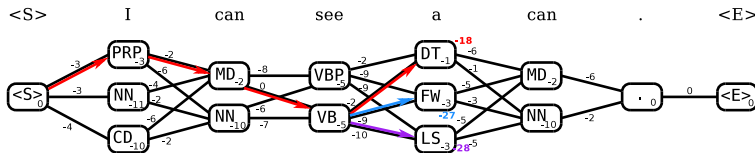
- ▶ in the first step, the list contains a single dummy link
 - ▶ the log probability is 0 (100% chance of being here!)
 - ▶ special start tag (for instance $\langle S \rangle$)
 - ▶ no previous link (use None or similar)

some hints: Viterbi (3)



- ▶ in the first step, the list contains a single dummy link
 - ▶ the log probability is 0 (100% chance of being here!)
 - ▶ special start tag (for instance <S>)
 - ▶ no previous link (use None or similar)
- ▶ let's assume we want to build the link for the DT tag
- ▶ find the previous link that maximizes the sum of
 - ▶ log probability in the previous link
 - ▶ log of transition probability to DT
 - ▶ log of emission probability of *a* for DT

some hints: Viterbi (4)



- ▶ in the last step, we end the sequence with another dummy link
- ▶ then we need to follow the links backwards to find the tag sequence
 - ▶ (see pseudocode in the Python file)

bootstrapping a confidence interval, pseudocode

- ▶ we have a test set T consisting of k sentences
- ▶ we compute a confidence interval by generating N random test sets and finding the interval where most estimates end up

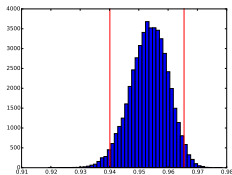
repeat N times

T^* = pick k sentences randomly from T

a = estimated accuracy of the tagger on T^*

store a in a list A

return 2.5% and 97.5% percentiles of A



bootstrapping a confidence interval in Python (part 1)

```
def bootstrap_ci(goldstandard, guess, N):  
  
    # for instance [ (8,9), (6,6), (12,14), ... ]  
    evals = [ evaluate_sentence(sen_gold, sen_guess)  
              for sen_gold, sen_guess in zip(goldstandard, guess) ]  
  
    # for instance [ 0.87, 0.92, 0.88, 0.89, ...]  
    A = [ random_testset_accuracy(evals) for _ in range(N) ]  
  
    lower = scipy.percentile(A, 2.5)  
    upper = scipy.percentile(A, 97.5)  
  
    return lower, upper
```

bootstrapping a confidence interval in Python (part 2)

```
def evaluate_sentence(sen_gold, sen_guess):  
    n_correct = sum(t1 == t2  
                    for t1, t2 in zip(sen_gold, sen_guess))  
    return n_correct, len(sen_gold)  
  
def random_testset_accuracy(evals):  
    # evals is for instance [ (8,9), (6,6), (12,14), ... ]  
  
    n_words = 0  
    n_correct = 0  
    for _ in range(len(evals)):  
        sentence_eval = random.choice(evals)  
        n_correct += sentence_eval[0]  
        n_words += sentence_eval[1]  
    return n_correct / n_words
```

bootstrapping for comparing to a fixed value

- ▶ is the tagging accuracy significantly greater than 0.94?
- ▶ we compute the p -value by checking how often the accuracy falls below 0.94

repeat N times

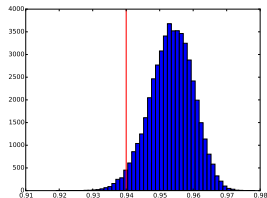
T^* = pick k sentences randomly from T

a = estimated accuracy of the tagger on T^*

if $a < 0.94$

increase counter C

return C/N



bootstrapping for comparing two taggers

- ▶ is tagger A significantly better than tagger B?
- ▶ we use a lot of randomly generated test sets, and compute the p -value by checking how often tagger B outperforms tagger A

repeat N times

T^* = pick k sentences randomly from T

a_A = estimated accuracy of tagger A on T^*

a_B = estimated accuracy of tagger B on T^*

if $a_B > a_A$

 increase counter C

return C/N

