

## The Viterbi Heuristic

In the previous segment, you learnt how to calculate the probability of a tag sequence given a sequence of words. The idea is to compute the probability of all possible tag sequences and assign the sequence having the maximum probability.

Although this approach can work in principle, it is computationally very expensive. For e.g. if you have just three POS tags - DT, JJ, NN, and you want to tag the sentence "The high cost", there are  $3^3 = 27$  possible tag sequences (each word can have three possible tags).

In general, for a sequence of  $n$  words and  $t$  tags, a total of  $t^n$  tag sequences are possible. The Penn Treebank dataset in NLTK itself has 36 POS tags, so for a sentence of length say 10, there are  $36^{10}$  possible tag sequences (that's about three thousand trillion!).

Clearly, computing trillions of probabilities to tag a 10-word sentence is impractical. Thus, we need to find a much more efficient approach to tagging.

String = "The high cost"

**PENN TREEBANK**

String = word1 word2 .... wordn

$36 \times 36 \times 36 \dots = 36^n$

## POS TAGGING

### Viterbi Heuristic

*W: Words and T: Tags*

$$\begin{aligned}
 & [P(W_1 | T_1) * P(W_2 | T_2) * P(W_3 | T_3) * P(T_1) * P(T_2 | T_1) * P(T_3 | T_2)] \\
 & = [P(W_1 | T_1) * P(T_1)] * [P(W_2 | T_2) * P(T_2 | T_1)] * [P(W_3 | T_3) * P(T_3 | T_2)]
 \end{aligned}$$

## VITERBI HEURISTIC

$P(\text{DT}|\text{start}) = 0.4$   
 $P(\text{VB}|\text{start}) = 0.3$   
 $P(\text{NN}|\text{start}) = 0.1$   
 $P(\text{JJ}|\text{start}) = 0.2$

## Tag For "The"

$P(\text{DT}|\text{start}) * P(\text{"The"}|\text{DT}) = 0.4 * 0.6 = 0.24$   
 $P(\text{NN}|\text{start}) * P(\text{"The"}|\text{NN}) = 0.1 * 0.2 = 0.02$   
 $P(\text{JJ}|\text{start}) * P(\text{"The"}|\text{JJ}) = 0.2 * 0.1 = 0.02$   
 $P(\text{VB}|\text{start}) * P(\text{"The"}|\text{VB}) = 0.3 * 0.1 = 0.03$

$P(\text{"The"}|\text{JJ}) = 0.1, P(\text{"high"}|\text{JJ}) = 0.4, P(\text{"cost"}|\text{JJ}) = 0.2,$   
 $P(\text{"The"}|\text{NN}) = 0.2, P(\text{"high"}|\text{NN}) = 0.2, P(\text{"cost"}|\text{NN}) = 0.3,$   
 $P(\text{"The"}|\text{DT}) = 0.6, P(\text{"high"}|\text{DT}) = 0.2, P(\text{"cost"}|\text{DT}) = 0.2,$   
 $P(\text{"The"}|\text{VB}) = 0.1, P(\text{"high"}|\text{VB}) = 0.2, P(\text{"cost"}|\text{VB}) = 0.3$

## Tag For "High"

$P(\text{DT}|\text{DT}) * P(\text{"high"}|\text{DT}) = 0.1 * 0.2 = 0.02$   
 $P(\text{NN}|\text{DT}) * P(\text{"high"}|\text{NN}) = 0.4 * 0.2 = 0.08$   
 $P(\text{JJ}|\text{DT}) * P(\text{"high"}|\text{JJ}) = 0.3 * 0.4 = 0.12$   
 $P(\text{VB}|\text{DT}) * P(\text{"high"}|\text{VB}) = 0.2 * 0.1 = 0.02$

## Tag For "Cost"

$P(\text{DT}|\text{JJ}) * P(\text{"cost"}|\text{DT}) = 0.2 * 0.2 = 0.04$   
 $P(\text{NN}|\text{JJ}) * P(\text{"cost"}|\text{NN}) = 0.5 * 0.3 = 0.15$   
 $P(\text{JJ}|\text{JJ}) * P(\text{"cost"}|\text{JJ}) = 0.2 * 0.2 = 0.04$   
 $P(\text{VB}|\text{JJ}) * P(\text{"cost"}|\text{VB}) = 0.1 * 0.3 = 0.03$

To summarise, the basic idea of the Viterbi algorithm is as follows - given a list of observations (words)  $O_1, O_2, \dots, O_n$  to be tagged, rather than computing the probabilities of all possible tag sequences, you assign tags sequentially, i.e. assign the most likely tag to each word using the previous tag.

More formally, you assign the tag  $T_j$  to each word  $O_i$  such that it maximises the likelihood:

$$P(T_j|O_i) = P(O_i|T_j) * P(T_j|T_{j-1}),$$

where  $T_{j-1}$  is the tag assigned to the previous word. Recall that according to the Markov assumption, the probability of a tag  $T_j$  is assumed to be **dependent only on the previous tag  $T_{j-1}$** , and hence the term  $P(T_j|T_{j-1})$ .

In other words, you assign tags sequentially such that the tag assigned to every word  $O_i$  maximises the likelihood  $P(T_j|O_i)$  locally, i.e. it is assumed to be dependent only on the current word and the previous tag. This algorithm does not guarantee that the resulting tag sequence will be the most likely sequence, but by making these simplifying assumptions, it reduces the computational time manifold (you'll see how much in the next lecture).

This is also why it is called a **greedy algorithm** - it assigns the tag that is most likely *at every word*, rather than looking for the overall most likely sequence. In the following video, Prof. Srinath will compare the computational costs of the two tagging algorithms - the brute force algorithm and the Viterbi algorithm.

In [ ]: