

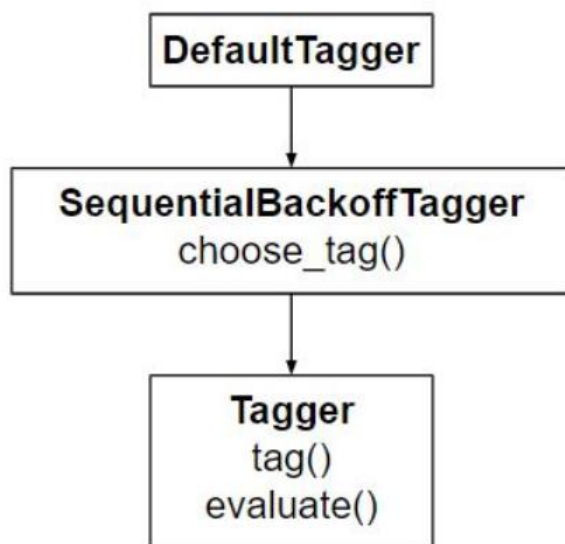
4. Explain POS tagging with HMM?

Parts of Speech Tagging (POS): It is a process of converting a sentence to forms – list of words, list of tuples (where each tuple is having a form (word, tag)). The tag in case of is a part-of-speech tag, and signifies whether the word is a noun, adjective, verb, and so on. reading a sentence and being able to identify what words act as nouns, pronouns, verbs, adverbs, and so on. All these are referred to as the part of speech tags.

According to Wikipedia, part-of-speech tagging (POS tagging or PoS tagging or POST), also called grammatical tagging or word-category disambiguation, is the process of marking up a word in a text as corresponding to a particular part of speech, based on both its definition and its context i.e. its relationship with adjacent and related words in a phrase, sentence, or paragraph.

Part of Speech	Tag
Noun	n
Verb	v
Adjective	a
Adverb	r

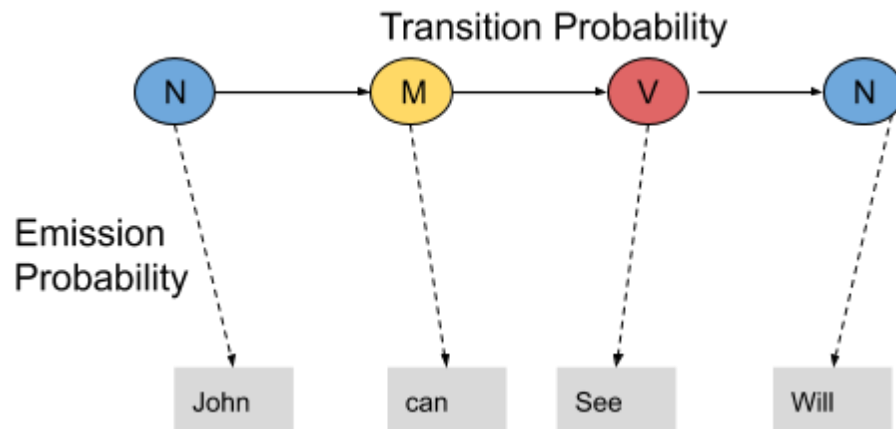
Default tagging is a basic step for the part-of-speech tagging. It is performed using the DefaultTagger class. The DefaultTagger class takes 'tag' as a single argument. NN is the tag for a singular noun. DefaultTagger is most useful when it gets to work with most common part-of-speech tag. that's why a noun tag is recommended.



POS tagging with Hidden Markov Model

HMM (Hidden Markov Model) is a Stochastic technique for POS tagging. Hidden Markov models are known for their applications to reinforcement learning and temporal pattern recognition such as speech, handwriting, gesture recognition, musical score following, partial discharges, and bioinformatics.

Let us consider an example proposed by Dr.Luis Serrano and find out how HMM selects an appropriate tag sequence for a sentence.



In this example, we consider only 3 POS tags that are noun, model and verb. Let the sentence “**Ted will spot Will**” be tagged as noun, model, verb and a noun and to calculate the probability associated with this particular sequence of tags we require their **Transition probability and Emission probability**.

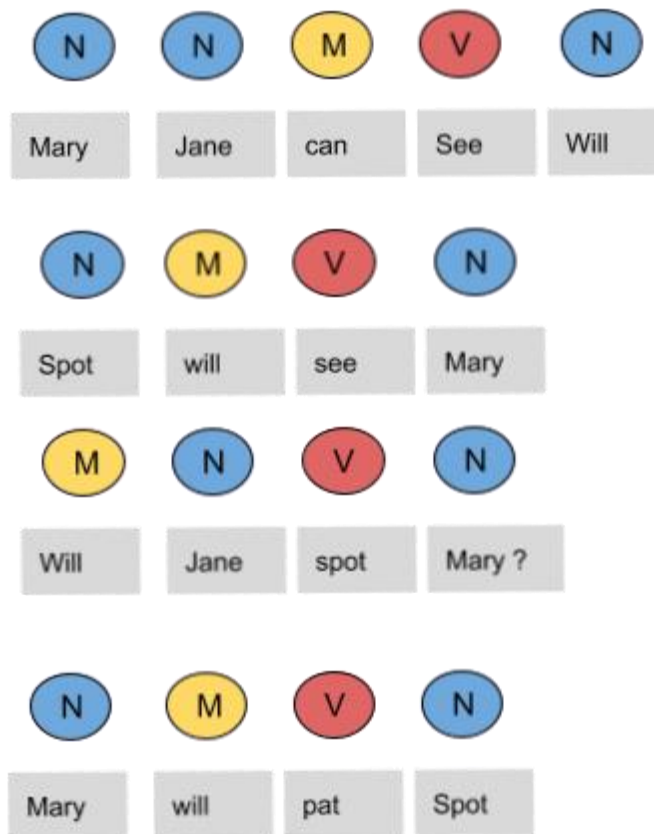
The **transition probability** is the likelihood of a particular sequence for example, how likely is that a noun is followed by a model and a model by a verb and a verb by a noun. This probability is known as Transition probability. It should be high for a particular sequence to be correct.

Now, what is the probability that the word Ted is a noun, will is a model, spot is a verb and Will is a noun. These sets of probabilities are **Emission probabilities** and should be high for our tagging to be likely.

Let us calculate the above two probabilities for the set of sentences below

- Mary Jane can see Will
- Spot will see Mary
- Will Jane spot Mary?
- Mary will pat Spot

Note that Mary Jane, Spot, and Will are all names.



In the above sentences, the word Mary appears four times as a noun. and see appears two times as a verb. we need to calculate the probability of a word appearing as noun, verb or model. to do this, we need to calculate the emission probabilities, which represented using below table.

Words	Noun	Model	Verb
Mary	4	0	0
Jane	2	0	0
Will	1	3	0
Spot	2	0	1
Can	0	1	0
See	0	0	2
pat	0	0	1

Now divide each column by the total number of their appearances .for example, ‘noun’ appears nine times in the above sentences, so divide each term by 9 in the noun column. and repeat the same for all remaining processes. We get the following table after this operation.

Words	Noun	Model	Verb
Mary	$4/9$	0	0
Jane	$2/9$	0	0
Will	$1/9$	$3/4$	0
Spot	$2/9$	0	$1/4$
Can	0	$1/4$	0
See	0	0	$2/4$
pat	0	0	1

From the above table, we can conclude that

The probability that Mary is Noun = $4/9$

The probability that Mary is Model = 0

The probability that Mary is Verb = 0

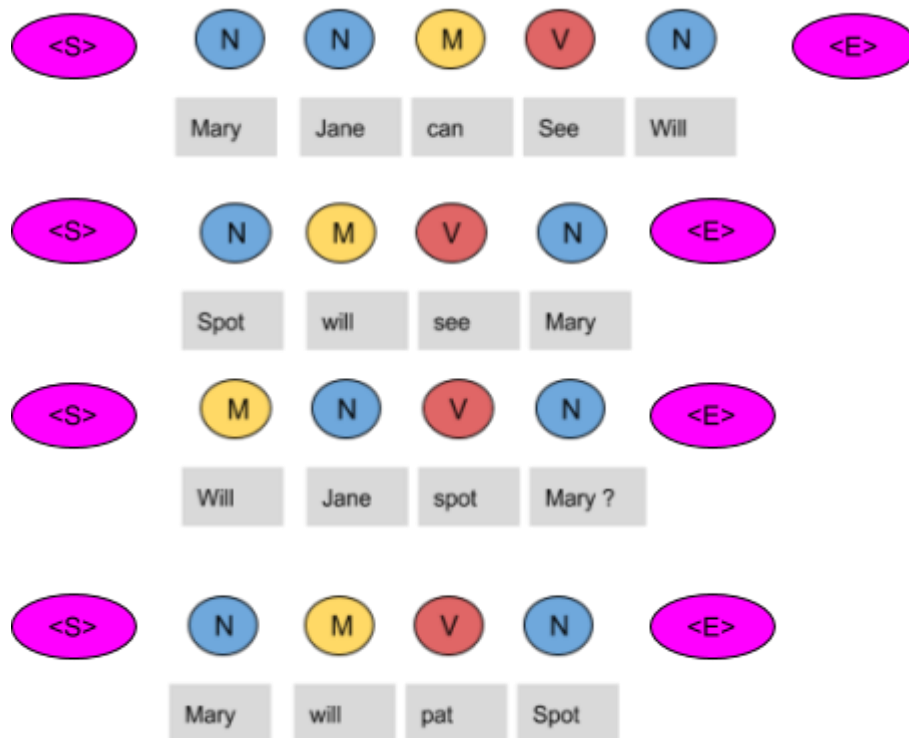
The probability that Will is Noun = $1/9$

The probability that Will is Model = $3/4$

In a similar manner, we can analyze rest of the probabilities. These are the **emission probabilities**.

Next, we have to calculate the transition probabilities, so define two more tags < S > and < E >. < S > is placed at the beginning of each sentence and < E > at the end as shown in the figure below.

since for first and last word there is no previous and next words, so we are adding extra dummy words. i.e < E > and < S >

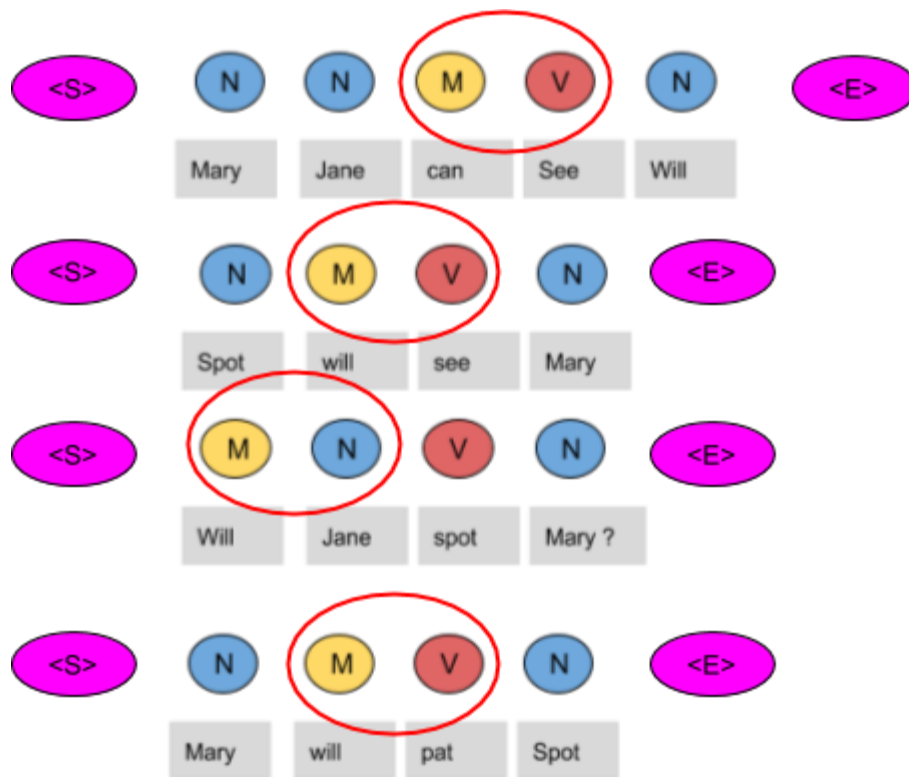


now we need to create a table and fill it with the co-occurrence counts of the tags.

	N	M	V	<E>
<S>	3	1	0	0
N	1	3	1	4
M	1	0	3	0
V	4	0	0	0

In the above figure, we can see that the < S > tag is followed by the N tag three times, thus the first entry is 3. The modal tag follows the < S > just once, thus the second entry is 1. In a similar manner, the rest of the table is filled.

Next, we divide each term in a row of the table by the total number of co-occurrences of the tag in consideration, for example, The Model tag is followed by any other tags four times (in total) as shown below, thus we divide each element in the third row by four.



the table is refined as below:

	N	M	V	<E>
<S>	$\frac{3}{4}$	$\frac{1}{4}$	0	0
N	$\frac{1}{9}$	$\frac{3}{9}$	$\frac{1}{9}$	$\frac{4}{9}$
M	$\frac{1}{4}$	0	$\frac{3}{4}$	0
V	$\frac{4}{4}$	0	0	0

These are the respective transition probabilities for the above four sentences.