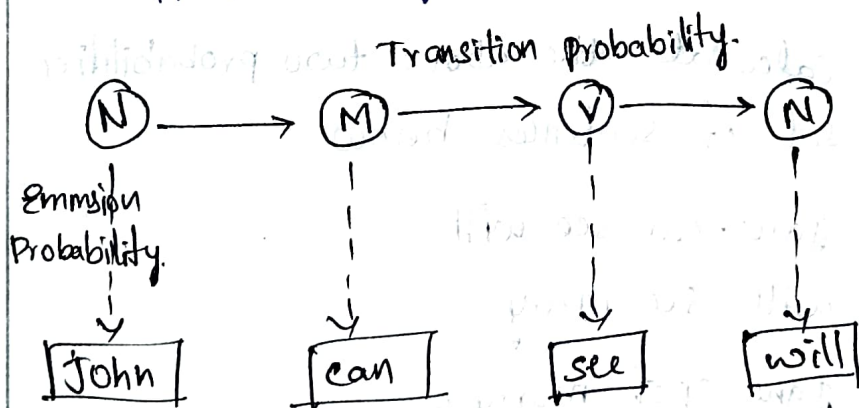4) Explain pos (parts of speech) with HMM?

Ans) HMM (Hidden Markov Model) is a stochastic technique for pos tagging

→ Hidden Markov models are known as their application to reinforcement learning and temporal pattern recognition such as speech, handwriting, gesture recognition, musical score following, partial discharge and bioinformatics.

Pos tagging with Hidden markov model

HMM (Hidden markor model) is a stachostic technique for Pos togging.

→ let us consider an example proposed by Dr. Lusserano and find out how HMM selects an appropriate tag sequence for a sentence.

Transition probability.



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, modal verb and a noun and to calculate the Probability associated with this particular Sequence of tags we require their transition Probability and Emmission probability.

→ The transition probability is the likelihood of a particular sequance for example how likely is that a noun is followed by a noun. This Probability is known as Transition Probability that the word. It should be high for a particular Sequence to be correct.

→ Now, what is the probability that the word led is a noun. These set of probabilities are emission probobabilites 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 marry?
* mary will Bat spot?

| (N) | (N) | (M) | (V) | (N) |
|-----|-----|-----|-----|-----|
| [mary] | [Jane] | [can] | [see] | [will] |

| (N) | (M) | (V) | (N) |
|-----|-----|-----|-----|
| [Spot] | [will] | [see] | [Mary?] |

| (M) | (N) | (V) | (N) |
|-----|-----|-----|-----|
| [will] | [Jane] | [spot] | [Mary?] |

| (N) | (M) | (V) | (N) |
|-----|-----|-----|-----|
| [Mary] | [will] | [pat] | [spot] |

→ The above sentences, the word mary appear four times as a noun to calculate the emission probabilities.

let us create a counting table in a similar manner;

| 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 let us divide each column by the total no. of their appearances, for example noun appears nine times in the above sentences so divide each term by a in the noun column 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 | 4/4 | 0 |
| see | 0 | 0 | 2/4 |
| Pat | 0 | 0 | 1 |

from the above table we infer that
The probability that mary noun = 4/9
The probability that mary is model = 0
the probability that will is noun = 1/9
The probability that will is model = 3/4

In a similar manner, we can configure out the rest of the probabilities, there 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 bigining of each sentence and <E> at the end as shown in the figure below.

Sequence 1: &lt;S&gt; — N (may) — N (Jane) — M (can) — V (see) — N (will) — &lt;E&gt;

Sequence 2: &lt;S&gt; — N (spot) — M (will) — V (see) — N (mary) — &lt;E&gt;

Sequence 3: &lt;S&gt; — M (will) — N (Jane) — V (spot) — N (mary) — &lt;E&gt;

Sequence 4: &lt;S&gt; — N (mary) — M (will) — V (pat) — N (spot) — &lt;E&gt;

Let us create a table and fill it with the co-occurance counts of tags.

|       | N | M | V | &lt;E&gt; |
|-------|---|---|---|-----------|
| &lt;S&gt; | 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 &lt;S&gt; tag is followed by the N tag three times, thus the first entry is 3. The model tag follows the &lt;S&gt; just once, thus the second entry is 1. In a similar manner, the rest of the table is filled.

→ Next we divide each item term in a row of the table by the total no. of co-occurence of the tag in consideration.

for example, The model tag is followed by any other tag four times as shown below, thus we divide each element in the third row by four.

⟨S⟩　Ⓝ　Ⓜ　Ⓜ Ⓥ　Ⓝ　⟨E⟩

　　mary　Jane　Car　See　will

⟨S⟩　Ⓝ　Ⓜ Ⓥ　Ⓝ　⟨E⟩

　　Spot　will　See　mary

⟨S⟩　Ⓜ　Ⓝ Ⓥ　Ⓝ　⟨E⟩

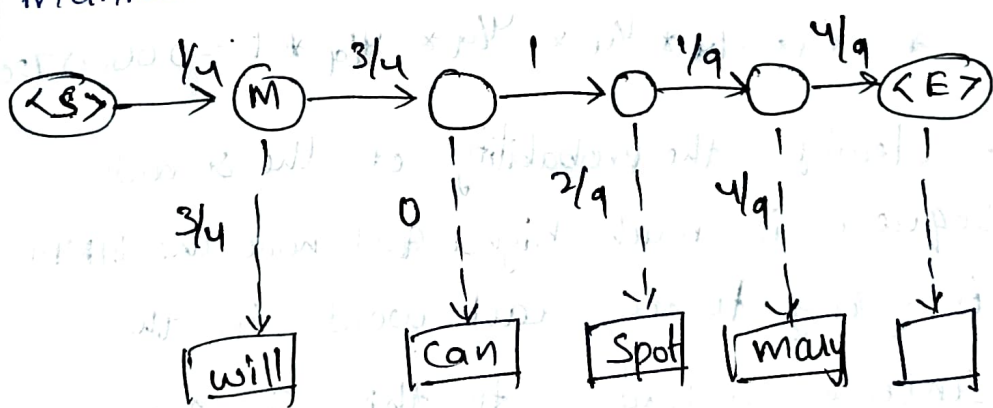　　will　Jane　Spot　mary

⟨S⟩　Ⓝ　Ⓜ Ⓥ　Ⓝ　⟨E⟩

　　mary　will　pat　spot

| | N | M | V | ⟨E⟩ |
|---|---|---|---|---|
| ⟨S⟩ | 3/4 | 1/4 | 0 | 0 |
| N | 1/9 | 3/9 | 1/9 | 4/9 |
| M | 1/4 | 0 | 3/4 | 0 |
| V | 4/4 | 0 | 0 | 0 |

→ These are the respective transition probabilities for the above four sentences. Now how does the HMM determine the appropriate sequences of tags for a particular sequence from the above tables ?
let us find it out

→ Take a new sentence and tag them with wrong tags, let the sentences
I will can spot mary' be tagged as

- will as a model
- can as a verb
- spot as a noun
- mary as a noun

Now we calculate the probability of this sequence being correct in the following manner.



→ The probability of the tag model (M) comes after the tag <S> is ¼ as seen in the table, Also the probability that the word will is a model is 3/4

→ Since the tags are not correct, the
Product is zero.

$$\frac{1}{4} * \frac{3}{4} * \frac{3}{4} * 0 * 1 * \frac{2}{a} * \frac{1}{a} * \frac{4}{9} * \frac{4}{9} = 0$$

when these words are correctly tagged
we get a probability greater than zero
as shown below.
Calculating the product of these terms
we get

$$\frac{3}{4} * \frac{1}{a} * \frac{3}{a} * \frac{1}{u} * \frac{3}{u} * \frac{1}{u} * 1 * \frac{4}{a} * \frac{4}{9} = 0.0002572 0164$$

$$<S> \rightarrow N \rightarrow M \rightarrow N \rightarrow N \rightarrow <E>$$

$$= \frac{3}{4} * \frac{1}{a} * \frac{3}{a} * \frac{1}{u} * \frac{4}{u} * \frac{2}{9} * \frac{1}{9} * \frac{1}{9} * \frac{4}{9}$$

$$= 0.000000846754$$

$$<S> \rightarrow N \rightarrow M \rightarrow N \rightarrow V \rightarrow <E> = \frac{3}{u} * \frac{1}{a} * \frac{3}{9}$$

$$* \frac{1}{u} * \frac{3}{u} * \frac{1}{u} *, \frac{4}{a} * \frac{4}{9} * 1 = 0.00025720164$$

→ clearly, the probability of the second
sequence is much higher and hence the HMM
is going to tag each word in the
Sentence according to this sequence.