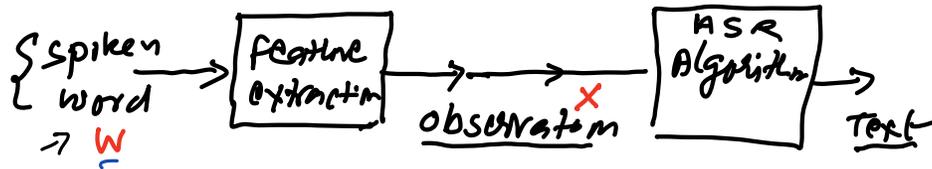


ASR - Automatic Speech Recognition.



The algorithm estimates (guesses, or computes based on probability) the spoken word.

What is probability of W given that we have observed the sequence X

$$P(W|X) = \frac{P(X|W)P(W)}{P(X)} \quad \text{--- (1)}$$

the word can be w_1, w_2, \dots, w_n . So we will select the word w_i for which the probability in RHS is maximum.

$$\underline{w^*} = \underset{w}{\text{argmax}} P(\underline{X}|\underline{w}) P(w) \quad \text{--- (2)}$$

w^* is value of w for which (2) is maximum.

Now we need to compute a) $P(\underline{X}|\underline{w})$: acoustic model

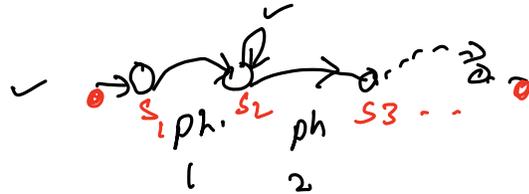
b) $P(w)$: language model

3) Find $P(x|w)$? i.e. probability of Feature vector x given that word is w ?

say, $w = \text{"Look"}$ (No. of phonemes = 2)



phoneme 1, phoneme 2, there are transitions.



$P(x|w)$ is probability

of transitions from state s_1 to s_2 , to s_3 , ... s_i

$$= \underbrace{\text{Prob. of } s_1 \text{ to } s_2} \times \underbrace{\text{Prob. of } s_2 \text{ to } s_3}$$

$$P(x|w) = \text{argmax}_{s_i} \prod P(x|s_i) * P(s_i)$$



In this case there is associated probability with each transition, and the probability of each sequence (of states) is different. Select that sequence of states which has max. probability

Let no. of states = N

and no. of observations = X

↓

$O(N^X)$

↑ complexity ✓

↑ exponential

$N = \{0, 1\}$
 If states are 0 and 1, and observations are 3

$\begin{matrix} 0 & 1 & 0 \\ 0 & 1 & 0 \end{matrix}$

∴ possible sequences
 $= 2^3$
 $= 8$

∴ we limit the no. of observations to reduce the complexity, e.g. a word spoken may be look, Look, look, ~~look~~, ~~look~~, ~~look~~
 $X = \{look, Look, look\}$

let the states are $s_0, s_1, \dots, s_k, s_{k+1}$



$$P(\underline{s_{k+1}} | \underline{s_0}, \underline{s_1}, \underline{s_2}, \dots, \underline{s_k}) = P(\underline{s_{k+1}} | \underline{s_k})$$

- ↑
- it simplifies the computation
- it is used in Markov model of speech recognition

