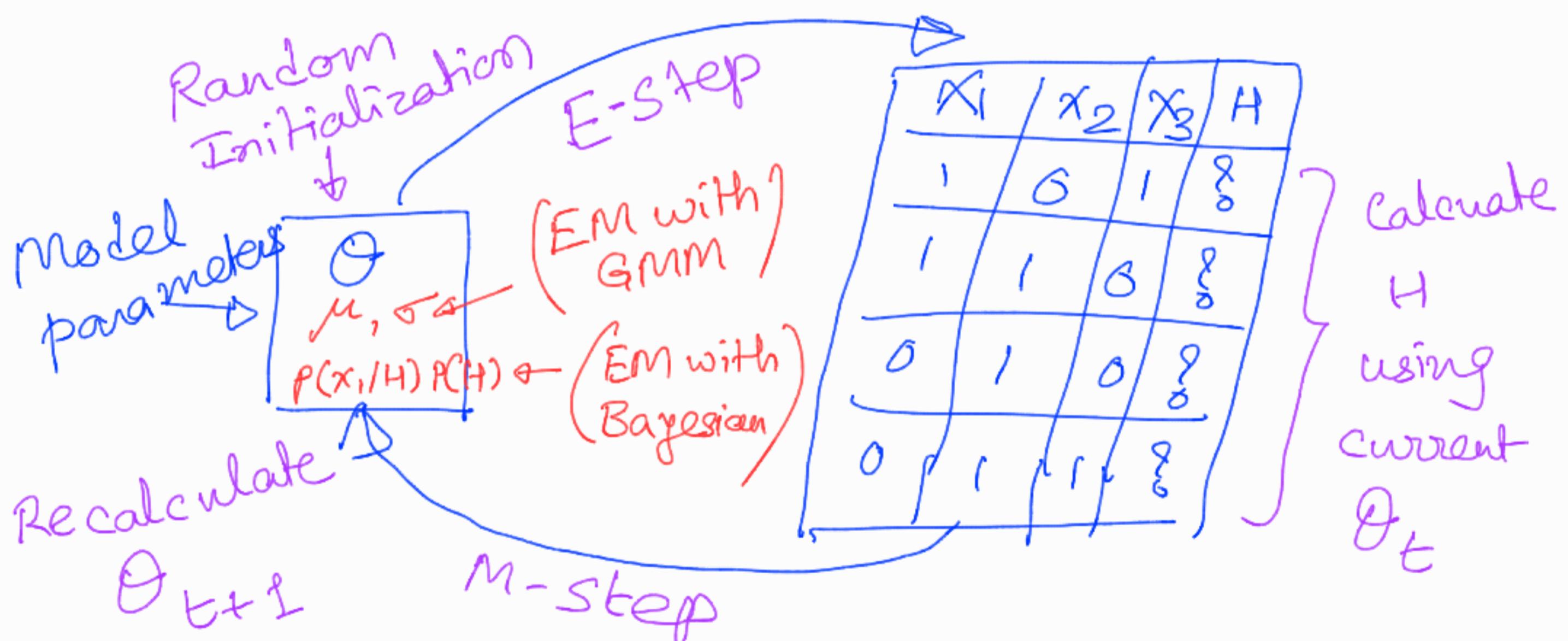


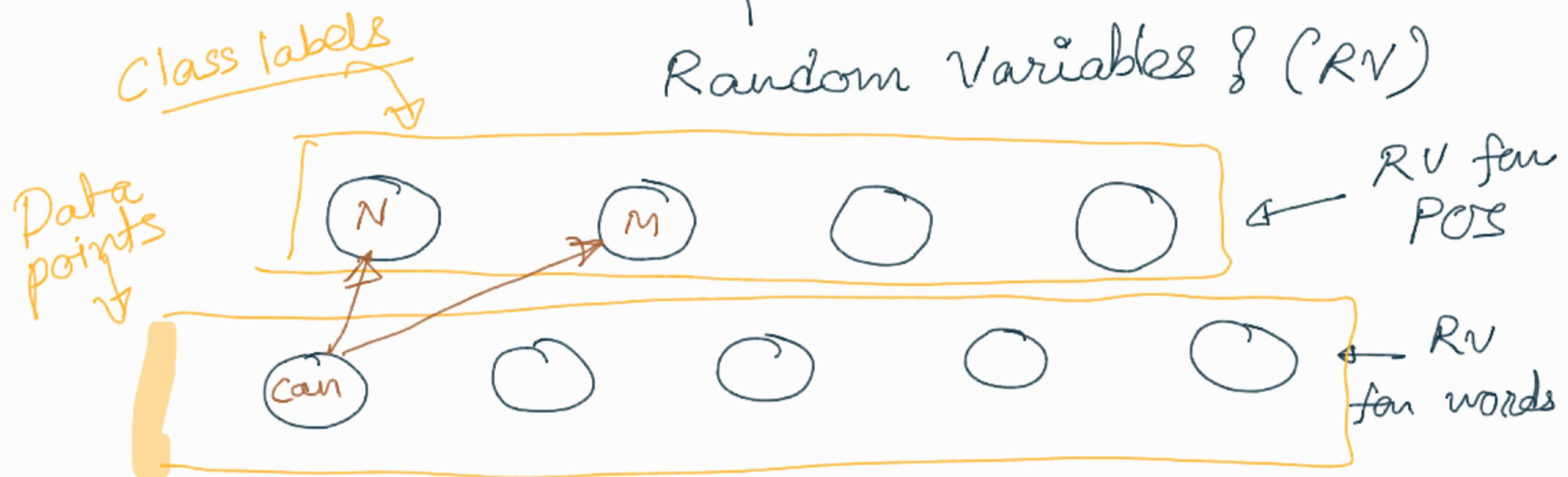
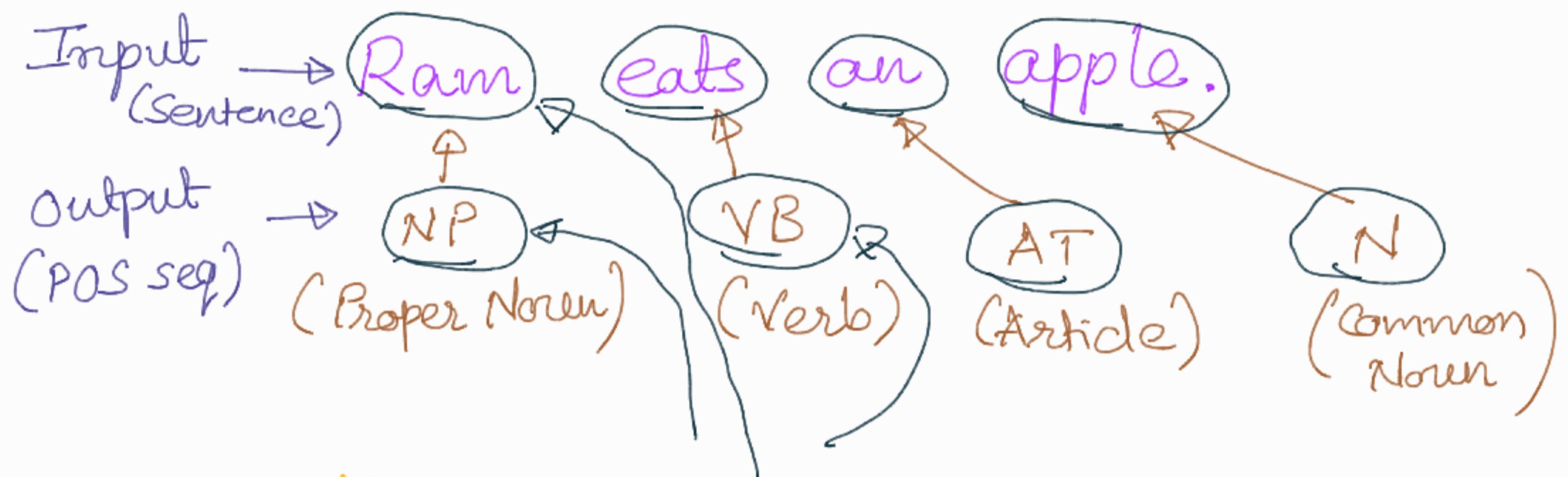
- Specialized Graphical Model
- Application of EM & Dynamic programming
- Part-of-Speech tagging
 - (Sequence Labelling task)

EM Algorithm

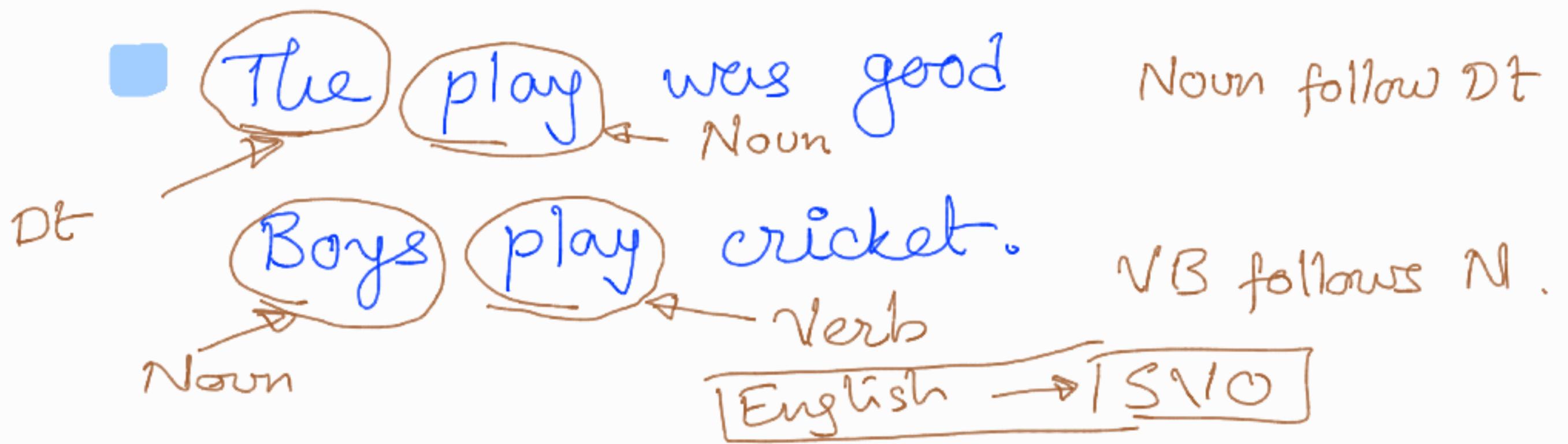


Hidden Markov Model (HMM)

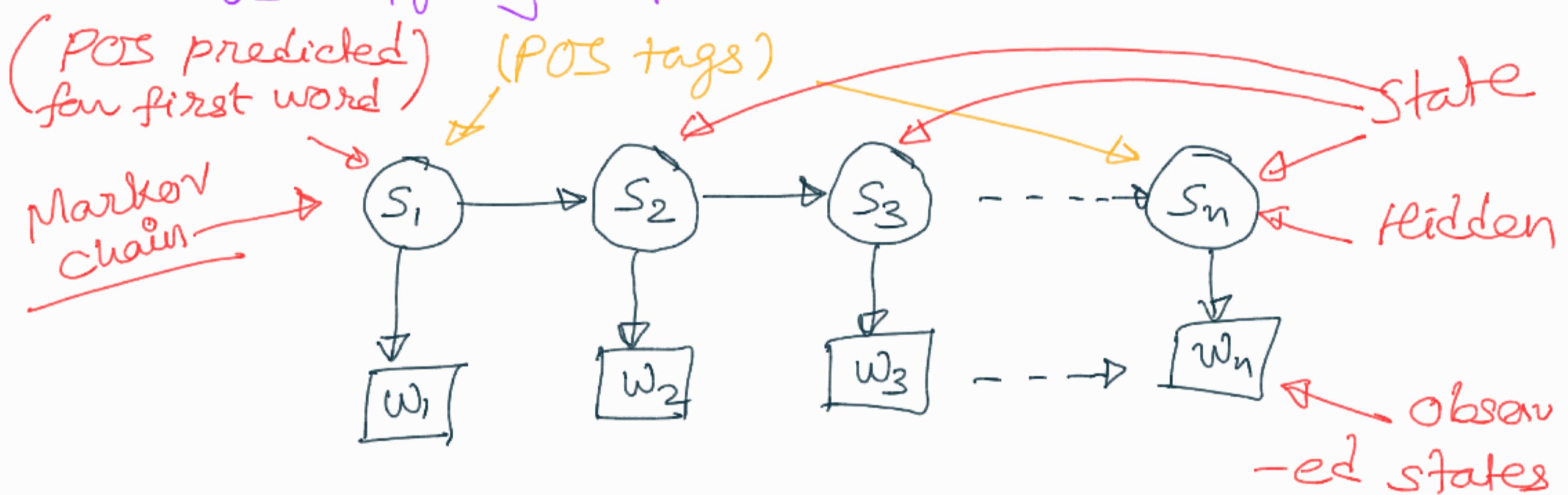
Part-of-Speech Tagging



- ① Set of RV for POS is finite
(AT, N, V, Adj, Adv, pre...)
- Set of word is infinite.
- ② POS as class labels, words as data points.
→ A word can be in more than one POS, will (N, M, auxverb)
can (N,
- ③ Given a sentence can we disambiguate POS of a word?



④ Tag of previous word helps in identifying POS of the current word



It only depends on (immediate) previous state, S_n depends on S_{n-1} only
(Markov property)

$P(S_{i+1} | S_i)$ ← Transition Probability

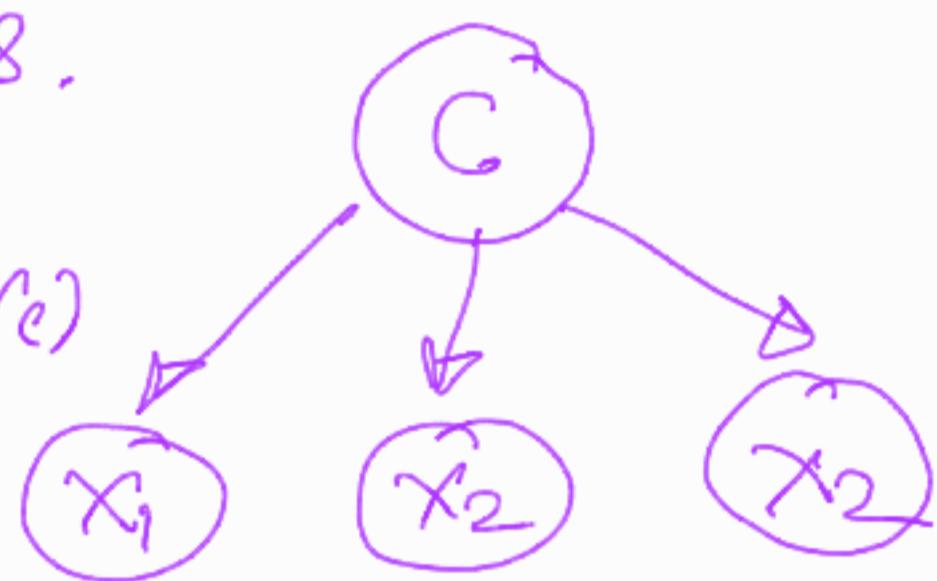
$P(w_i | S_i)$ ← Emission probability

$$P(S_i, w_1, w_2 \dots w_{i-1}, S_1, S_2 \dots S_{i-1})$$

$$\approx \left[\prod P(w_i | S_i) \cdot P(S_i | S_{i-1}) \right] P(S_0)$$

⑤ - w_s are observed variables.
 s are hidden variables.

$$P(x_1, x_n, c) = P(x_1|c) \cdot P(x_2|c) \dots P(x_n|c) \cdot P(c)$$

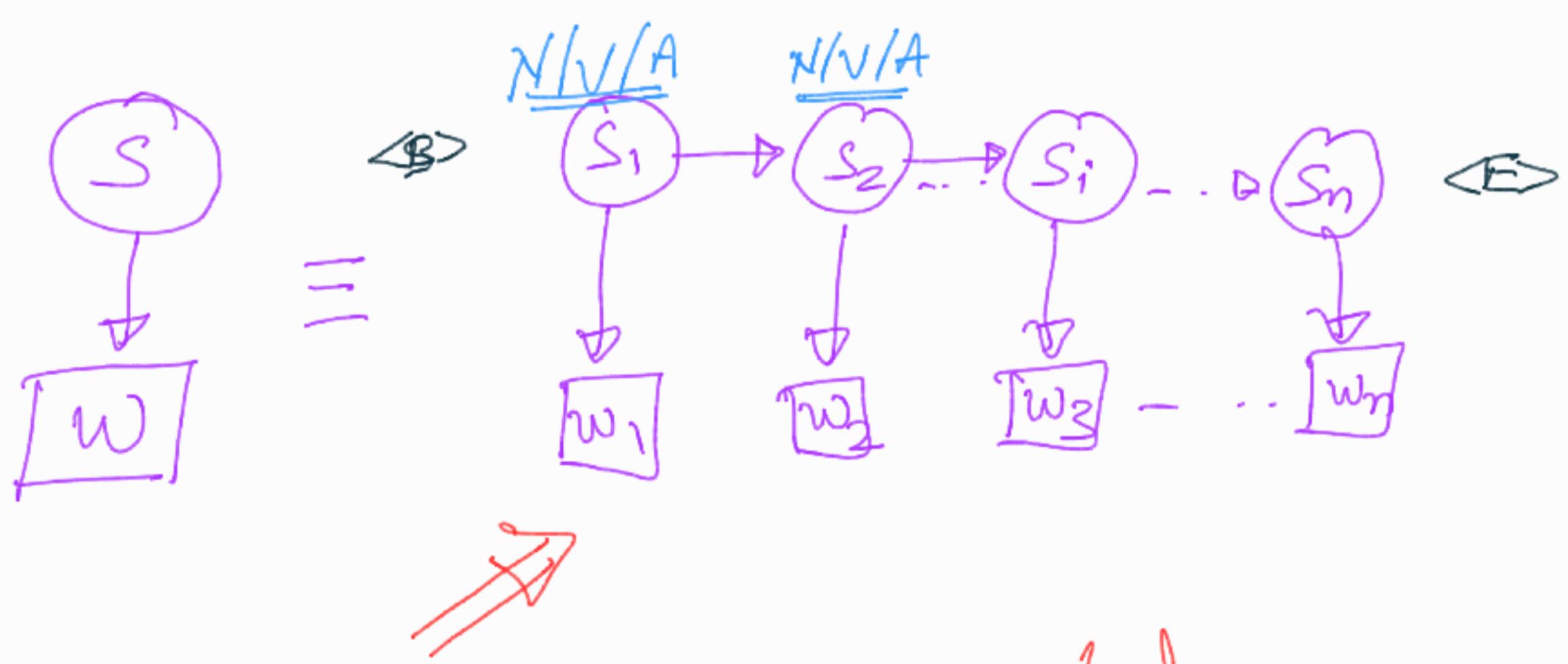


⑥ s are states, (Different from traditional RVs)

$$S \rightarrow \{N, VB, ADJ \dots\}$$

$$w \rightarrow \{a, cat, dog \dots\}$$

⑦ n (number of states) = No. of words
 in the sentence



$$\pi(s_0) = p(s_0) \xleftarrow{\text{Initial state}}$$

Transition Pro: $p(s_i | s_{i-1})$

Emission Pro: $p(w_i | s_i)$

Supervised Learning:

Text corpus with every tagged.

Ex:

1. The-DT boys-N Play-VB Cricket-NP
2. Ram-NP likes-VB the-DT play-N

Transition Prob.

$$p(s_{i+1} | s_i)$$

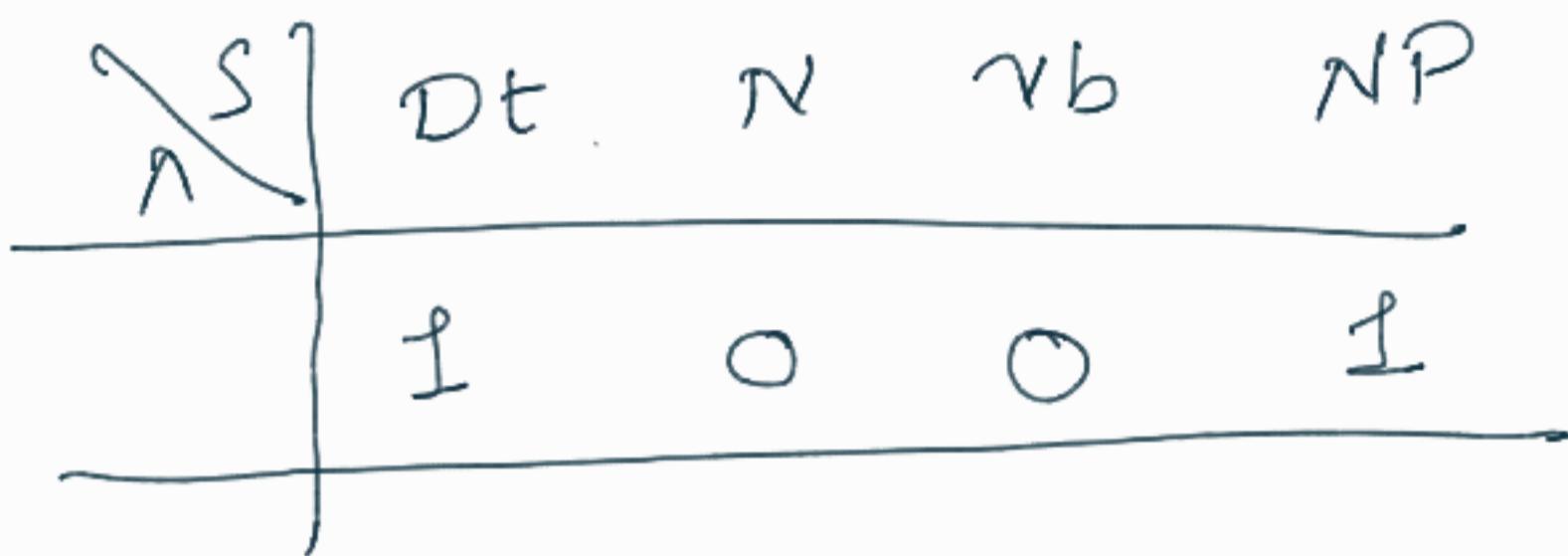
s_i	s_{i+1}	DT	N	VB	NP
DT	0 1	2	0	0	0
N	0	0	1	0	0
VB	1	0	0	1	0
NP	0	0	1	0	0

Emission Prob

$$p(w_i | s_i)$$

	T	b	P	C	R	L
	h	o	g	i	a	e
DT	2	0	0	0	0	0
N	0	1	1	0	0	0
VB	0	0	1	0	0	1
NP	0	0	0	1	1	0

$P(S_0)$



Brown Corpus

Inference :

Input: Ram → ? Plays → ? Cricket → ?

α_{i-1}
Ram =

α_i Plays

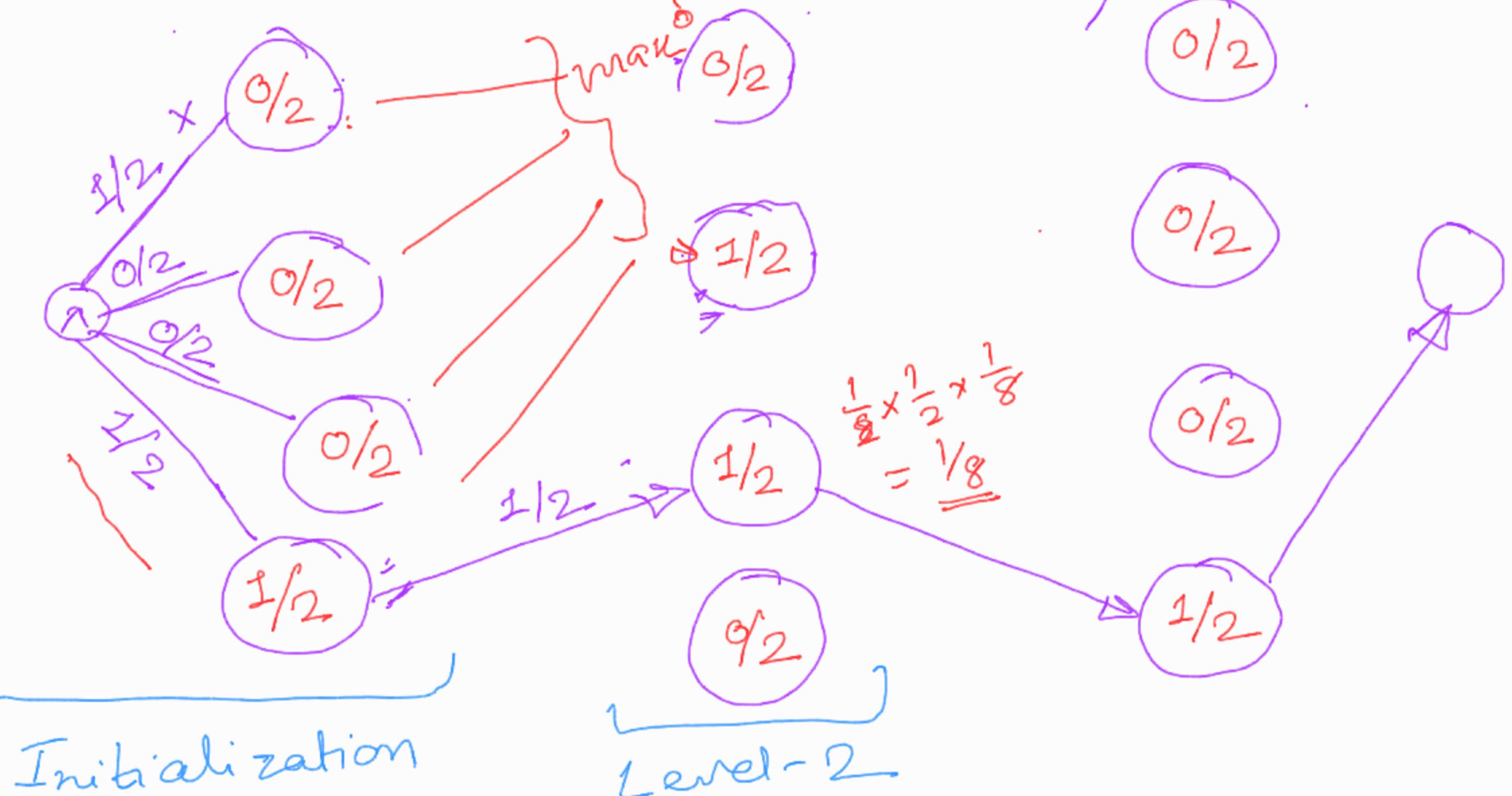
Cricket

DT

N

vb

NP



Dynamic Programming :

Viterbi Algorithm (Decoding)

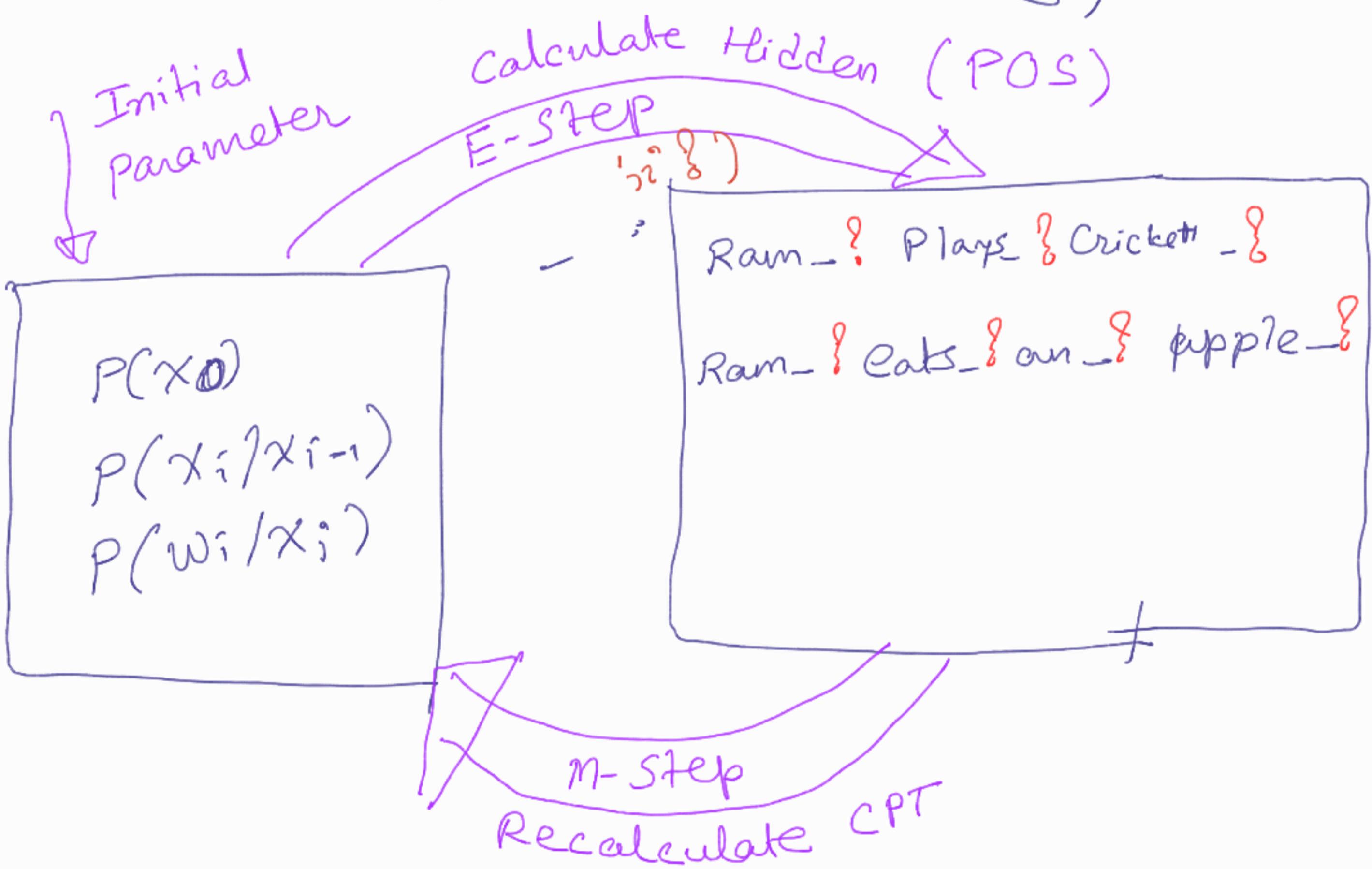
Outline

- Finish HMM and discussion on probabilistic graphical model
- Begin with discriminative classifier
- Neural Network basics
 - Perceptron

"Part-of-Speech is not known"

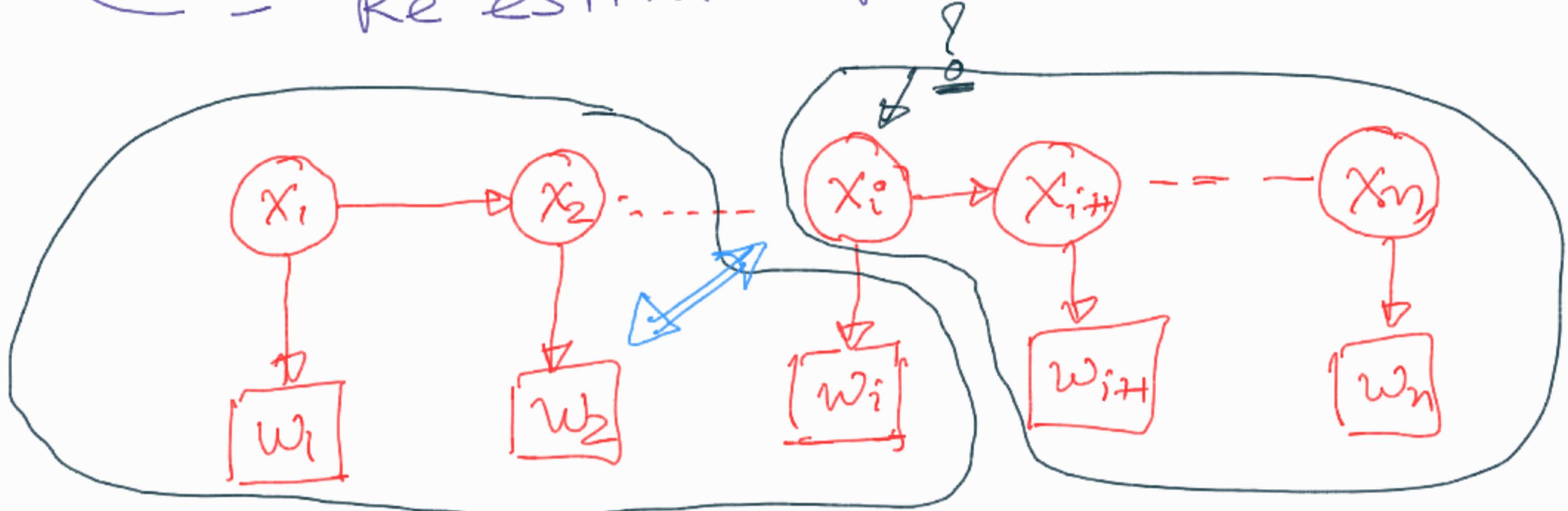
unlabelled data

(Unsupervised Learning)



Baum-welch Algorithm

- Fine tuning parameters
- Special case of EM
- start with initial parameter
- Based on parameter calculate hidden variable
- Re estimate parameter



When deciding x_i^*

- All observed words w_1, \dots, w_n
- All previous labels x_1, \dots, x_{i-1}

$$P(x_i^* | x_1, \dots, x_{i-1}, w_1, \dots, w_n)$$

$$\alpha = P(w_1, \dots, w_i, x_i^*) \Leftarrow \text{observed before deciding } x_i^*$$

$$\beta = P(w_{i+1}, \dots, w_n | x_i^*) \Leftarrow \text{All the forth coming observations}$$

$\alpha = x_i$ - depends on
 w_i, w_{i-1}, \dots, w_1

$\approx P(x_i/x_{i-1}) \cdot P(w_i/x_i)$
x (whatever value up to)
 x_{i-1}

Forward calculation

$$\alpha(i) = \sum \alpha(i-1) \cdot P(x_i/x_{i-1})$$

$\circ P(w_i/x_i)$

Forward

sum over all possible
labels at $\alpha(i-1)$

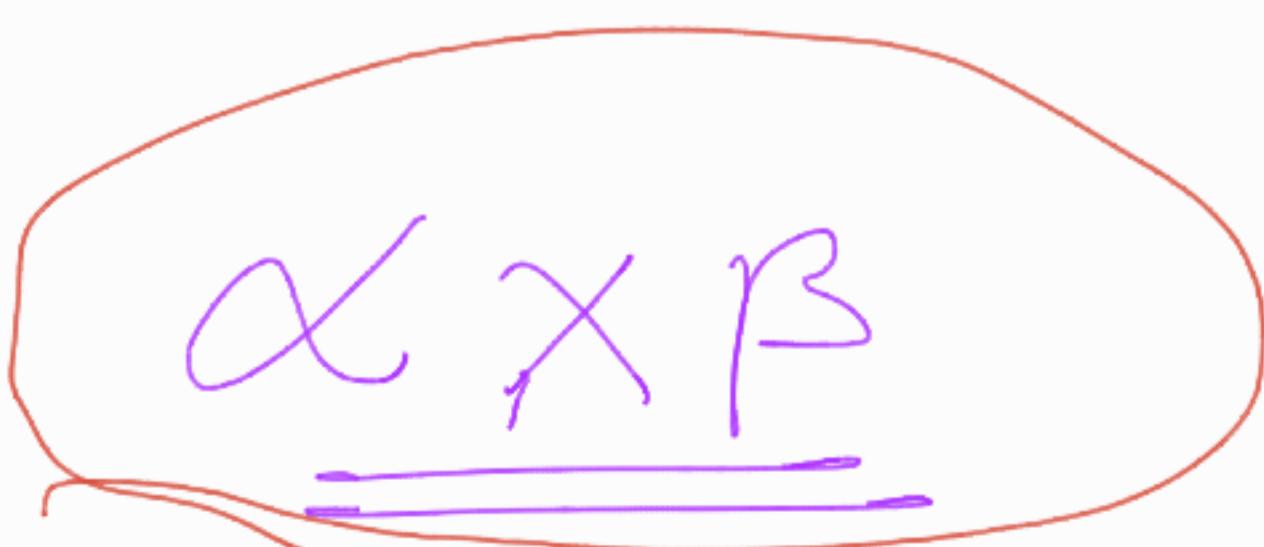
Backward computation

$$\beta(i) = \sum \beta(i+1) \cdot P(x_{i+1}/x_i)$$

$\circ P(w_{i+1}/x_{i+1})$

Backward

complete
distri \approx



Forward- Backward Algo.

$$\boxed{\beta_n = 1}$$

HMM

↳ Application to POS

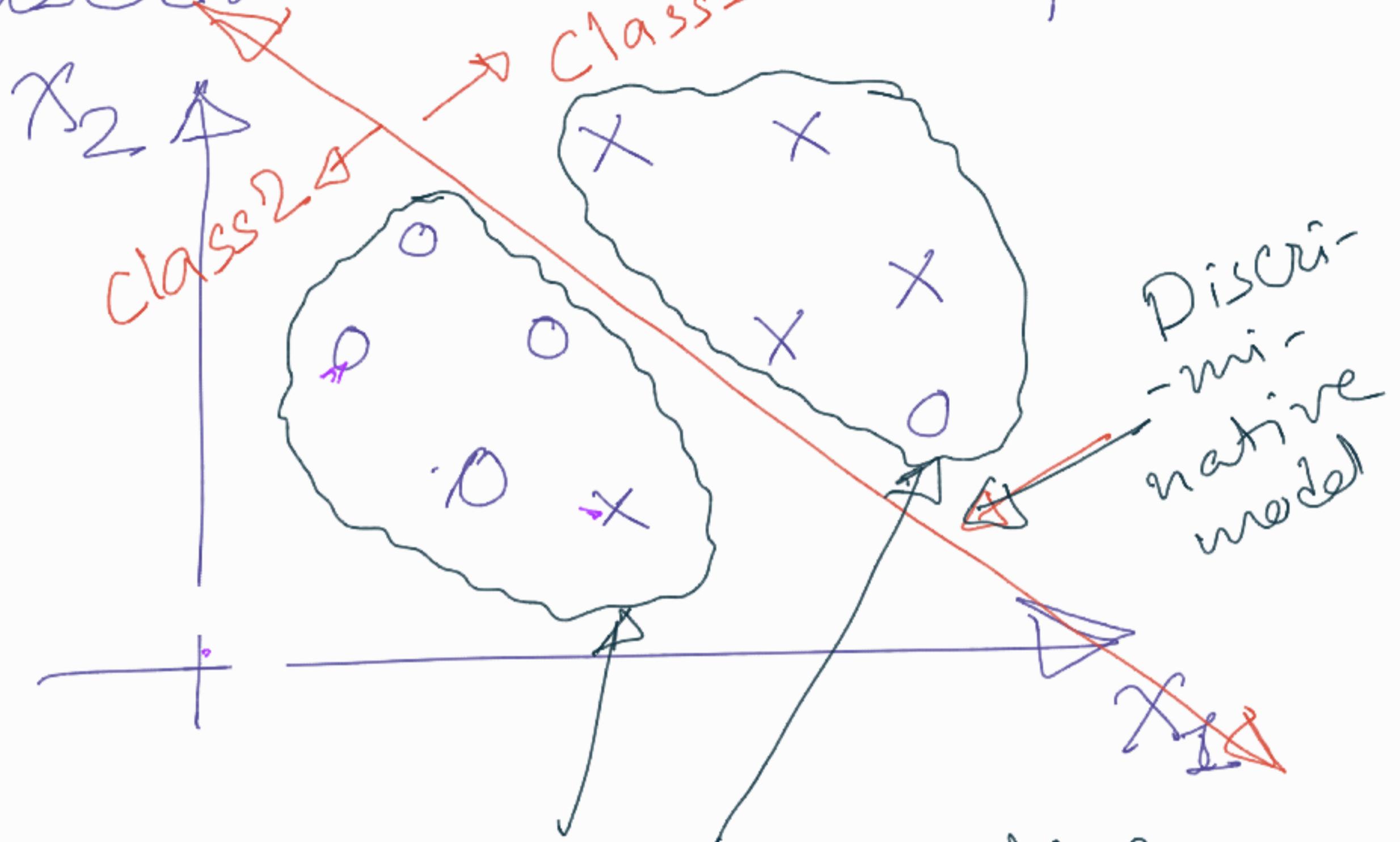
created interest in
application of ML
to NLP.

↳ Speech processing
/ Recognition

↳ Sequence labeling

- "Every feature is a
vector in N -dimensional
space."

- Discriminative classifiers



Generative \rightarrow Learn distribution of data points

Neural Network

\hookrightarrow Perceptron