06 | Jan | 2023.

.
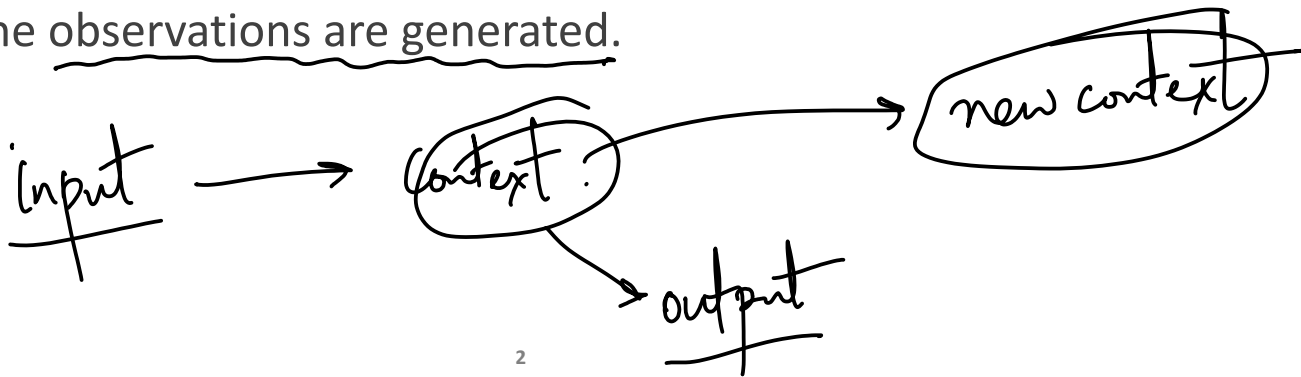
# Hidden Markov Models

# Markov Model

- Let's say there is a process which generates a sequence:

  - Sequence of written words    *text*

  - Sequence of spoken words (speech)

  - DNA, RNA protein's amino acid sequence, etc.

- The observed sequence is governed by some context which itself keeps updating as the observations are generated.

*input* → *context* → *output* → *new context*

$$w_1 \; w_2 \; \cdots \quad w_n$$

$P(obs \mid past\ history)$

- The probability of a sequence can be written as

$$P(w_{1,n}) = P(w_1)P(w_2 \mid w_1)p(w_3 \mid w_{1,2})\ldots P(w_n \mid w_{1,n-1})$$

2    3    n variable    → past history.

$(w_1 \; w_2 \; w_3 \cdots w_{n-1} \; \; w_n \; \; (w_{1:n}) \; (w_{1,n})$ entire sequence with indices 1 to n

- This is the chain rule of probability.

- Modeling the sequence of words $P(w_{1,n})$ requires learning each of the

component terms.

$P(w_3 \mid w_1, w_2)$

| $w_1$ | $w_2$ | $P(w_3)$ |
|---|---|---|
| — | — | |
| — | — | |

- We can do something smarter by modelling the context in a better way.

$P(w_{1,n})$    We want to model. $P(w_1) \; P(w_2 \mid w_1)$

CPT    CPT needs to be learned

$P(w_1 \ w_2 \ \cdots \ w_n)$ parameterized CPT

Valid sentence $w_1 \cdots w_2$      Prob high.

Correct

$\downarrow$

has errors      Prob $\downarrow$.

Language Model

$w_{1,n} \longrightarrow$ Conforms to a language
well formed sentence

- The context sequence can be

  - Observable

  - Hidden

- An example of observable context model is the n-gram model of text.

- n-Gram model: Only the previous n-1 words have any effect on the probabilities of the next word.

Context $\longrightarrow$ observation

$$P\left(W_K \Big| W_{1:K-1}\right)$$

$$= \frac{}{context}$$

HMMs

Limiting the context

unigram $P(W$

trigram $P\left(W_3 \Big| W_{1,2}\right)$

$P\left(W_5 \Big| W_{34}\right)$

$\dfrac{W_1}{}\Big| \underline{n-1}$ previous words

Bigram $\left(P(W_2 | W_1)\right)$

$P\left(W_1 | W_0\right)$

4

- For a trigram model:

$$p(w_n \mid w_1, \ldots, w_{n-1}) = p(w_n \mid w_{n-2}, w_{n-1})$$

*Past history* *last two words.*

*Due to limited context (assumption)*

- Applying the trigram model gives:

*LHS*

*Joint distribution of $w_{1,n}$*

1  2  3 . . . . . 3 variables.

$$P(w_{1,n}) = P(w_1)P(w_2 \mid w_1)p(w_3 \mid w_{1,2})\ldots P(w_n \mid w_{n-2,n-1})$$

*CPTs*

$$= P(w_1)P(w_2 \mid w_1) \quad \prod_{i=3}^{n} P(w_i \mid w_{i-2}, w_{i-1})$$

$w_0$
$w_{-1}$

$$\prod_{i=1}^{n} P($$

Corpus.

- To create such a trigram model, we require the probability of every possible valid trigram.

$$\to C(w_{i-2} \; w_{i-1} \; w_i)$$

$$C(w_{i-2} \; w_{i-1})$$

- The estimate of the conditional probability

$$P_e(w_i | w_{i-2}, w_{i-1}) = \frac{C(w_{i-2,i})}{C(w_{i-2,i-1})}$$

Training dataset

e: estimate

Count

| $w_1$ | $w_2$ | $P(w_3)$ | $P(w_3=)$ --
|---|---|---|
| $=$ | $=$ | 0 | 
| $=$ | $=$ | 0 | 
| $=$ | $=$ | 0 | 
|  |  | 0 | 

- "To create such a model we simply go through some training text"

Corpus.

$$P_e(\text{'model'} | \text{'such'}, \text{'a'}) = \frac{C(\text{'such' 'a' 'model'})}{C(\text{'such' 'a'})} = \frac{1}{1}$$

6

- After all the hard work, we now have a model which can generate the probability of any string. $P(w_{1,n})$

- This model is called a Markov chain. → Context is observable.

- It can be represented as a graph of nodes and arcs

- Nodes are the states → Context

- Arcs represent transitions from one state to another.

  - An arc is labelled with the observation (emitted symbol) and the probability of the observed symbol given the state
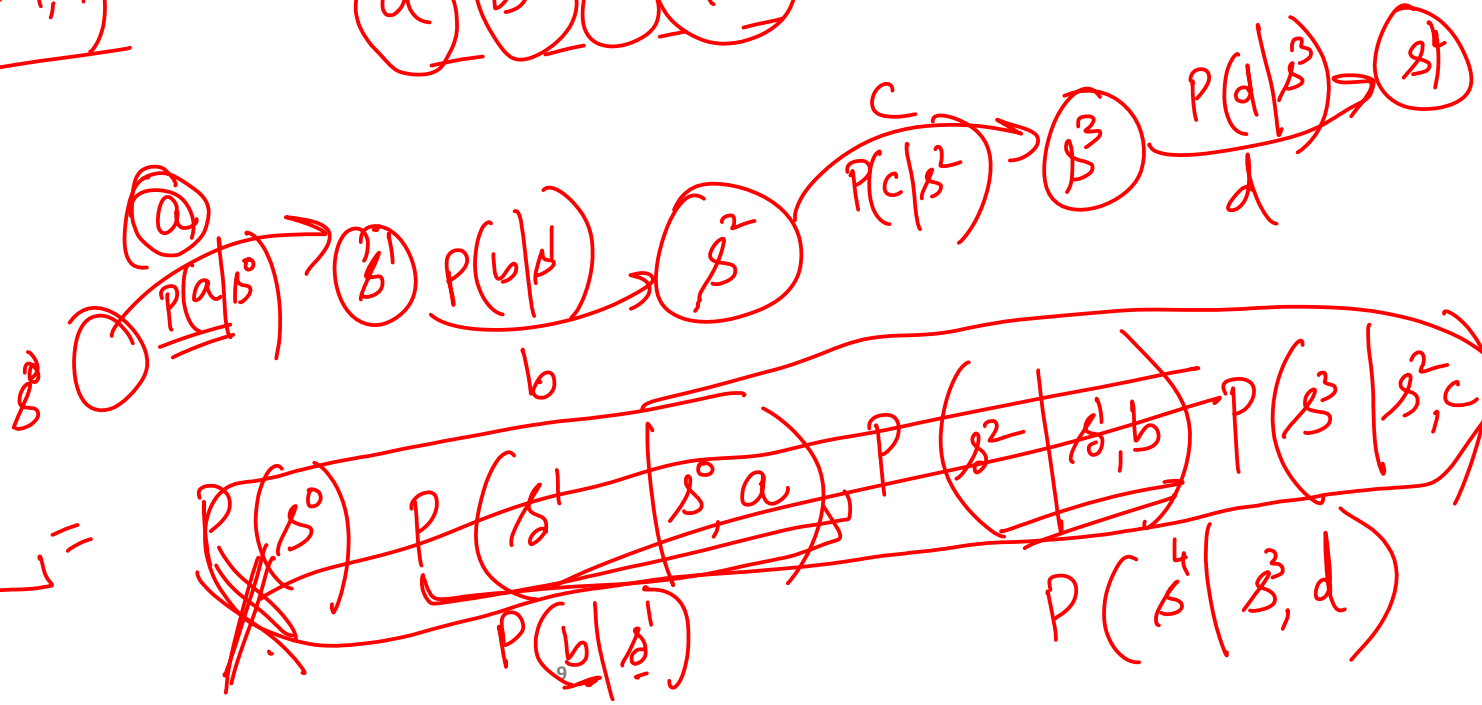
- Why is the state obvious here?

- Given a Markov chain, the probability of an observed sequence can be computed as the product of the probabilities of each transition in the path.

$$P(w_{1,n})$$

(a) (b) (c) (d)

$$s^0 \xrightarrow{a} s^1$$

$$P(s^1 \mid s^0, a) = 1$$

$$P(a \mid s^0)$$

(a) $\xrightarrow{P(a \mid s^0)}$ (s^1) $\xrightarrow{P(b \mid s^1)}$ (s^2) $\xrightarrow{P(c \mid s^2)}$ (s^3) $\xrightarrow{P(d \mid s^3)}$ (s^4)

$$s^0$$

b

c

d

$$P(abcd) = P(s^0) \, P(s^1 \mid s^0, a) \, P(s^2 \mid s^1, b) \, P(s^3 \mid s^2, c)$$

$$P(b \mid s^1)$$

$$P(s^4 \mid s^3, d)$$
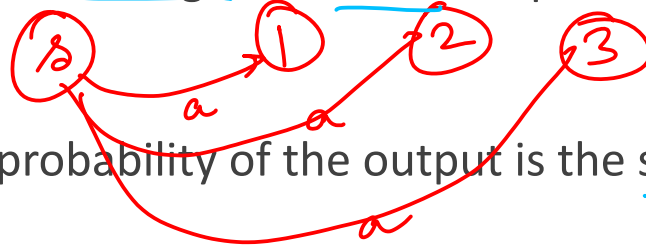
# What if the states are hidden?

- Now consider the scenario when the states are non-obvious.

- That means, we are uncertain about the state sequence that generated the given observation sequence.

- That means,

  - there can be multiple state sequences that can lead to a given observation sequence.

  - a given state sequence can generate multiple observation sequences

  - for a given state, we can make transition to multiple possible states, even while generating the same output.

# Hidden Markov Models

- So now, transition to multiple states is possible from a given state, even while emitting the same observation (symbol).

  *|| Difference from a Markov chain.*

- HMMs are a generalisation of Markov Chains in which a given state may have several transitions out of it, all with the same label (symbol).

  *arc*

  - Recall that in a Markov chain, given a state, the next state is certain (fixed) once you observe a given symbol.



$$P\left( s^{k}_{j} \mid s^{i}_{l}, a \right).$$

- Since more than one successor state may have the same output, in general there may be several paths through an HMM that produce the same output. *sequence.*

- In such cases, the probability of the output is the sum of the probabilities of all possible paths.

- The Markov chain cannot be used now.

    - What we need is a state transition diagram to represent an HMM.

# Moving from a Markov Chain to an HMM

- Consider that we are learning a trigram model.

$P(w_3 | w_1 \, w_2)$    CPT

- If the training corpus is sparse, it may happen that there is a trigram in the next text that never appeared in the training corpus.

  - In that case, a 0 probability will be assigned to the trigram, i.e. probability of the third word, given the first two words.

- One solution to this problem is to smooth the probabilities by also using the bigram and unigram probabilities.

$$P(w_n | w_{n-2,n-1}) = \lambda_1 P_e(w_n) + \lambda_2 P_e(w_n | w_{n-1}) + \lambda_3 P_e(w_n | w_{n-2,n-1})$$

$\lambda_1 + \lambda_2 + \lambda_3 = 1$

$\lambda_3 > \lambda_2 > \lambda_1$

trigram

uni    0.1    Path 1

Bi    0.3    Path 2

Trigram    0.6    Path 3

0.6  0.3  0.1

- If we have missing trigram then the bigram and unigram probabilities take over.

- If we have missing trigram and bigram, then we fall back on the unigram probability.

- The probability of a transition has become an addition of multiple components (arcs) between the two states.

  - An HMM will allow multiple paths between two states.

# HMM state transition diagram for a trigram model

- Consider that we have seen the symbol ab.

- The next symbol in the sequence is a.

- So the next state is ba

- The trigram $P(a|ab)$ itself can be represented as an HMM

$$P(a|ab) = {}_1P_e(a) + {}_2P_e(a|b) + {}_3P_e(a|ab)$$

- An HMM is a 4 tuple

$$s^1, S, W, E$$

- $S$: the set of states

- $s^1 \quad S$ is the initial state of the model

- $W$ : the set of output symbols generated/emitted/accepted

- $E$ : the set of edges or transitions.

- For each of the sets $S$, $W$, $E$ we assume a canonical ordering of the elements

$$S = \quad s^1, s^2, \ldots, s$$

$$W = \quad w^1, w^2, \ldots, w$$

$$E = \quad e^1, e^2, \ldots, e$$

- The starting state of the HMM is the first element in the ordering of states.

- A transition is a 4-tuple $s^i, s^j, w^k, p$ where $s^i$ is the state from which the

  transition starts, $s^j$ is the state where the transition ends, $w^k$ is the output

  symbol generated by the model, $p$ is the probability of taking the transition

- $\rho$ is the probability of the transition $\quad s^i \quad w^k \quad s^j$

- No two transitions can have the same starting and ending states as well as the same output value.

- We can leave out the zero probability paths from the graphical representation of the HMM.

- Note that a state $s$ can be the starting state for several transitions that have the same output symbol but go to different ending states.

- It is not possible to know what state the machine has gone into simply by looking at the output.

- Thus, the state sequence followed by an HMM is not deductible from the input. It is hidden.

- There are n+1 states for n <u>outputs</u>

- The probability $p$ associated with a transition $s^i \xrightarrow{w_k} s^j$

$$P(s^i \xrightarrow{w_k} s^j) = P(S_{t+1} = s^j, W_t = w^k | S_t = s^i)$$

$$= P(s^j, w^k | s^i)$$

- When writing $P(s^j, w^k | s^i)$, it is understood that $s^i$ is the prior state and $s^j$ is the next state.

- Markov models assume that the only information affecting the probability of an output, or of the next state, is the prior state.

- As per this Markov Assumption

$$P(w_n, s_{n+1} \mid w_{1,n-1}, s_{1,n}) = P(w_n, s_{n+1} \mid s_n) = P(s^i \xrightarrow{w^k} s^j)$$

- The probability of a sequence $w_{1,n}$ is the probability of all possible paths through the HMM that can produce this sequence.

In other words:

$$P(w_{1,n}) = \sum_{s_{1,n+1}} P(w_{1,n}, s_{1,n+1})$$

# The Markov Assumption helps

$$P(w_{1,n}) = \sum_{s_{1,n+1}} P(w_{1,n}, s_{1,n+1})$$

$$P(w_{1,n}) = \sum_{s_{1,n+1}} P(s_1)P(w_1, s_2 | s_1)P(w_2, s_3 | w_1, s_{1,2}) \ldots P(w_n, s_{n+1} | w_{1,n-1}, s_{1,n})$$

- The Markov Assumption helps here

$$P(w_{1,n}) = \sum_{s_{1,n+1}} \prod_{i=1}^{n} P(w_i, s_{i+1} | s_i)$$

# Exercise: Graphically visualize the simplification offered by Markov Assumption

# An application of HMM: Part of Speech Tagging

- English words can be in more than one PoS class.

- PoS tags can be assigned using HMM.

- But to use HMM we need to phrase the problem of assigning PoS tags to the words as one of assigning probabilities to the input text.

- We assume that the HMM generates outputs which are the words of the corpus.

- We assume that there is some connection between the states and the transitions of the tags.

- Earlier we had for Language Model (LM)

$$P(w_{1,n}) = \sum_{s_{1,n+1}} P(w_{1,n}, s_{1,n+1})$$

- For PoS tagging, we assume correspondence between states and the tags.

- So,

$$P(w_{1,n}) = \sum_{t_{1,n+1}} P(w_{1,n}, t_{1,n+1})$$

- Here, $t_{1,n+1}$ is a sequence of $n + 1$ parts of speech or tags.

- The last tag $t_{n+1}$ is the tag achieved after emitting $w_n$.

- So $t_{n+1}$ is a tag which corresponds to the non-existent word $w_{n+1}$

- We define the problem of PoS tagging as finding

$$\arg\max_{t_{1,n}} P(t_{1,n} \mid w_{1,n}) = \arg\max_{t_{1,n}} \frac{P(w_{1,n}, t_{1,n})}{P(w_{1,n})} = \arg\max_{t_{1,n}} P(w_{1,n}, t_{1,n})$$