# Hidden Markov Models

Hidden Markov Models are probabilistic models

They are used to find probabilities of sequences of events.

Based on the sequences, questions can be answered

| Day | Temperature | Icecreames eaten |
|---|---|---|
| Monday | Hot | 2 |
| Tuesday | Hot | 1 |
| Wednesday | Cold | 0 |

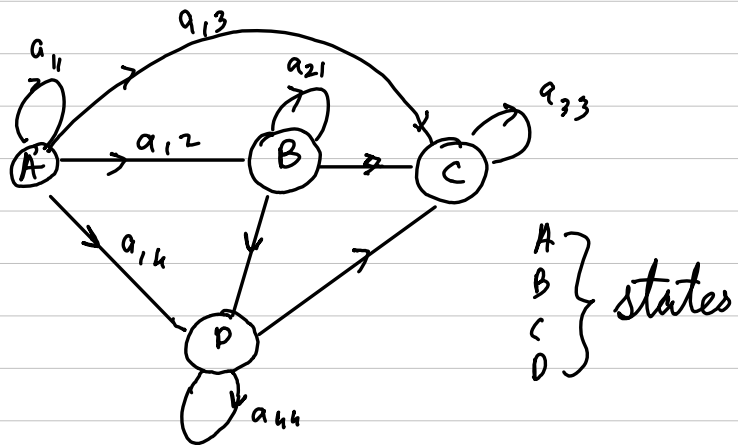what is the probability that I will eat 2 icecreames tomorrow?

What is the temperature on thursday & friday if I ate 2 icecreames on Friday?

HMMs are kind of system called Finite or Discrete Markov model.

Markov model is a finite state machine
with N distinct states begins at $t = 1$ in
Initial state

Moves from one state to another state
according to probabilities associated with
current state

Number of states are finite



$a_{ij}$ is probability of moving from
state $a_i$ to $a_j$

$$\sum_{i=1}^{N} a_{ik} = 1 \qquad \forall k$$

is sum of all outgoing arrows $= 1$

# Hidden Markov Model

HMM is a stastical model in which the system is assumed to be markov process with unobserved hidden states
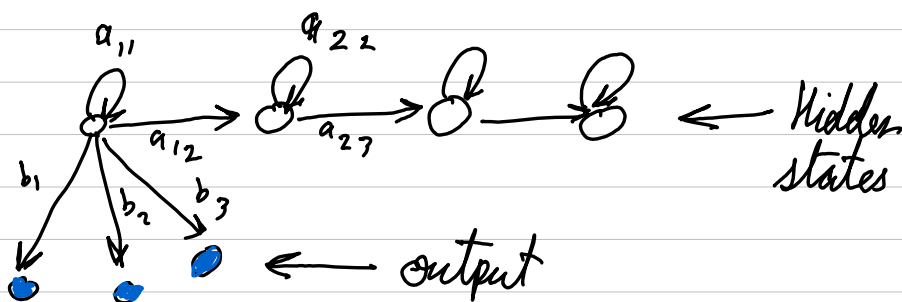
Consists of states $\dot{S}, S_2, S_3 \ldots \ldots$

$$P(S_{ij} | S_{i1}, S_{j2}, \ldots S_{ik-1}) = P(S_{ik} | S_{ik-1})$$
(Markov property)

$A \rightarrow$ Set of transition probabilities

$B \rightarrow$ set of output probabilities

$\Pi \rightarrow$ initial probabilities



$\leftarrow$ Hidden states

$\leftarrow$ output

$$b_{11} + b_{12} + b_{13} + b_{14} = 1$$
$$b_{21} + b_{22} + b_{23} + b_{24} = 1$$

output probability

$$M = \{A, B, \Pi\}$$

Markov property : The current state of system depends only on the previous state of system

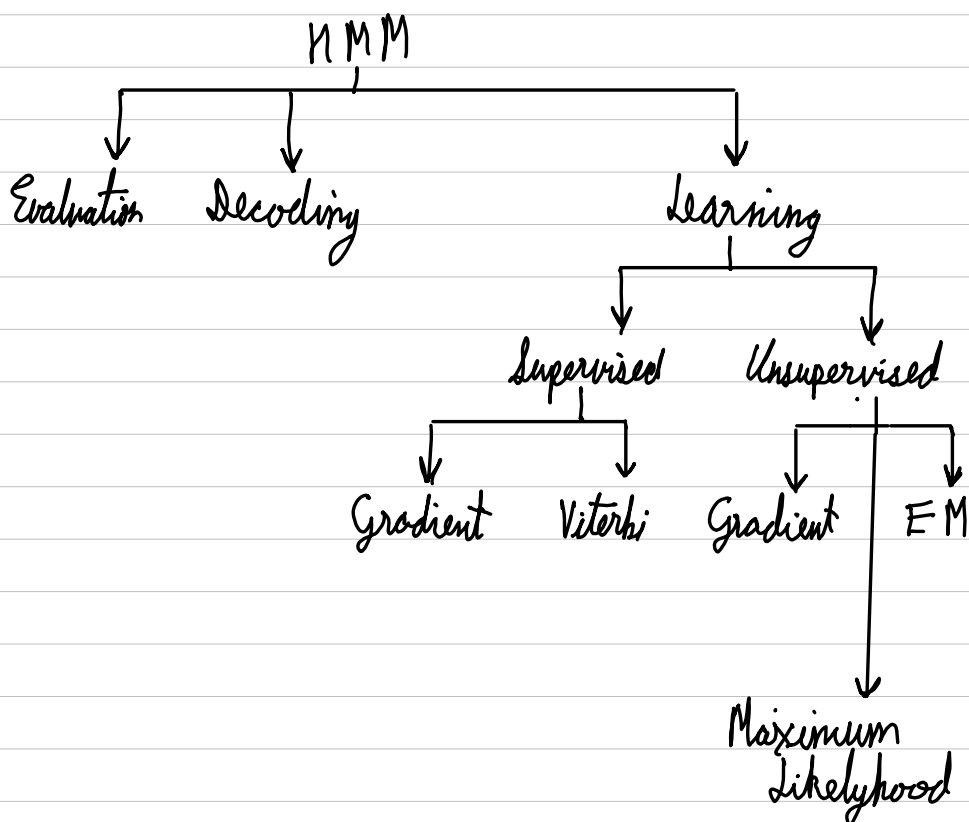state of time $T_{+1}$ depends on state at $T$

Markov models obey the Markov property

Given $H(A, B, \pi)$ & a sequence $O$
we can do solve 3 problems →

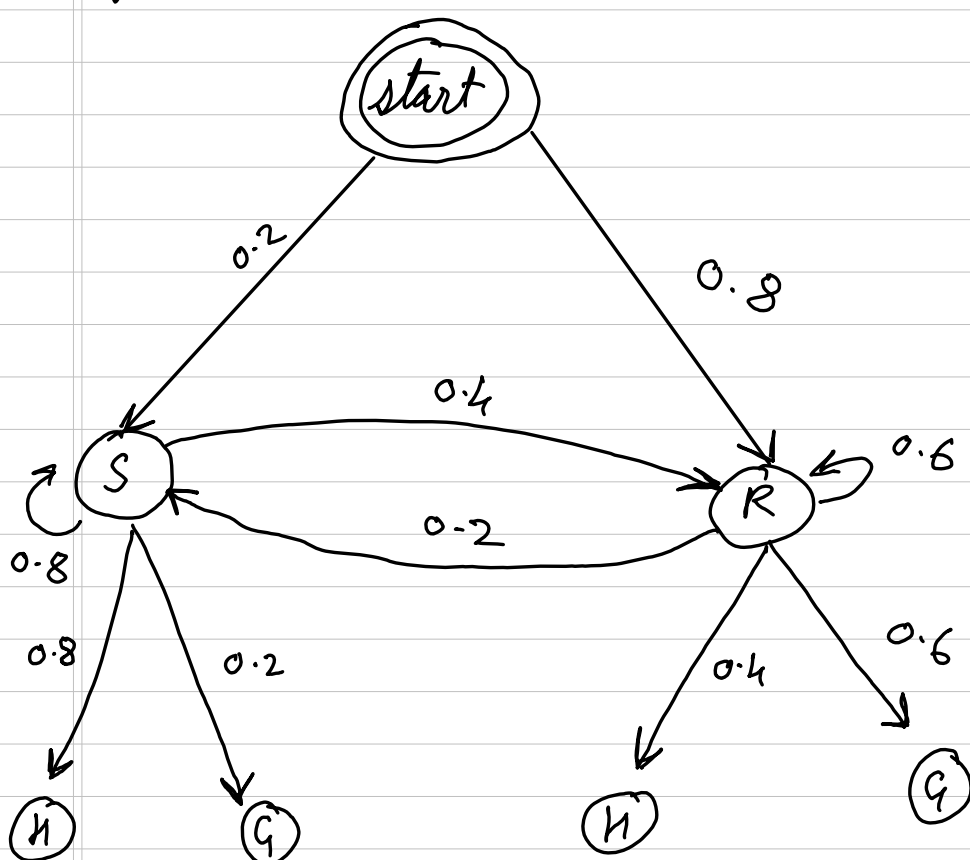$O$ is a sequence (No of icecreames
      eaten this week
      1, 2, 0, 0, 2, 2, 1)

states are hidden states traversed
    (Temperature of week
      hot, Hot, cold ....)

① Evaluation problem → calculate the
   probability that model generates $O$

② Decoding problem → calculate the most
   likely sequence of states visited
      for $O$

③ Learning problem → Determine HMM
   parameters that fit $O$

     HMM

Evaluation Decoding    Learning

        Supervised Unsupervised

     Gradient Viterbi Gradient EM

           Maximum
           likelyhood

① Given Markov Model



S- Sunny Day       H- Happy
R- Rainy day       G- Grumpy

Given Emotion recorded on 6 days

| Day | 1 | 2 | 3 | 4 | 5 | 6 |
|-----|---|---|---|---|---|---|
| Emotion | H | H | G | G | G | H |

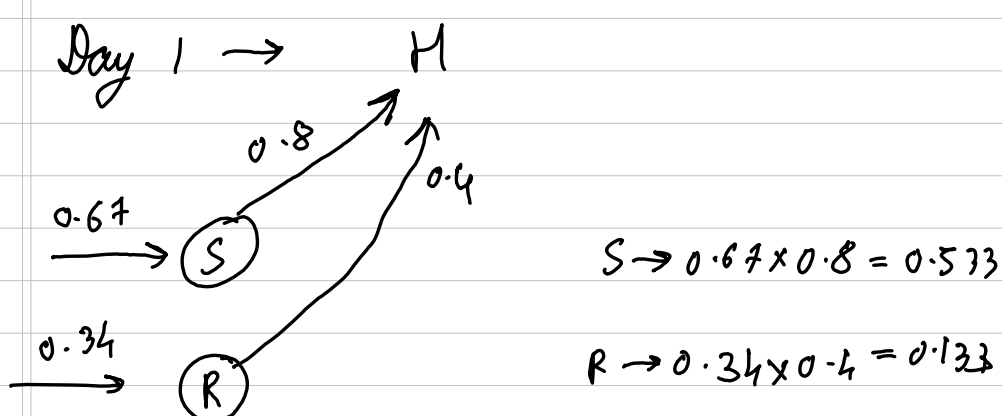Find the most likely climate on the 6 days

→

The Viterbi Algorithm is a dynamic programming algorithm.

It tracks the maximum probability and corresponding state sequence.

There may be many paths that lead to the following emotions

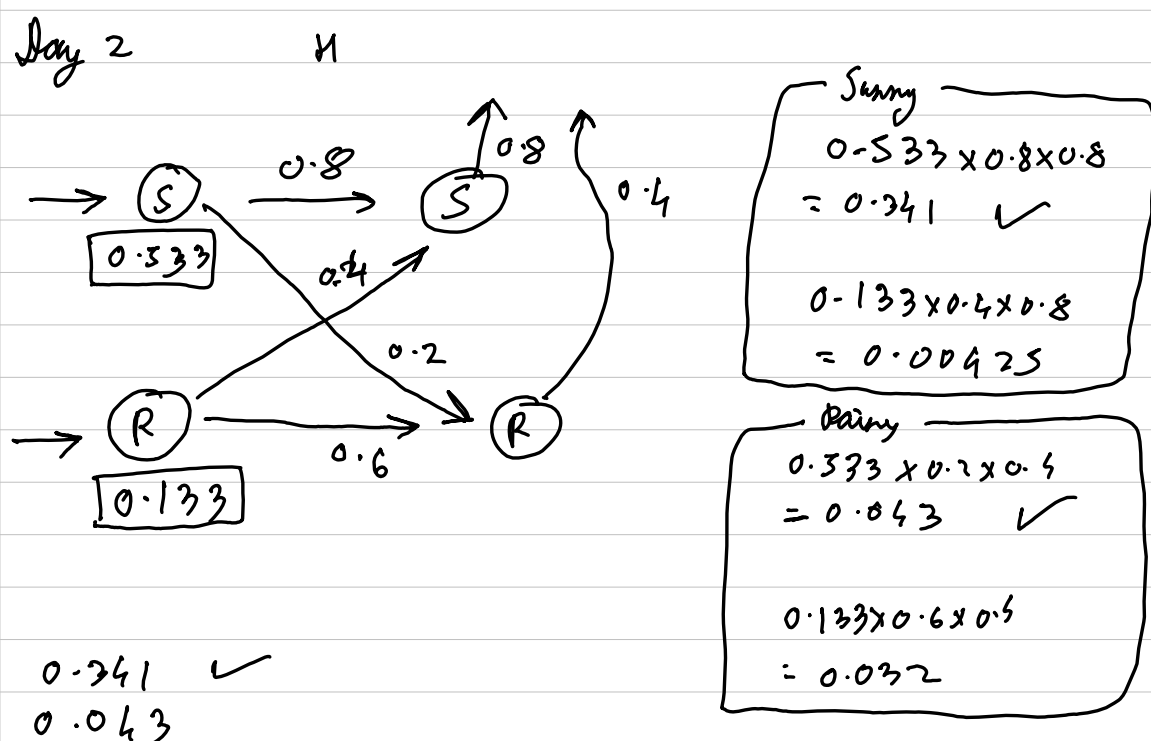eg every day might be sunny yet these emotions might be there.

But we need to find the most likely path

Day 1 →    H



$S → 0.67 \times 0.8 = 0.533$

$R → 0.34 \times 0.4 = 0.133$

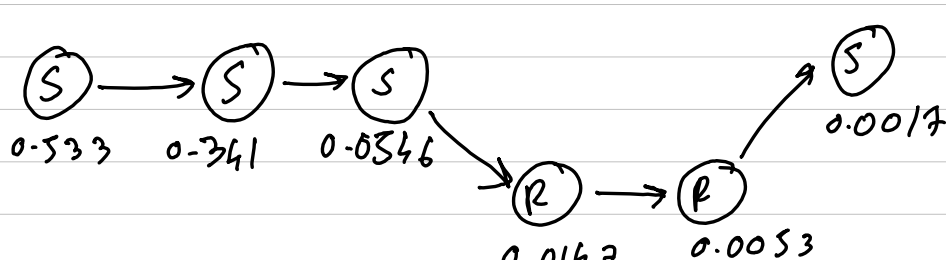Here the probability that the day will be sunny & emotion will be happy is 0.533

It is higher as opposed to day will be rainy and emotion will be happy

So we can say that Day 1 might be Sunny

Day 2        H



Sunny
$0.533 \times 0.8 \times 0.8$
$= 0.341$ ✓

$0.133 \times 0.4 \times 0.8$
$= 0.00425$

Rainy
$0.533 \times 0.2 \times 0.4$
$= 0.043$ ✓

$0.133 \times 0.6 \times 0.5$
$= 0.032$

0.341 ✓
0.043

∴ Sunny day

Finally we get

①     $S = \{ \text{Hot, cold} \}$    Day type    

$V = \{ V_1, V_2, V_3 \}$    No of icecreames    output
                              consumed    states

Example sequence   
$\left.\begin{array}{l} x_1 = V_2 \\ x_2 = V_3 \\ x_3 = V_1 \\ x_4 = V_2 \end{array}\right\}$ 4 days data

$A =$

| | H | C |
|---|---|---|
| H | 0.7 | 0.3 |
| C | 0.4 | 0.6 |

Transmission Matrix

$B =$

| | $V_1$ | $V_2$ | $V_3$ |
|---|---|---|---|
| H | 0.1 | 0.4 | 0.5 |
| C | 0.7 | 0.2 | 0.1 |

emission Matrix

$\pi =$

| H | C |
|---|---|
| 0.6 | 0.4 |

initial state

Find probability that sequence $x$ will be recorded

→



Day 0



$\left.\begin{array}{l} \text{(H)} \\ 0.6 \\ \\ \text{(C)} \\ 0.4 \end{array}\right\}$ Prior probabilities $\pi$

Day 1



For first sequence    $x_1 = V_2$

H → $(0.6 \times 0.7 \times 0.4) + (0.4 \times 0.4 \times 0.4) = 0.232$

C → $(0.4 \times 0.6 \times 0.2) + (0.6 \times 0.3 \times 0.2) = \underline{0.084}$

                                   $0.316$

∴ 0.316 is probability of buying 2 icecreames on day 1

Day 2



$x_2 = V_3$

∴      H = $0.232 \times 0.7 \times \cdots + 0.232 \times \cdots$

       C = $0.084 \times 0.6 \times \cdots$    $+ 0.084 \times \cdots$

                                    $\Sigma = \overline{\overline{0.11}}$

and so on

We get    $V_1 = 0.316$
           $V_2 = 0.11$
           $V_3 = 0.03296$
           $V_4 = 0.00966$

This is the probability of the sequences day wise

finally   Probability of sequence is
     $P(x_1) \times P(x_2) \times \cdots$

$= 0.316 \times 0.11 \times \cdots$

This is the probability that sequence $V_1 V_2 V_3 V_4$ will be recorded

"Evaluation problem"

# Note

We are <u>NOT</u> finding a random walk
(unlike Markov chains)

In random walk, probabilities are
calculated only from the past inputs

In evaluation problem, we are given
the X dataset (eg no of items eaten)

Probability is of day | X . Hence the

probability of outputs is also included

In decoding problem, X is given also
the path is to be found out

Hence the most likely paths are kept
the paths are not added.

Evaluation $\rightarrow$ chosen

Decoding $\rightarrow$ Added

# HMM Learning (Training)

HMM can be used in Supervised as well as unsupervised scenerios

Supervised → Data consists of sequences of observations along with the corresponding sequences of hidden states

   Temperature, icecreames → known

Unsupervised → No info about hidden states only observed sequences

       only icecreames    known

There are many methods for HMM learning

① Maximum likelyhood estimation → (Unsupervised)

      estimate parameters that maximize the likelyhood of the observed data sequence

      Baum Welch algorithms are used

② Expectation Maximization algorithm based training ( Unsupervised)

   Initialize    A, B, π

   for an observed sequence O

   E:  Calculate probabitly at being in state
         S at time t given O

   M:  Update transition probabilities

③ Viterbi training — when true state sequence is known or can be estimated (Supervised)

④ Gradient based optimizations

      Adjust parameters based on gradient descent on target functions

          (Supervised , Unsupervised)

Advantages → ① Flexibility

② Efficiency (low cost)

③ Interpretable

Disadvantages → ① Markov property may not hold always

② Overfitting

③ Not as robust as NN

Applications → ① Gene Prediction

② NLP (POS tagging)

③ SLAM (Robotics)

④ Speech recognition