

Unit 4

Syllabus: Sequential Data - Markov Models, HMM - Maximum likelihood for the HMM, The forward & Backward algorithm, the sum product algorithm, scaling factors, Viterbi algo., linear dynamical system.

Sequential data - Sequential data machine learning is a type of ML that uses models to process data that is ordered into sequences, & where the order of data point is important.

Example of Sequential Data - Textual data, Time series data, video streams, audio signals etc

Tasks that can be modelled using sequential data -

Text classification

Language Translation

Sentiment classification

Time series forecasting ie stock price prediction

* RNNs are well-known method in sequence models.

Markov Model

Markov models are also called Markov chains or markov processes.

These are probabilistic models that capture the dynamics & dependencies of a sequence of events.

Use of Markov model -

To generate sequence of data such as text, music or time series

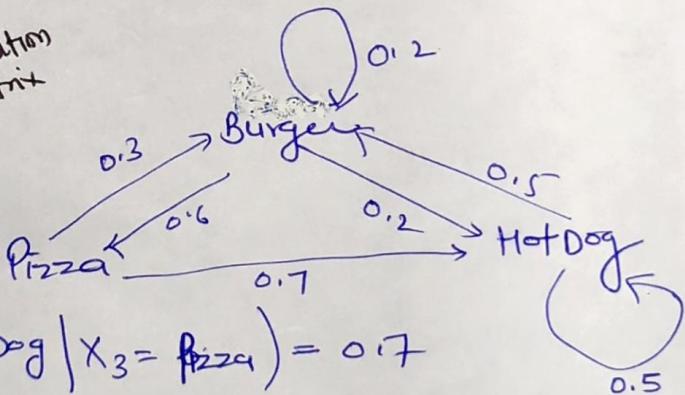
Markov properties -

① future state depends on current state only not on the previous states

$$P(X_{n+1} = x \mid X_n = x_n)$$

② Sum of the weights of outgoing edge from any state is equal to 1

$$\begin{matrix} & B & P & H \\ B & 0.2 & 0.6 & 0.2 \\ P & 0.3 & 0 & 0.7 \\ H & 0.5 & 0 & 0.5 \end{matrix} = A \quad \text{Transition Matrix}$$



$$P(X_4 = \text{HotDog} \mid X_3 = \text{Pizza}) = 0.7$$

Let 1st day one person ate Pizza

$$\pi_0 = [0 \ 1 \ 0]$$

$$\pi_0 A = [0 \ 1 \ 0] \begin{bmatrix} 0.2 & 0.6 & 0.2 \\ 0.3 & 0 & 0.7 \\ 0.5 & 0 & 0.5 \end{bmatrix} = \begin{bmatrix} 0.3 & 0 & 0.7 \\ | & | & | \\ P \rightarrow B = 0.3 & P \rightarrow P = 0 & P \rightarrow H = 0.7 \end{bmatrix}$$

future probability
of eating hot dog
state

$$\pi_1 A = [0.3 \ 0 \ 0.7] \begin{bmatrix} 0.2 & 0.6 & 0.2 \\ 0.3 & 0 & 0.7 \\ 0.5 & 0 & 0.5 \end{bmatrix} = [0.41 \ 0.18 \ 0.41]$$

Perform same process until we get equilibrium or stationary state i.e. output row vector = input row vector
ie.

$$\boxed{\pi_1 A = \pi_1} \text{ Eigen vector eqn}$$

$$Av = \lambda v$$

$$\pi = [0.35211 \quad 0.21127 \quad 0.43662]$$

Notes

- In Markov model, a sequence of events are represented by a series of states, and the transition b/w states are governed by probability
- Markov Model is represented as directed graph.

Order of Markov model - No of previous states considered when predicting the next state.

1st order markov model is also called markov chain.

Training Markov model - Estimating the transition probabilities from a given dataset.

Given a sequence of events, the model calculates the probabilities of transitioning from one state to another based on the observed frequency in the data.

Hidden Markov Model

A HMM is a statistical model that can be used to describe the evolution of observable events that depends on internal factors, which are not directly observable.

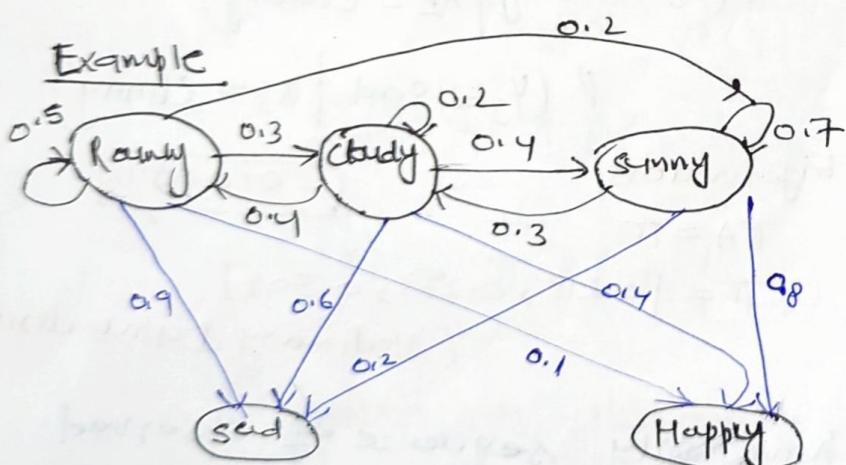
An HMM consists of two stochastic processes, an invisible process of hidden state & a visible process of hidden observable symbols.

HMM assumes Markov property

Application of HMM

Reinforcement learning

Temporal pattern Recognition ie speech handwriting



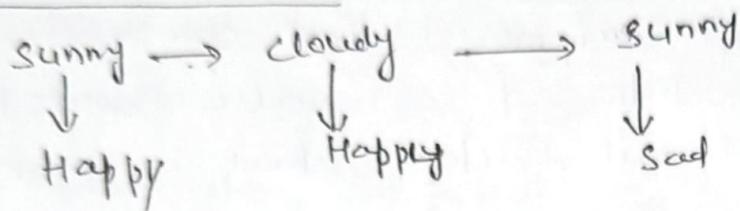
$$\begin{bmatrix} R & C & S \\ \begin{matrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{matrix} \end{bmatrix}$$

Transition matrix

$$\begin{bmatrix} R & C & S \\ \begin{matrix} \text{Sad} & \text{Happy} \\ \begin{bmatrix} 0.9 & 0.1 \\ 0.6 & 0.4 \\ 0.2 & 0.8 \end{bmatrix} \end{matrix} \end{bmatrix}$$

EMMM Matrix

Let us consider



WHAT is the probability of occurrence of above scenario?

or

$$P(Y = \text{Happy Happy Sad}, X = \text{sunny cloudy sunny}) ?$$

Calculation can be performed as follows

$$P(X_1 = \text{sunny}) \quad P(Y_1 = \text{Happy} | X_1 = \text{sunny})$$

$$P(X_2 = \text{cloudy} | X_1 = \text{sunny}) \quad P(Y_2 = \text{Happy} | X_2 = \text{cloudy})$$

$$P(X_3 = \text{sunny} | X_2 = \text{cloudy})$$

$$P(Y_3 = \text{sad} | X_3 = \text{sunny})$$

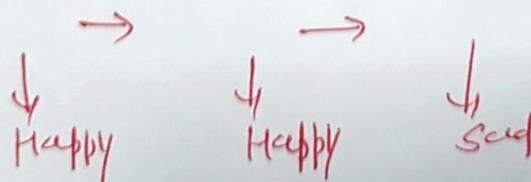
Using left Eigen vectors

$$\pi A = \pi$$

$$\pi = [0.218, 0.273, 0.509]$$

(Stationary Distribution)

But if we have only sequence of observed variables then what is the most likely weather sequence for the observed mood sequence given below?



Solution To find most likely sequence we have to find/ calculate the probability for each of the possible weather sequences to find the sequence with MAX probability.

$$P(Y = \text{Happy Happy Sad}, X = \text{sunny sunny cloudy})$$

Y = Observed variable
 X = Markov chain

$$\arg \max_{X=x_1, x_2, \dots, x_n} P(X=x_1, x_2, x_3, \dots, x_n | Y=y_1, y_2, \dots, y_n)$$

Using Bayes Theorem

$$\arg \max_{X=x_1, x_2, \dots, x_n}$$

$$\frac{P(Y|X) P(X)}{P(Y)}$$

Ignore this

$$P(Y|X) = \prod P(Y_i | X_i)$$

$$P(X) = \prod P(X_i | X_{i-1})$$

for first term use stationary distribution vector
∴ Eqn 1 can be written as -

$$\arg \max_{X=x_1, x_2, \dots, x_n} \prod P(Y_i | X_i) P(X_i | X_{i-1})$$

Forward Algorithm for HMM

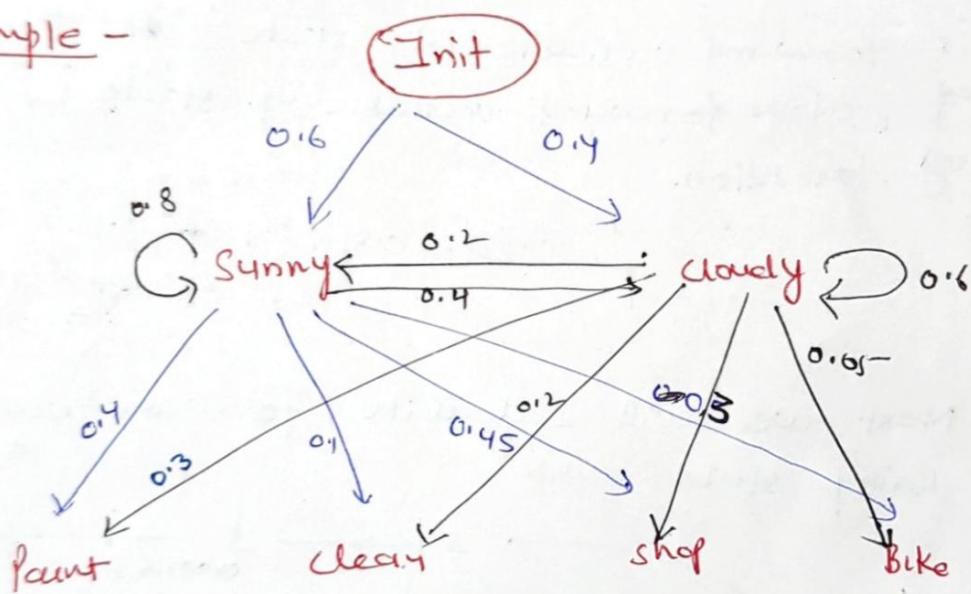
Forward Algorithm comprised of three steps -

- ① Initialization $\alpha_1(i) = \pi_i \cdot b_i(O_1)$
- ② Recursion
- ③ Termination

$$\text{Recursion } \alpha_{t+1}(j) = \sum_{i=1}^N \alpha_t(i) \cdot a_{ij} b_j(O_{t+1})$$

$$\text{Termination } P(O|A) = \sum_{i=1}^N \alpha_T(i)$$

Example -



Hidden states - sunny & rainy

Observable states - paint, clear, shop & bike.

Initial probability $= [0.6, 0.4]$

Transition probability

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$B = \begin{matrix} \text{Emission Probability} \\ \begin{bmatrix} b_{11} & b_{12} & b_{13} & b_{14} \\ b_{21} & b_{22} & b_{23} & b_{24} \end{bmatrix} \end{matrix}$$

$$B = \begin{bmatrix} b_{11} & b_{12} & b_{13} & b_{14} \\ b_{21} & b_{22} & b_{23} & b_{24} \end{bmatrix}$$

Question - what is the likelihood that observation sequence $O = \{O_1, O_2, O_3, O_4\} = \{\text{paint}, \text{clear}, \text{shop}, \text{bike}\}$ can be derived from HMM λ ? $P(O|\lambda) = ???$

Forward Algorithm for HMM.

Initialization

Initial probability of sunny = 0.6
" of sunny to point = 0.4

\therefore Initial forward variable for sunny = 0.24

$$\underline{\text{Rainy state}} = 0.4 \times 0.3 = 0.12$$

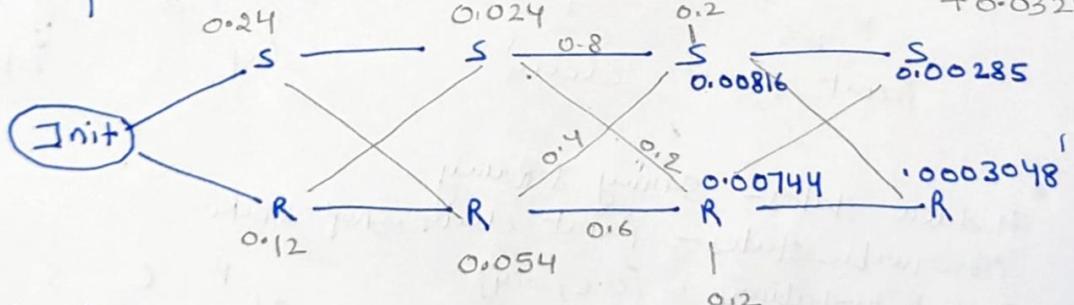
Recursion

$$\alpha_{t+1}(j) = \sum_{i=1}^N \alpha_t(i) \cdot a_{ij} \cdot b_j(O_{t+1})$$

i.e forward variable of state j is the product
 of previous forward variable of state i , multiplied
 by transition.

$$0.24 \times 0.8 \times 0.1 + 0.12 \times 0.4 \times 0.1 \\ = 0.192 + 0.0048 \\ = 0.1968$$

Next we will calculate forward variable for Rainy state.



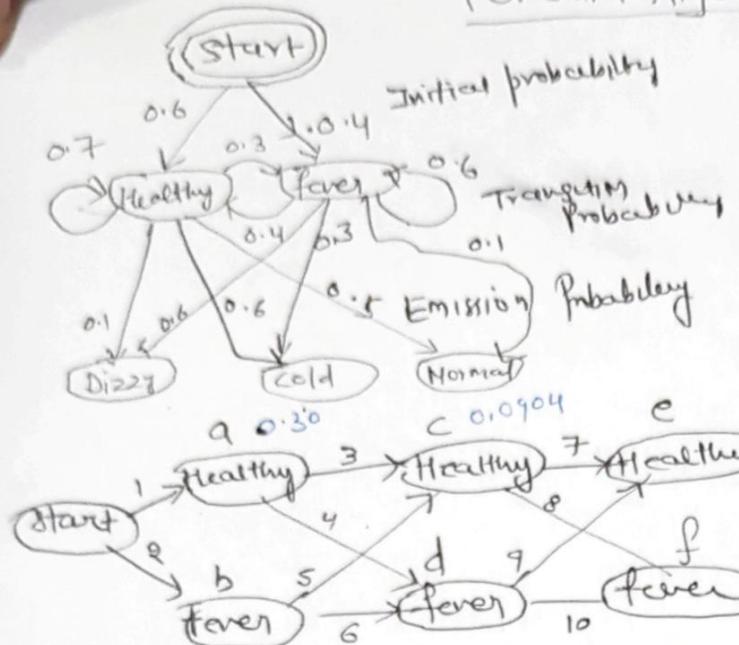
Termination

$$\rho(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$$

$$= 0.0028512 + 0.0003048$$

$$= 0.003156$$

Forward Algorithm for HMM



	D	C	H
H	• 1	• 4	• 5
F	• 6	• 3	• 1

H	• 7	• 3
F	• 4	• 6

Observed sequence
= Normal cold dizzy

$$\begin{aligned} \textcircled{1} &\Rightarrow 0.6 \times 0.5 = 0.3 \quad (\text{start} \rightarrow \text{Healthy} \rightarrow \text{Normal}) \\ \textcircled{2} &\Rightarrow 0.4 \times 0.1 = 0.04 \quad (\text{start} \rightarrow \text{Fever} \rightarrow \text{Normal}) \end{aligned}$$

(Mode A)
(Mode B)

Node C

$$0.3 \times 0.7 \times 0.4 = 0.084 \quad (\text{start} \rightarrow \text{Healthy} \rightarrow \text{Healthy} \rightarrow \text{Cold})$$

$$\text{Mode } 0.04 \times 0.4 \times 0.4 = 0.0064 = 0.0904$$

Node D

$$0.3 \times 0.3 \times 0.3 = 0.027 + = 0.0342$$

$$0.04 \times 0.6 \times 0.3 = 0.0072$$

Node E

$$0.0904 \times 0.7 \times 0.1 = 0.006328$$

$$0.0342 \times 0.4 \times 0.1 = 0.001368 + = 0.007696$$

Node F

$$0.0904 \times 0.3 \times 0.6 = 0.016272 + \Rightarrow 0.028584$$

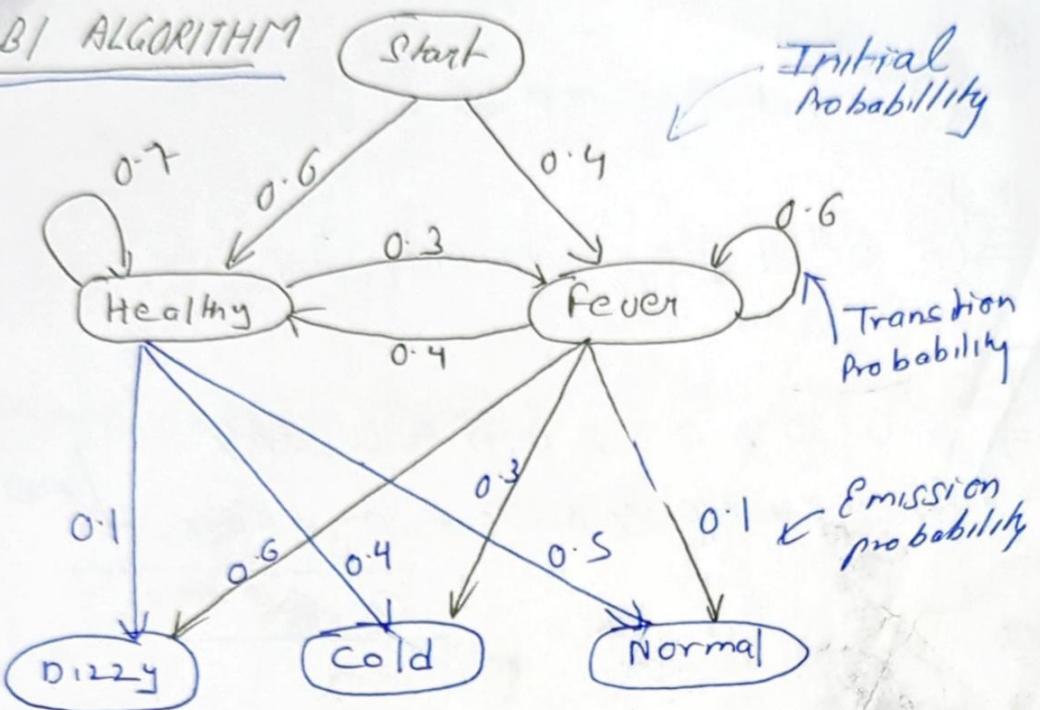
$$0.0342 \times 0.6 \times 0.6 = 0.012312 +$$

for sequence
Healthy Healthy fever the

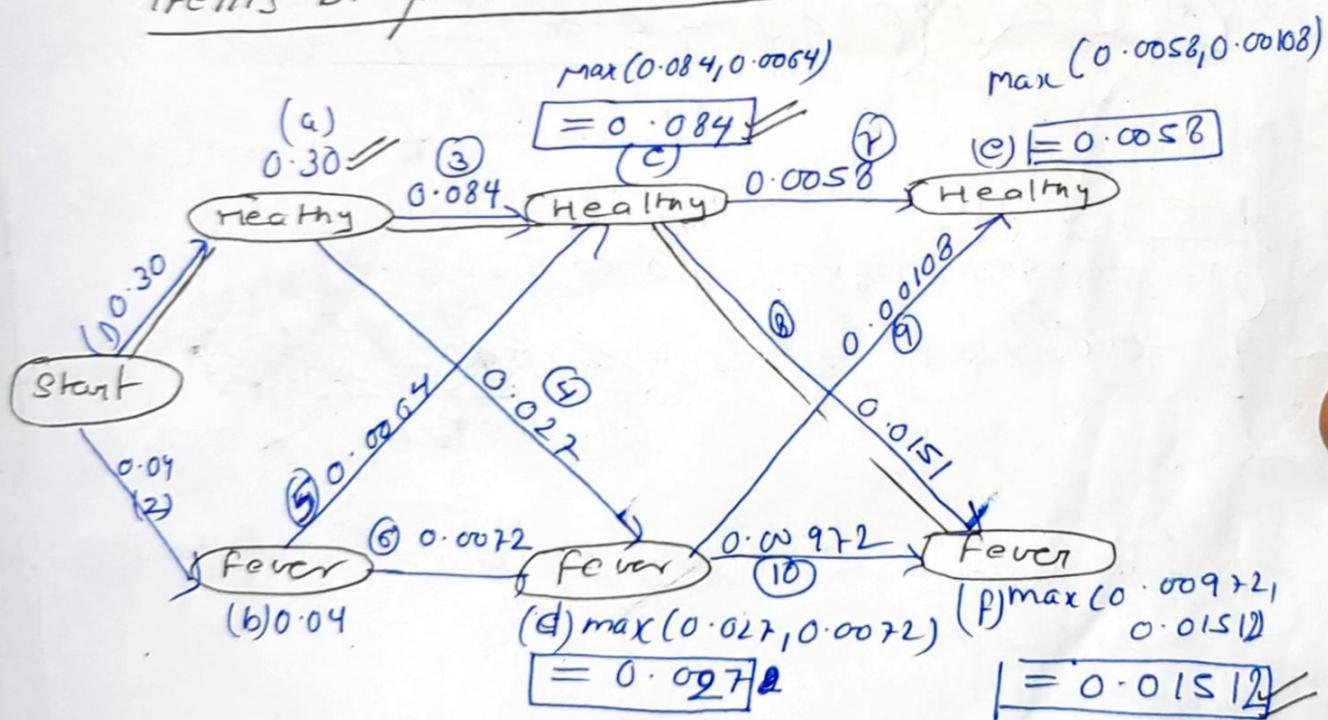
Observation Normal cold ~~dizzy~~ will be
most suitable observation.

for given observation H HF one most likely
transition.

TERBI ALGORITHM



Trellis Diagram (Time Based Diagram)



```
graph LR; A[Normal] --> B[Cold]; B --> C[Dizzy]
```

Healthy → Healthy → Fever

~~Q 9~~

$$\rightarrow \textcircled{1} 0.0 + 0.5 = 0.30$$

~~b~~

$$\rightarrow \textcircled{2} 0.4 \times 0.1 = 0.04$$

~~c~~

$$\rightarrow \textcircled{3} 0.30 \times 0.7 \times 0.4 = 0.084$$

$$\rightarrow \textcircled{5} 0.04 \times 0.4 \times 0.4 = 0.0064 \quad]^{\max}$$

$$\underline{0.084}$$

~~d~~

$$\rightarrow \textcircled{4} 0.30 \times 0.3 \times 0.3 = 0.027, \quad]^{\max}$$

$$\rightarrow \textcircled{6} 0.04 \times 0.6 \times 0.3 = 0.0072 \quad]^{\max}$$

$$\underline{0.027}$$

~~e~~

$$\rightarrow \textcircled{7} 0.094 \times 0.7 \times 0.1 = 0.006328 \quad]^{\max}$$

$$\rightarrow \textcircled{9} 0.0342 \times 0.4 \times 0.1 = 0.001368 \quad]^{\max}$$

$$\underline{0.006328}$$

And similarly for R_1, F_1, g find max value

$$\rightarrow \textcircled{7} 0.084 \times 0.2 \times 0.1 = 0.00588 \quad \}^{\max}$$

$$\rightarrow \textcircled{9} 0.027 \times 0.4 \times 0.1 = \frac{0.0108}{0.00588}$$

$$\underline{0.00588}$$

$$\rightarrow \textcircled{8} 0.84 \times 0.3 \times 0.6 = 0.1512 \quad]^{\max} = \underline{0.1512}$$

$$\rightarrow \textcircled{10} 0.017 \times 0.6 \times 0.6 = 0.0097 \quad]^{\max} = \underline{0.0097}$$

$$A = \begin{matrix} & \text{B} & \text{P} & \text{H} \\ \text{B} & \left[\begin{matrix} 0.2 & 0.6 & 0.2 \end{matrix} \right] \\ \text{P} & \left[\begin{matrix} 0.3 & 0 & 0.7 \end{matrix} \right] \\ \text{H} & \left[\begin{matrix} 0.5 & 0 & 0.5 \end{matrix} \right] \end{matrix}$$

(Transition Matrix)

π → denotes probability distribution of the states

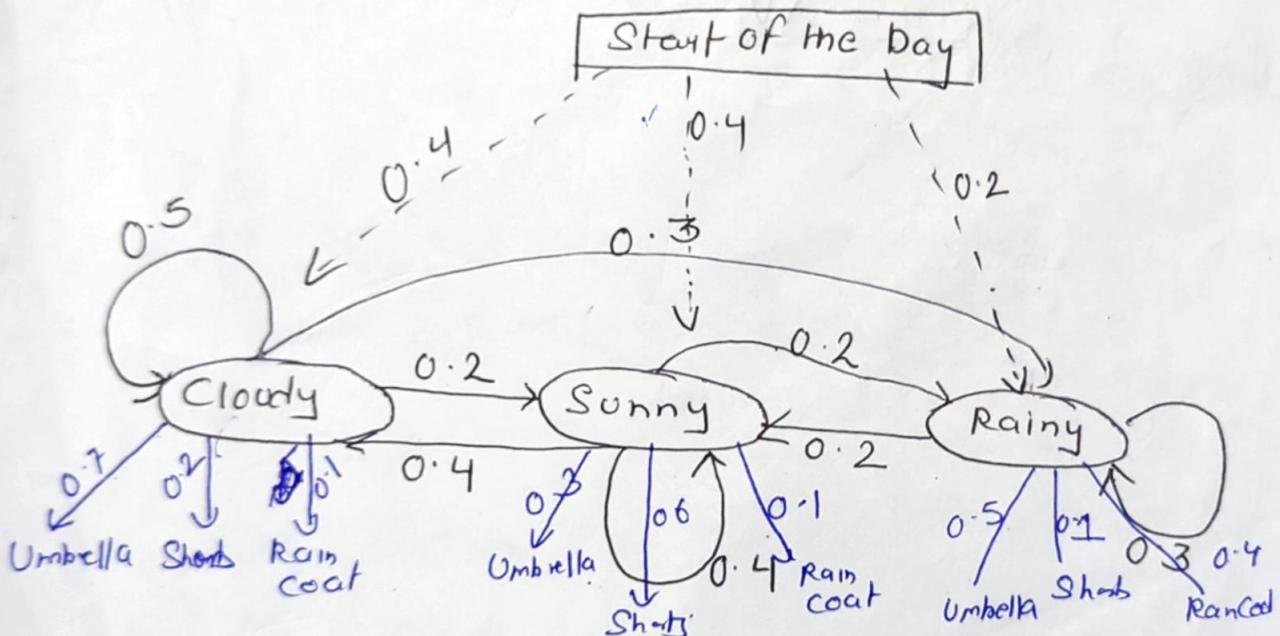
Suppose on beginning we are on Pura day $\pi_0 = [0 \ 1 0]$

$$\pi_0 A = [0 \ 1 0] \left[\begin{matrix} 0.2 & 0.6 & 0.2 \\ 0.3 & 0 & 0.7 \\ 0.5 & 0 & 0.5 \end{matrix} \right] = [0.34 \ 0.25 \ 0.41]$$

$$\pi A = \pi$$

$$\boxed{\pi A v = \lambda v} \text{ Eigen Vector Equation .}$$

HMM



	Cloudy	Sunny	Rainy	
Cloudy	0.5	0.2	0.3	= 1
Sunny	0.4	0.4	0.2	= 1
Rainy	0.3	0.2	0.5	

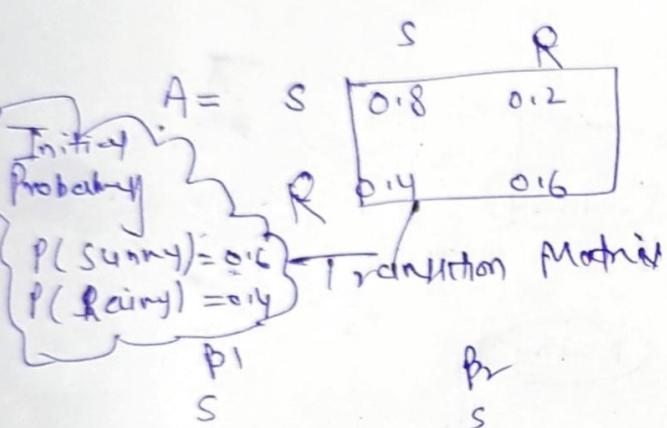
(Transition Probability)

Backward HMM Algorithm

Question for Observation sequence

Paint (P) Clean (C) Shop (S) Bike (B)

Calculate backward probability after applying
Backward HMM algorithm.



Emission Matrix

	P	C	S	B
S	0.4	0.1	0.2	0.3
R	0.3	0.45	0.2	0.05

$$\beta_3(S) = T(S \rightarrow S) B(S|S) \times \beta_4(S) + T(S \rightarrow R) B(S|R) \times \beta_4(R)$$

$$\begin{aligned} \beta_2(S) &= T(S \rightarrow S) B(S|S) \times \beta_3(S) \\ &\quad + T(S \rightarrow R) B(S|R) \times \beta_3(R) \end{aligned}$$

↑
Paint Clean Shop Bike

$$\beta_3(\text{Sunny}) = A(S \rightarrow S) \times B(S|S) \times \beta_4(S)$$

$$+ A(S \rightarrow R) \times B(S|R) \times \beta_4(R)$$

$$\begin{aligned} \beta_3(\text{Rainy}) &= A(R \rightarrow S) B(S|S) \times \beta_4(S) \\ &\quad + A(R \rightarrow R) B(S|R) \times \beta_4(R) \end{aligned}$$

Backward Algorithm

Observation = {Opaint, Oclean, Oshop, Obike}

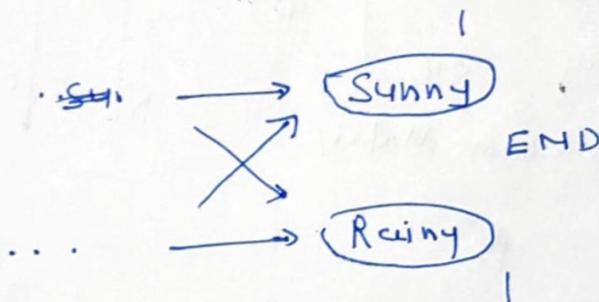
Backward algorithm runs backward in time

It also has 3 steps

- ① Initialization
- ② Recursion
- ③ Termination

Initialization

$$\beta_T(i) = 1$$



calculate the backward variable of previous state sunny.

= transition probability from s to s times emission probability of bike from sunny state

$$0.8 \times 0.3 \times 1 + 0.2 \times 0.05 \times 1 = 0.01 \\ = 0.25$$

calculate backward Variable for previous state
rainy

$$= 0.6 \times 0.05 \times 1 + 0.2 \times 0.05 \times 1 = 0.25$$

Backward Algorithm

It goes backward in time

- Steps
- ① Initialization
 - ② Recursion
 - ③ Termination

Initialization

$$\beta_T(i) = 1$$

At time T (at the end of the observation sequence) the backward variable of every state is equal to 1

Recursion

$$\beta_t(i) = \sum_{j=1}^M a_{ij} \cdot b_j^{\text{Emission Probability}}(O_{t+1}) \cdot \beta_{t+1}(j)$$

↓
Transition probability

$\beta_4(\text{sunny}) = 1$
 $\beta_4(\text{Rainy}) = 1$

Backward variable

For sunny at $t=3$:

$$\begin{aligned} P_3(\text{sunny}) &= A(\text{sunny} \rightarrow \text{sunny}) \times B(\text{Bike} | \text{sunny}) \times \\ &\quad \beta_4(\text{sunny}) \\ &+ A(\text{sunny} \rightarrow \text{Rainy}) \times B(\text{Bike} | \text{Rainy}) \times \beta_4(\text{Rainy}) \end{aligned}$$

$$\beta_3(\text{sunny}) = 0.8 \times 0.3 \times 1 + 0.2 \times 0.05 \times 1 = 0.25$$

$$\begin{aligned} P_3(\text{Rainy}) &= A(\text{Rainy} \rightarrow \text{sunny}) \times B(\text{Bike} | \text{Rainy}) \times \beta_4(\text{Rainy}) \\ &+ A(\text{Rainy} \rightarrow \text{Rainy}) \times B(\text{Bike} | \text{Rainy}) \times \beta_4(\text{Rainy}) \\ &= 0.4 \times 0.3 \times 1 + 0.6 \times 0.05 \times 1 = 0.15 \\ \beta_3(\text{Rainy}) &= 0.15 \end{aligned}$$

For sunny at $t=2$

$$\beta_2(\text{sunny}) = A(\text{sunny} \rightarrow \text{sunny}) \times B(\text{shop}|\text{sunny}) \times \beta_3(\text{sunny}) \\ + A(\text{sunny} \rightarrow \text{Rainy}) \times B(\text{shop}|\text{Rainy}) \times \beta_3(\text{Rainy})$$

for Rainy at $t=2$ =

$$\beta_2(\text{Rainy}) = 0.8 \times 0.2 \times 0.25 + 0.2 \times 0.2 \times 1.5 = 0.046$$

$$\beta_2(\text{Rainy}) = A(\text{Rainy}|\text{sunny}) \times B(\text{Shop}|\text{sunny}) \times \beta_3(\text{sunny}) \\ + A(\text{Rainy}|\text{Rainy}) \times B(\text{Shop}|\text{Rainy}) \times \beta_3(\text{Rainy}) \\ = 0.4 \times 0.2 \times 0.25 + 0.6 \times 0.2 \times 0.15$$

for sunny at $t=1$

$$0.038$$

$$\beta_1(\text{sunny}) = A(\text{sunny} \rightarrow \text{sunny}) \times B(\text{clean}|\text{sunny}) \times \beta_2(\text{sunny}) \\ + A(\text{sunny} \rightarrow \text{Rainy}) \times B(\text{clean}|\text{Rainy}) \times \beta_2(\text{Rainy}) \\ = 0.8 \times 0.1 \times 0.046 + 0.2 \times 0.45 \times 0.038 \\ = 0.0071$$

for Rainy at $t=1$

$$\beta_1(\text{Rainy}) = A(\text{Rainy} \rightarrow \text{sunny}) \times B(\text{clean}|\text{sunny}) \times \beta_2(\text{sunny}) \\ + A(\text{Rainy} \rightarrow \text{Rainy}) \times B(\text{clean}|\text{Rainy}) \times \beta_2(\text{Rainy}) \\ = 0.4 \times 0.1 \times 0.046 + 0.6 \times 0.45 \times 0.038$$

Termination

$$= 0.0121$$

$$P(O/\lambda) = \sum_{i=1}^N \pi_i \cdot b_i(O_1) \cdot \beta_1(i)$$

π_i = initial probability of being in state (sunny or Rainy)
 $b_i(O_1)$ is the emission probability of first observation O_1 ,
 β_1 is the backward probability at time $t=1$

$$\pi(\text{sunny}) = 0.6 \quad \pi(\text{Rainy}) = 0.4$$

Emission probabilities for point

$$B(\text{Point}|\text{sunny}) = 0.4 \quad B(\text{Point}|\text{Rainy}) = 0.3$$

$$\beta_1(\text{sunny}) = 0.0071 \quad \beta_1(\text{Rainy}) = 0.0121$$

considering sunny state :

$$P(O) = \pi(\text{sunny}) \times B(\text{Paint} | \text{sunny}) \times \beta_1(\text{sunny}) \\ = 0.6 \times 0.4 \times 0.0071 = 0.001704$$

Considering rainy state

$$P(O) = \pi(\text{Rainy}) \times B(\text{Paint} | \text{Rainy}) \times \beta_1(\text{Rainy}) \\ = 0.4 \times 0.3 \times 0.0121 = 0.001452$$

Final value

$$= P(O|A) = 0.001704 + 0.001452 \\ = 0.003156$$

1. SCALING

What is scaling

In Hidden Markov Models (HMMs), scaling is a technique used to prevent numerical underflow during the computation of probabilities in algorithms like the forward and backward algorithms.

Why Scaling is needed

When dealing with long sequences of observations, the probabilities calculated in HMMs can become extremely small, leading to numerical underflow (i.e., values so small that they can't be represented accurately by the computer). This happens because the probabilities are products of many small values, and multiplying many such small values together results in a very small number.

How Does Scaling Work?

To avoid this problem, scaling factors are introduced at each time step to normalize the probabilities. The scaling procedure typically involves the following steps:

1. Calculate the scaling factor: For each time step, compute a scaling factor that normalizes the forward or backward probabilities. This factor is usually the sum of the probabilities at that time step.
2. Scale the probabilities: Divide the probabilities at each time step by the scaling factor. This normalization step keeps the values within a manageable numerical range.
3. Store the scaling factors: Keep track of the scaling factors for each time step so that the original probabilities can be recovered if needed (e.g., for calculating the likelihood of the observed sequence).
4. Benefits of Scaling

1. Prevents numerical underflow: Keeps the values of probabilities from becoming too small to represent.

2. Improves numerical stability: Ensures that the calculations remain accurate even for long sequences.

Methods for scaling in data

Scaling in data preprocessing is essential in machine learning and statistical modelling to ensure that different features contribute equally to the learning process. Different methods of scaling are used based on the data type and algorithm. Here are some common scaling techniques:

1. Min-Max Scaling (Normalization)

- This technique transforms the data to a fixed range, typically [0, 1] or [-1, 1].
- The formula is: $X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$

Where X is original value and, X_{\min} and X_{\max} are minimum and maximum value of features.

- **Use Case:** Recommended when the data distribution does not follow a Gaussian distribution and when you need to bound the data values to a certain range (e.g., image pixel values).

2. Standardization (Z-score normalization)

- This method centers the data around the mean with a unit standard deviation.
- The formula is: $X' = \frac{X - \mu}{\sigma}$

- Where X is the original value, μ is the mean of the feature, and σ is the standard deviation.
- **Use Case:** Commonly used in algorithms like Support Vector Machines, k-means clustering, and Principal Component Analysis (PCA) where the data distribution should be normalized to a mean of 0 and standard deviation of 1.

3. Max Abs Scaling

- This technique scales each feature by the maximum absolute value of the feature. The resulting values range between -1 and 1.
- The formula is: $X' = X / \max(|X|)$
- **Use Case:** Useful when the data has a sparse representation and you want to preserve the sign of the original data.

4. Robust Scaling

- This method uses the median and the interquartile range (IQR) for scaling.
- The formula is: $X' = (X - \text{median}) / \text{IQR}$
- Where IQR is the difference between the 75th and 25th percentiles.
- **Use Case:** Suitable when the data contains outliers, as it is robust to the presence of extreme values.

5. Log Transformation

- It applies a logarithmic function to the data to reduce skewness.
- The formula is: $X' = \log(X + 1)$
- **Use Case:** Used when the data is highly skewed, as it helps to stabilize variance and make the data more normally distributed.

6.. Power Transformation (Box-Cox or Yeo-Johnson)

- A more advanced technique used to stabilize variance and make data more normal distribution-like.
- The transformation can be done using a power function (Box-Cox) or an optimized method (Yeo-Johnson).
- **Use Case:** Applied when the data needs to be transformed to approximate a normal distribution.

7.. Unit Vector Scaling (Normalization to unit norm)

- This scales the feature vector to have a unit norm (e.g., 1).
- The formula is: $X' = X / \|X\|$
- **Use Case:** Commonly used in text classification or clustering where you want to normalize the feature vectors' length.

Why is scaling important when training machine learning models?

Scaling is an important preprocessing step when training machine learning models because it ensures that features contribute equally to the learning process. Here are the key reasons why scaling is essential:

1. Prevents Features from Dominating Others

- When features have different scales, those with larger ranges may dominate the learning process and skew the results.

- For example, if one feature ranges from 0 to 1 (e.g., "age") and another ranges from 1,000 to 100,000 (e.g., "income"), the model might prioritize the feature with the larger range, even if it's not more important.
- Scaling puts all features on the same scale, allowing the model to treat them with equal importance.

2. Improves Convergence Speed in Gradient Descent

- Algorithms like Gradient Descent, which are used in many optimization-based machine learning algorithms (e.g., linear regression, neural networks), are sensitive to the scale of the features.
- When features are on different scales, the optimization process may take longer to converge because it has to navigate uneven gradients.
- Scaling helps to make the gradient smoother, thus speeding up the convergence.

3. Ensures Fair Distance Calculations

- Distance-based algorithms such as k-Nearest Neighbors (k-NN), Support Vector Machines (SVM), and clustering algorithms like k-means are highly sensitive to the scale of data.
- These algorithms rely on distance metrics like Euclidean distance, which will be affected if one feature has a much larger range than others.
- Scaling ensures that each feature contributes equally to the distance calculation, leading to better performance of these models.

4. Helps with Regularization Techniques

- Regularization techniques like Lasso (L1) and Ridge (L2) regression add a penalty term to the loss function based on the magnitude of feature coefficients.
- If features are on different scales, the penalty may disproportionately affect some coefficients, leading to suboptimal regularization.
- Scaling helps to ensure that the regularization term affects all features equally.

5. Avoids Numerical Instability

- Some algorithms are sensitive to the numerical ranges of the data, which can cause numerical instability during computation.
- For instance, neural networks may suffer from exploding or vanishing gradients when the input features have very large or very small values.
- Scaling features to a smaller range (e.g., [0, 1] or [-1, 1]) helps avoid these issues.

6. Improves Model Interpretability

- In some cases, having features on a common scale makes it easier to interpret the results.
- For instance, when visualizing decision boundaries in SVMs or clusters in k-means, scaled data provides a more meaningful representation.

2. Linear Dynamical System (LDS)

3. What is LDS

Linear Dynamical Systems (LDS) in machine learning are models used to represent systems that evolve over time in a linear manner. LDS is used to model time series data, where observations are made over time and there is an underlying latent state that evolves linearly with some noise. LDS is particularly useful in tasks like tracking, filtering, and predicting future states in a sequence.

Components of Linear Dynamical Systems

1. State Vector (x_t):
 - o Represents the hidden or latent state of the system at time t .
 - o This vector captures the underlying factors that drive the evolution of the system over time.
2. Observation Vector (y_t):
 - o Represents the observed data at time t , which is a noisy measurement of the state.
 - o In many cases, it is a linear transformation of the state vector with added noise.
3. State Transition Matrix (A):
 - o Defines the linear relationship between the state vector at time t and the state vector at time $t+1$.
 - o It determines how the system evolves from one state to the next.
4. Observation Matrix (C):
 - o Maps the latent state vector to the observed data.
 - o This matrix defines the relationship between the hidden states and the observations.
5. Process Noise (w_t):
 - o Represents random fluctuations in the state evolution.
 - o It accounts for the uncertainty in the state transition.
6. Observation Noise (v_t):
 - o Represents the noise or error in the observations.
 - o It accounts for the measurement uncertainty in the observed data.
7. Initial State Distribution (x_0):
 - o Specifies the distribution of the initial state of the system.
 - o Usually assumed to be a Gaussian distribution with a mean and covariance.

Example of Linear Dynamical Systems

An example of LDS is tracking the position and velocity of an object moving in one dimension. In this case:

- The state vector (x_t) could represent the position and velocity of the object.
- The observation vector (y_t) could represent the observed position of the object, possibly with some measurement noise.
- The state transition matrix (A) would model how the position and velocity change over time, while the observation matrix (C) could map the state (position and velocity) to the observed position.

How can the state transition process in an HMM be modeled as a linear dynamical system?

The state transition process in a Hidden Markov Model (HMM) can be extended to be modeled as a Linear Dynamical System (LDS) by using continuous-valued hidden states and linear relationships for transitions and observations. Here's how this process works:

Key Differences and Modifications

1. State Representation:

- In an HMM, states are discrete, where each state represents a specific condition or class.
- In an LDS, the states are continuous-valued vectors, allowing the system to represent a range of values rather than discrete labels.

2. State Transition Process:

- In an HMM, the state transition is governed by a transition probability matrix, where each element represents the probability of transitioning from one state to another.
- In an LDS, the state transition is modeled as a linear equation