

# Probability Course

## Lecture 6

Created by : Israa Abdelghany

LinkedIn : [www.linkedin.com/in/israa-abdelghany-4872b0222](https://www.linkedin.com/in/israa-abdelghany-4872b0222)

GitHub : <https://github.com/IsraaAbdelghany9>

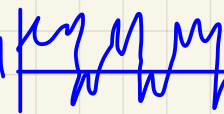
# Session 6

we will discuss:-

- Markov process / Markov chain
- Entropy (Information Entropy)

## Random Process (stochastic process)

↳ vs Random variable

Process  $\Rightarrow$  changes with time, not only value, signal   
variable  $\Rightarrow$  number / value.  
 $x(\xi, t)$   
Random var.

Example of process is temperature in different places  $\Rightarrow$  Rand. var

measuring the temperature over the day (time)  $\Rightarrow$  Random process



ensemble of different / possible  
Realizations of process

Images: Random process  
"spatial domain"  
space

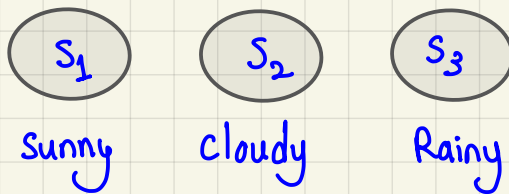
Markov Random process :- special type of Random processes



Discrete → Continuous will not be discussed.

depend on divide my process into some possible states.

نستخدم في حساب  
و توقع للاقتصاد  
و البوابة  
signal processing



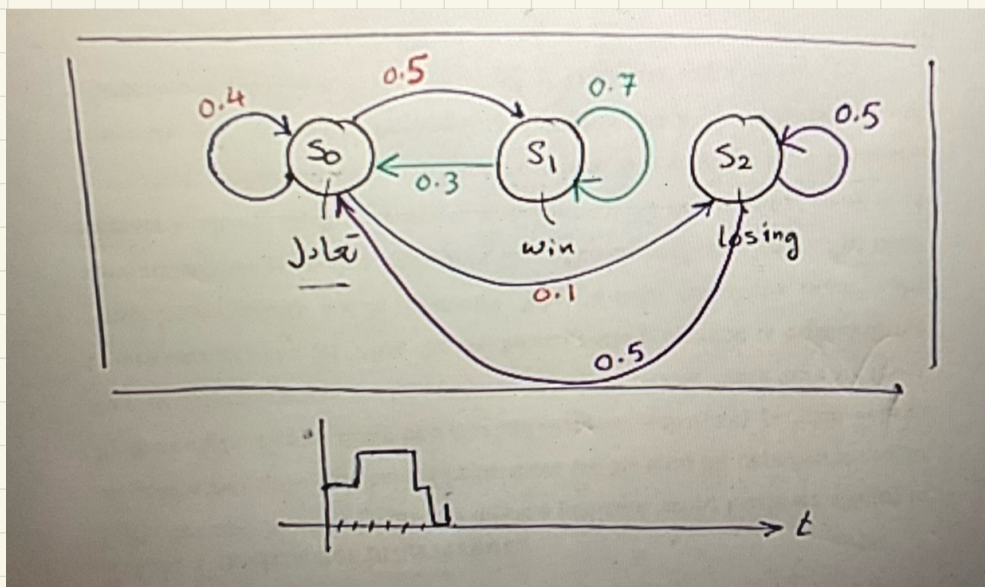
but note that not all processes that can be divided into states is a markov Random process

⇒ it has a characteristic:-

- ↳ memory less
  - ↳ next state don't depend on last state
- ↳ short term memory

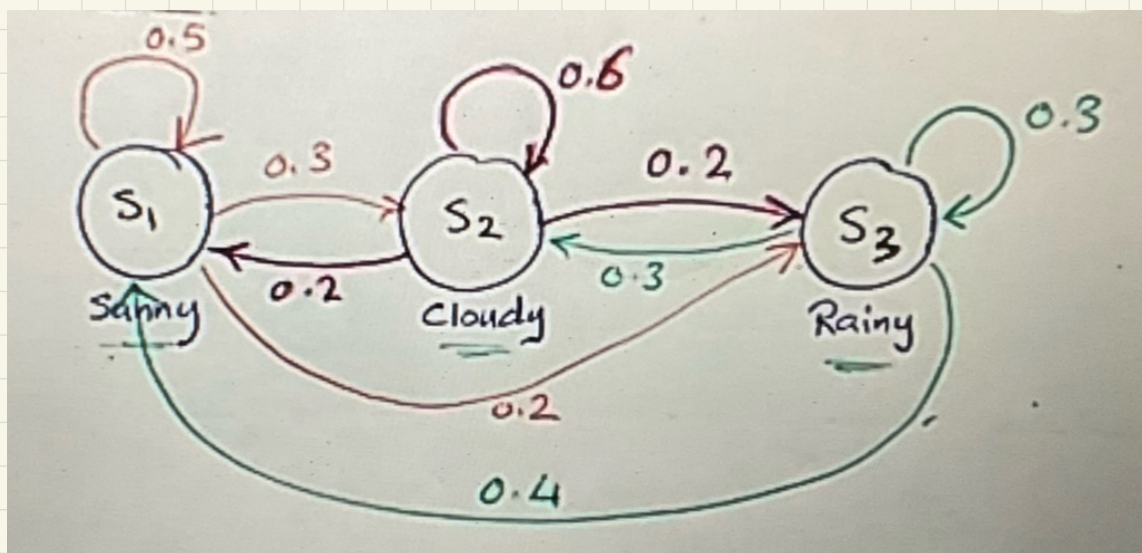
depend or has a value of prediction depend on the current state  
not affected by the future

Example 2:- match



$$\begin{bmatrix} 0.4 & 0.5 & 0.1 \\ 0.3 & 0.7 & 0 \\ 0.5 & 0 & 0.5 \end{bmatrix}$$

Example 2:-



next state

	$S_1$	$S_2$	$S_3$	
$S_1$	0.5	0.3	0.2	= 1 ← transition probability matrix
$S_2$	0.2	0.6	0.2	
$S_3$	0.4	0.3	0.3	

current state

Some books or websites transpose this matrix

predictions using markov chain:-

$$\vec{x}^+ = A \vec{x}$$

transition matrix  $A$   $\vec{x}$  current state matrix  
next state matrix

next state vector:-


$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

initial state  $s_0$   
only (فقط)

$$x^{(1)} = \begin{bmatrix} 0.4 & 0.5 & 0.1 \\ 0.3 & 0.7 & 0 \\ 0.5 & 0 & 0.5 \end{bmatrix}$$

### Notes:-

Note that transition probability is like conditional probability

0.4  given  $s_0$  the probability to get it again (فقط)  
is 0.4

$$(s_0 | s_0) = 0.4$$

$$(s_1 | s_0) = 0.5 \text{ and so on}$$

$\Rightarrow$  markov assumes the transition probabilities is constant along the process

$$P(B|A) = P(BA)/P(A)$$

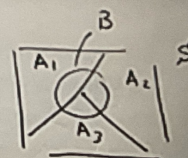
$$P(\underline{BA}) = P(B/A) P(A)$$

$$P(B) = P(BA_1) + P(BA_2) + P(BA_3)$$

$$P(\underline{B}) = \begin{matrix} \left( \begin{matrix} P(B/A_1)P(A_1) \\ + P(B/A_2)P(A_2) \\ + P(B/A_3)P(A_3) \end{matrix} \right) \end{matrix}$$

$$P(\underline{B}) = \sum_i P(B/A_i) P(A_i)$$

$$P(s_j^+) = \sum_i P(s_j^+/s_i) P(s_i)$$



$s_1^{(0)}$   
 $s_1^{(1)}$   
 $s_1^{(2)}$   
 $\vdots$

$$[P_{0,0} \dots \dots] \begin{bmatrix} s_0 \\ s_1 \\ \vdots \end{bmatrix}$$

$$P_{0,0} \equiv P(s_0^+/s_0)$$

$$P_{0,1} \equiv P(s_1^+/s_0)$$

$$P_{i,j} \equiv P(s_j^+/s_i)$$





$$P(B) = \sum_i P(B|A_i) P(A_i)$$

$$P(s_j^+) = \sum_i P(s_j^+ / s_i) P(s_i)$$

$$P(s_j^+ / s_0) P(s_0) + P(s_j^+ / s_1) P(s_1) + P(s_j^+ / s_2) P(s_2)$$

$$\begin{bmatrix} P_{0,0} & P_{0,1} & P_{0,2} \\ P_{1,0} & P_{1,1} & P_{1,2} \\ P_{2,0} & P_{2,1} & P_{2,2} \end{bmatrix} \begin{bmatrix} P(s_0) \\ P(s_1) \\ P(s_2) \end{bmatrix}$$

notations of Markov transition probabilities.

$$= P_{0,j} P(s_0) + P_{1,j} P(s_1) + P_{2,j} P(s_2) = P(s_j^+)$$

$$P = \begin{bmatrix} P_{0,0} & P_{0,1} & P_{0,2} \\ P_{1,0} & P_{1,1} & P_{1,2} \\ P_{2,0} & P_{2,1} & P_{2,2} \end{bmatrix}$$

(4)

$$P(s_0^+) = P_{0,0} P(s_0) + P_{1,0} P(s_1) + P_{2,0} P(s_2)$$

$$P(s_0^+) = P_{0,0} \times P(s_0) + P_{1,0} \times P(s_1) + P_{2,0} \times P(s_2)$$

$$\vec{X}^+ = \begin{bmatrix} P(s_0) & P(s_1) & P(s_2) \end{bmatrix} \begin{bmatrix} P_{0,0} & P_{0,1} & P_{0,2} \\ P_{1,0} & P_{1,1} & P_{1,2} \\ P_{2,0} & P_{2,1} & P_{2,2} \end{bmatrix}$$

next state «column number»

(5)

current state «row number»

$$P = \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.2 & 0.6 & 0.2 \\ 0.4 & 0.3 & 0.3 \end{bmatrix}$$

$$P = [\leftarrow P_{i,j} \rightarrow]$$

Initially

$$\vec{X}_0 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\vec{X}^{(1)T} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.2 & 0.6 & 0.2 \\ 0.4 & 0.3 & 0.3 \end{bmatrix}$$

$$\vec{X}^{(1)T} = \begin{bmatrix} 0.5 & 0.3 & 0.2 \end{bmatrix}$$

$$\vec{X}^2 = \vec{X}^{(1)T} [P] = \begin{bmatrix} 0.5 & 0.3 & 0.2 \end{bmatrix} \begin{bmatrix} \phantom{0.5} & \phantom{0.3} & \phantom{0.2} \\ \phantom{0.5} & \phantom{0.3} & \phantom{0.2} \\ \phantom{0.5} & \phantom{0.3} & \phantom{0.2} \end{bmatrix}$$

$\begin{bmatrix} x & x & x \end{bmatrix}$

# "steady state concept"

⇒ if I stayed doing this experiment for a long time will it be const or it will stay differ

⇒ steady state is about it will be const in some point  
< the probability will stay the same "vector" >

current step  
↓

next step →

$$P = \begin{bmatrix} \underbrace{P_{0,0}} & \underbrace{P_{1,0}} & \underbrace{P_{2,0}} \\ \underbrace{P_{0,1}} & \underbrace{P_{1,1}} & \underbrace{P_{2,1}} \\ \underbrace{P_{0,2}} & \underbrace{P_{1,2}} & \underbrace{P_{2,2}} \end{bmatrix}$$

1                  1                  1

• This is the transpose of the transition matrix

$$x^{(1)} = \begin{bmatrix} \end{bmatrix} = \begin{bmatrix} \end{bmatrix} \quad \begin{bmatrix} \end{bmatrix} \begin{bmatrix} \end{bmatrix} \begin{bmatrix} \end{bmatrix} \xrightarrow{x^{(0)}}$$

$$x^{(n)} = P^n x^{(0)}$$

↙ p x p x p · p n times

if p is diagonalizable.

$$P = V \Lambda V^{-1} = \begin{bmatrix} \uparrow \\ \checkmark \\ \downarrow \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_i \end{bmatrix} \begin{bmatrix} \uparrow \\ \checkmark \\ \downarrow \end{bmatrix}^{-1}$$



$$P^n \vec{x}^{(0)} = \begin{bmatrix} \uparrow \\ \sqrt{} \\ \downarrow \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} \uparrow \\ \sqrt{} \\ \downarrow \end{bmatrix}^{-1} \quad x^{(0)} = \vec{x}^{(n)}$$

steady state? will the probability remains const. after period of time?

~~means~~

$$\vec{x}^{(n)} = \underbrace{\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}}_P \vec{x}^{(n-1)} = \vec{x}^{(n-1)}$$

أول لقيم احتمالات مهما ضربت فيها الناتج هيبقى  
مثلاً الفوز 75% دائماً أو 50% أو إلى آخره

$$P x = \frac{1}{n} x \Rightarrow \text{as if a transformation matrix}$$

~ if I did eigendecomposition

$$\lambda_i = 1 \rightarrow \vec{v}_i \text{ is the steady state}$$

# Information Entropy

Claude Shannon considered as the father of modern Communication discussed the Entropy of Information

set 1  $\rightarrow$  010101010011

set 2  $\rightarrow$  0000010001000001

In our example set 1 usually 1, 0 has same amount of information because 0 & 1 have same occurrence probability but in set 2 1 has more because its probability is much less than Zero

$$I(x) \propto \frac{1}{f(p(x))}$$

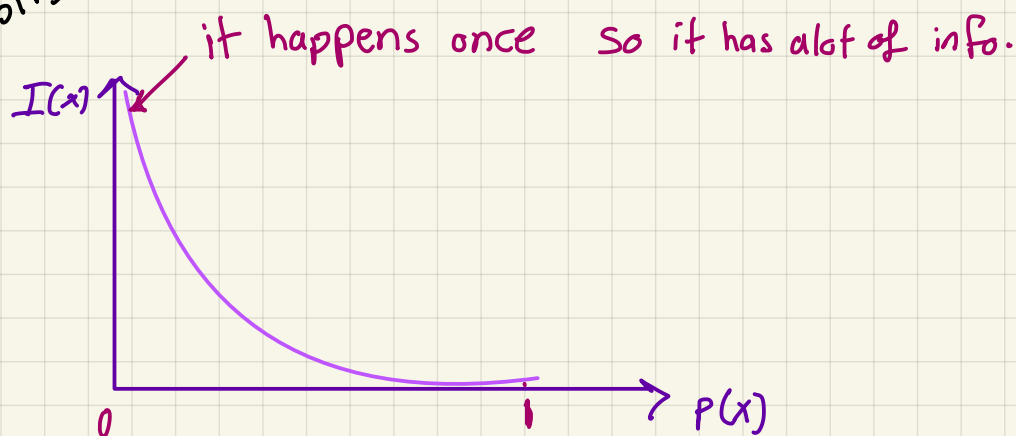
amount of information

$$\text{total} = I(x_1) + I(x_2) \\ (x_1, x_2)$$

$$I(x) = \log_2 \frac{1}{p(x)} = \text{bits} = -\log_2 p(x)$$

to let the output be in bits

probability



# Information theory

Information entropy - topic in Information theory

$$P(H) = 0.5$$

$$P(T) = 0.5$$

$$I(H) = -\log_2(0.5)$$

$$I(T) = -\log_2(0.5)$$

$$\text{average information} = \sum_i \frac{I(x_i) P(x_i)}{1}$$

$$\text{Entropy} \sim H(x) = \sum_i (-\log_2 P(x_i)) P(x_i)$$

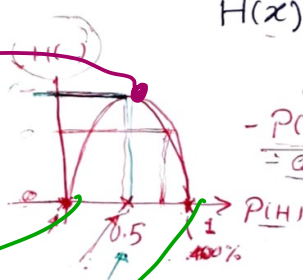
$$H(x) = -\sum_i P(x_i) \log P(x_i)$$

$$= -P(H) \log P(H) - P(T) \log P(T)$$

$$= -0.5 \log 0.5 - 0.5 \log 0.5$$

$$= -0.5 \times (-1) - 0.5 \times (-1)$$

$$= 0.5 + 0.5 = 1$$



0.5 is the avg

each time on tail

each time on head

both information = zero

because nothing new

if coin was biased the value will appear in the curve

Loss less compression

Information gain

## Cross-Entropy loss function (Classification)

$$H(P, q) = - \sum_i P(x_i) \log q(x_i)$$

نصف المقصود  
Entropy

true Prob.

ground truth

prediction prob.



$$\rightarrow \frac{0.6}{0.4} \leftarrow$$

normalize error more than L1 norm & L2 norm  
(mean square distance & mean absolute distance)