

## Information Theory

## Chapter I: Discrete information sources

# Block diagram of a communication system

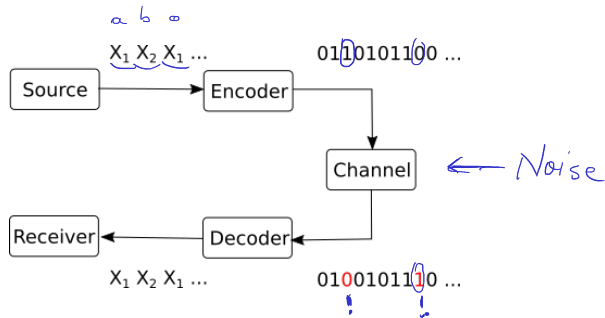
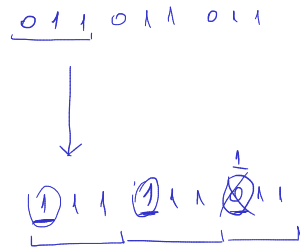


Figure 1: Block diagram of a communication system

- Source: creates information messages  $X_1 X_2 X_3$
- Encoder: converts messages into symbols for transmission (i.e bits)
- Channel: delivers the symbols, introduces errors
- Decoder: detects/corrects the errors, rebuilds the information messages



# What is information?

Example:

- ▶ Consider the sentence: “your favorite football team lost the last match”
- ▶ Does this message carry information? How, why, how much?
- ▶ Consider the following facts:
  - ▶ the message carries information only when you don't already know the result
  - ▶ if you already known the result, the message is useless (brings no information)
  - ▶ if the result was to be expected, there is little information. If the result is highly unusual, there is more information in this message (think betting)

# Information and events

- ▶ We define the notion of information for a probabilistic event
  - ▶ the happening of a probabilistic event = creation of information
- ▶ Information brought by an event depends on the probability of the event
- ▶ Rule of thumb: if you can guess something most of the times, it has little information
- ▶ Questions:
  - ▶ does a sure event ( $p = 1$ ) bring any information?  $\rightarrow i = 0$
  - ▶ does an almost sure event (e.g.  $P = \underline{0.9999}$ ) bring little or much information?  $\rightarrow i = \text{small}$
  - ▶ does a rare event (e.g.  $P = \underline{0.0001}$ ) bring a little or much information?  $\rightarrow i = \text{large}$

# Information

- ▶ The information attached to a particular event (known as “message”)  $s_i$  is rigorously defined as:

$$i(s_i) = -\log_2(p(s_i))$$

$$i(\text{something}) = -\log_2(P(\text{something}))$$

- ▶ Properties:

- ▶  $i(s_i) \geq 0$
- ▶ lower probability (rare events) means higher information
- ▶ higher probability (frequent events) means lower information
- ▶ a certain event brings no information:  $-\log(1) = 0$
- ▶ an event with probability 0 brings infinite information (but it never happens...)
- ▶ for two independent events, their information gets added

$$p(111) = 0.0001$$

$$i(111) = -\log_2(0.0001) = 13.2 \text{ bits}$$

$$-\log_2\left(\frac{1}{36}\right) = -\log_2\frac{1}{6} + \log_2\frac{1}{6}$$

$$i(p(s_i) \cdot p(s_j)) = i(s_i) + i(s_j)$$

# The choice of logarithm

- ▶ Any base of logarithm can be used in the definition.
- ▶ Usual convention: use binary logarithm  $\log_2()$
- ▶ In this case, the information  $i(s_i)$  is measured in bits
- ▶ If using natural logarithm  $\ln()$ , it is measured in nats.
- ▶ Logarithm bases can be converted to/from one another:

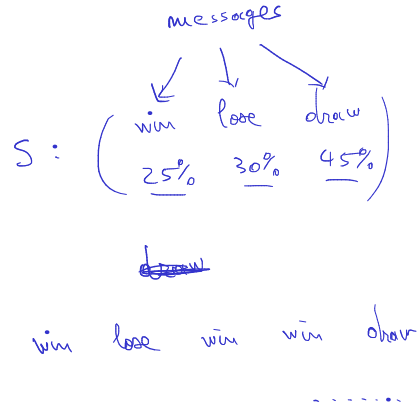
$$\log_b(x) = \frac{\log_a(x)}{\log_a(b)}$$

- ▶ Information defined using different logarithms differ only in scaling:

$$i_b(s_i) = \frac{i_a(s_i)}{\log_a(b)}$$

# Information source

- ▶ A probabilistic event is always part of a set of multiple events (options)
  - ▶ e.g: a football team can win/lose/draw a match (3 possible events)
  - ▶ each event has a certain probability. All probabilities are known beforehand
  - ▶ at a given time, only one of the events can happen
- ▶ An **information source** = the set of all events together with their probabilities
- ▶ One event is called a message
- ▶ Each message carries the information that **it** happened, the quantity of information is dependent on its probability





# Sequence of messages

- ▶ An information source creates a **sequence of messages**
  - ▶ e.g. like throwing a coin or a dice several times in a row
- ▶ The probabilities of the messages are known and fixed
- ▶ Each time, a new message is randomly selected according to the probabilities

# Discrete memoryless source

- ▶ A discrete memoryless source (DMS) is an <sup>information</sup> information source which produces a sequence of independent messages
  - ▶ i.e. the choice of a message at one time does not depend on the previous messages
- ▶ Each message has a fixed probability. The set of probabilities is the distribution of the source:

$$S : \begin{pmatrix} s_1 & s_2 & s_3 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \end{pmatrix}$$

i.e. = "that is"  
e.g. = "for example"

# Discrete memoryless source

$$S : \begin{pmatrix} s_1 & s_2 & s_3 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \end{pmatrix} \leftarrow \text{discrete set}$$

## ► Terminology:

- Discrete: it can take a value from a discrete set (“alphabet”)
  - Complete:  $\sum p(s_i) = 1$
  - Memoryless: successive values are independent of previous values (e.g. successive throws of a coin)
- A message from a DMS is also called a random variable in probabilistics.

## Examples

- ▶ A coin is a discrete memoryless source (DMS) with two messages:

$$S : \begin{pmatrix} heads & tails \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

- ▶ A dice is a discrete memoryless source (DMS) with six messages:

$$S : \begin{pmatrix} s_1 & s_2 & s_3 & s_4 & s_5 & s_6 \\ \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \end{pmatrix}$$

- ▶ Playing the lottery can be modeled as DMS:

$$S : \begin{pmatrix} s_1 & s_2 \\ 0.9999 & 0.0001 \end{pmatrix}$$

# Examples

- ▶ An extreme type of DMS containing the certain event:

$$S : \begin{pmatrix} s_1 & s_2 \\ 1 & 0 \end{pmatrix}$$

- ▶ Receiving an unknown *bit* (0 or 1) with equal probabilities:

$$S : \begin{pmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix} \rightarrow H(S) = 1 \text{ bit}$$

## Sequence of messages from DMS

- ▶ A DMS produces a sequence of messages by randomly selecting a message every time, with the same fixed probabilities
  - ▶ throwing a dice several times in a row you can get a sequence  
4, 2, 3, 2, 1, 6, 1, 5, 4, 5, . . . . .
- ▶ If the sequence is very long (has  $N$  messages,  $N$  very large), each message  $s_i$  appears approximately  $p(s_i) * N$  times in the sequence
  - ▶ gets more precise as  $N \rightarrow \infty$

# Entropy of a DMS

- ▶ We usually don't care about a single message. We are interested in long sequences of messages (think millions of bits of data)
- ▶ We are interested in the *average* information of a message from a DMS
- ▶ Definition: the **entropy** of a DMS source  $S$  is **the average information of a message**:

$$H(S) = \sum_k p(s_k) i(s_k) = - \sum_k p(s_k) \log_2(p_k)$$

where  $p(s_k)$  is the probability of message  $k$

$$(0.6 \cdot 6_1 + 0.4 \cdot 6_2)$$

$$S : \begin{pmatrix} \Delta_1 & \Delta_2 \\ 3/4 & 1/4 \end{pmatrix}$$

$$i(\Delta_1) = -\log_2\left(\frac{3}{4}\right) = 0.41 \text{ b}$$

$$i(\Delta_2) = -\log_2\left(\frac{1}{4}\right) = 2 \text{ b}$$

$$\Delta_1 \Delta_2 \Delta_1 \Delta_1 \Delta_2 \Delta_1 \Delta_2 \Delta_2 \Delta_1 \Delta_1 \dots$$

bits/message

Average info. of one message:

$$\bar{i} = \frac{3}{4} \cdot 0.41 + \frac{1}{4} \cdot 2$$

$$= -\sum p(\Delta_i) \cdot \log_2(p(\Delta_i))$$

# Entropy of a DMS

- ▶ Since information of a message is measured in bits, entropy is measured in bits (or **bits / message**, to indicate it is an average value)
- ▶ Entropies using information defined with different logarithms differ only in scaling:

$$H_b(S) = \frac{H_a(S)}{\log_a(b)} \leftarrow \text{constant}$$



# Examples

$$-\log \alpha = \log \frac{1}{\alpha}$$

$$\text{Coin: } \begin{pmatrix} \overset{0}{H} & \overset{1}{T} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

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

$$\log_2 \frac{1}{2} = \log_2 2^{-1} = -1$$

$$H(\text{coin}) = -\frac{1}{2} \underbrace{\log_2 \frac{1}{2}}_{-1} - \frac{1}{2} \underbrace{\log_2 \frac{1}{2}}_{-1} = \frac{1}{2} + \frac{1}{2} = 1 \text{ bit}$$

► Coin:  $H(S) = 1 \text{ bit/message}$

► Dice:  $H(S) = \log(6) \text{ bits/message} \rightarrow$

► Lottery:  $H(S) = -0.9999 \log(0.9999) - 0.0001 \log(0.0001) = \text{very small}$

► Receiving 1 bit:  $H(S) = 1 \text{ bit/message}$  (hence the name!)

$$S: \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \end{pmatrix}$$



$$\begin{aligned} H(S) &= -\frac{1}{6} \log \frac{1}{6} - \frac{1}{6} \log \frac{1}{6} - \dots - \frac{1}{6} \log \frac{1}{6} \\ &= -\cancel{1} \cdot \frac{1}{\cancel{6}} \log \frac{1}{6} = -\log \frac{1}{6} = \log 6 = \underline{\underline{2.58}} \end{aligned}$$

# Interpretation of the entropy

$$S: \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \end{pmatrix}$$

All the following interpretations of entropy are true:

- ▶  $H(S)$  is the average uncertainty of the source  $S$
  - ▶  $H(S)$  is the average information of the messages from source  $S$
  - ▶ A long sequence of  $N$  messages from  $S$  has total information  $\approx N \cdot H(S)$
- ▶  $H(S)$  is the minimum number of bits (0,1) required to uniquely represent an average message from source  $S$

$$\Delta_1 \Delta_2 \Delta_6 \Delta_2 \Delta_1 \dots \dots$$

$\underbrace{\hspace{10em}}_{N \text{ messages}}$

$$I =$$

# Properties of entropy

$$\log_2 x' = \frac{\ln x'}{\ln 2}$$

$$H(S) = - \sum_k p_k \cdot \underbrace{\log(p_k)}_{\leq 0} \geq 0$$

$$p_k \in [0, 1]$$

We prove the following **properties of entropy**:

1.  $H(S) \geq 0$  (non-negative)

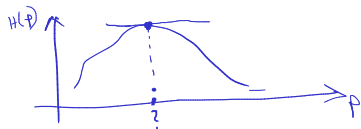
Proof: via definition

2.  $H(S)$  is maximum when all  $n$  messages have equal probability  $\frac{1}{n}$ . The maximum value is  $\max H(S) = \log(n)$

Proof: only for the case of 2 messages, use derivative in definition

3. Diversification of the source always increases the entropy

Proof: compare entropies in both cases



$$S: \begin{pmatrix} \Delta_1 & \Delta_2 & \Delta_3 \\ p_1 & p_2 & p_3 \\ 1/3 & 1/3 & 1/3 \end{pmatrix} \quad p_1 + p_2 + p_3 = 1$$

$$S: \begin{pmatrix} \Delta_1 & \Delta_2 \\ p & 1-p \end{pmatrix}$$

$$H(S) = -p \log p - (1-p) \log(1-p)$$

Want maximum  $\uparrow$

$$\frac{dH(S)}{dp} = 0 \quad \text{TODO next week}$$

$$-\log p - \cancel{\frac{p}{p \cdot \ln 2}} + \log(1-p) + \cancel{(1-p) \frac{1}{(1-p) \ln 2}} = 0$$

$$-\log p + \log 1-p = 0 \Leftrightarrow \log_2 \frac{1-p}{p} = 0 \Leftrightarrow \frac{1-p}{p} = 1 \Leftrightarrow 1-p = p$$

$$\log_2 \frac{1-p}{p} = 2^0$$

$$\Leftrightarrow p = \frac{1}{2}$$

# The entropy of a binary source

- Consider a general DMS with two messages:

$$S: \begin{pmatrix} s_1 & s_2 \\ p & 1-p \end{pmatrix}$$

- It's entropy is:

$$H(S) = -p \cdot \log(p) - (1-p) \cdot \log(1-p)$$

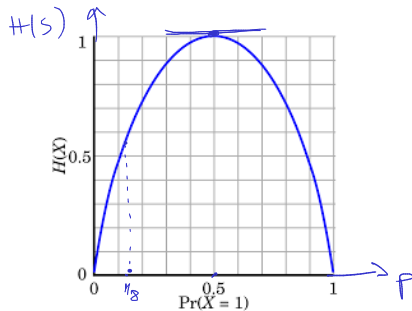


Figure 2: Entropy of a binary source

Prove 3<sup>rd</sup> property:

$$3) S: \begin{pmatrix} s_1 & s_2 & \dots & s_M \\ p_1 & p_2 & \dots & p_M \end{pmatrix}$$

$$H(S_d) \geq H(S)$$

$$S_d: \begin{pmatrix} s_1 & s_2 & \dots & s_{M+1} & s_{M+2} \\ p_1 & p_2 & \dots & p_{M+1} & p_{M+2} \end{pmatrix}$$

$$p_M = p_{M+1} + p_{M+2}$$

$$H(S) = -p_1 \log_2 p_1 - \dots$$

$$-p_M \log p_M$$

$$H(S_d) = -p_1 \log_2 p_1 - \dots$$

$$-p_{M+1} \log p_{M+1} - p_{M+2} \log p_{M+2}$$

$$H(S_d) - H(S) = -p_{M+1} \log p_{M+1} - p_{M+2} \log p_{M+2} + (p_{M+1} + p_{M+2}) \log (p_{M+1} + p_{M+2})$$

$$= p_{M+1} \cdot (\log(p_{M+1} + p_{M+2}) - \log p_{M+1}) + p_{M+2} \cdot (\log(p_{M+1} + p_{M+2}) - \log p_{M+2})$$

$$= \underbrace{p_{M+1} \cdot \log \frac{p_{M+1} + p_{M+2}}{p_{M+1}}}_{\geq 0} + \underbrace{p_{M+2} \cdot \log \frac{p_{M+1} + p_{M+2}}{p_{M+2}}}_{\geq 0} \geq 0 \quad \text{g.e.d.}$$

## Example - Game

$$S: \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1/8 & 1/8 & 1/8 & 1/8 & 1/8 & 1/8 & 1/8 & 1/8 \end{pmatrix}$$

$$H(S) = \cancel{8} \log_2 \frac{1}{8} = \log_2 8 = 3 \text{ bits}$$

Game: I think of a number between 1 and 8. You have to guess it by asking yes/no questions.

- ▶ How much uncertainty does the problem have?
- ▶ How is the best way to ask questions? Why?
- ▶ What if the questions are not asked in the best way?
- ▶ On average, what is the number of questions required to find the number?

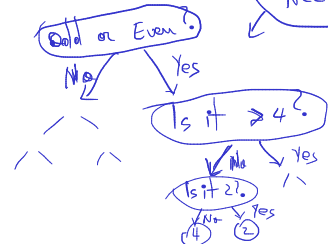
1) Is it odd or even?

$$\text{Answer: } \begin{pmatrix} \text{Yes} & \text{No} \\ 1/2 & 1/2 \end{pmatrix} \rightarrow H(\text{Answer}) = 1 \text{ bit}$$

Decision Tree

Is it 1?

$$\text{Answer: } \begin{pmatrix} \text{Yes} & \text{No} \\ 1/8 & 7/8 \end{pmatrix} \Rightarrow H(\text{Answer}) < 1$$



## Example - Game v2

- Suppose I choose a number according to the following distribution:

$$H(S) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - 2 \frac{1}{8} \log \frac{1}{8}$$

$$= \frac{1}{2} + \frac{1}{2} + \frac{3}{4} = \frac{7}{4} \text{ bits}$$

$$= 1.75$$

$$S: \begin{pmatrix} s_1 & s_2 & s_3 & s_4 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{8} & \frac{1}{8} \end{pmatrix}$$

Is it > 2  
 Yes No  
 6/8 2/8

Optimal Decision Tree?

Is it 1?

No 1/2

Yes 1/2  
 1

Is it 2?

No 1/2

Yes 1/2

Is it 3?  
 2

No

Yes

4

3

$$S: \begin{pmatrix} s_1 & s_2 & s_3 & s_4 \\ 0.14 & 0.29 & 0.4 & 0.17 \end{pmatrix}$$

Is it 1 or 3

Yes

No

Yes

No

Yes

No

- In general:

- What distribution makes guessing the number the most difficult?
- What distribution makes guessing the number the easiest?

# Efficiency and redundancy

- ▶ Efficiency of a DMS:

$$\eta = \frac{H(S)}{H_{max}} = \frac{H(S)}{\log(n)}$$

- ▶ Absolute redundancy of a DMS:

$$R = H_{max} - H(S)$$

- ▶ Relative redundancy of a DMS:

$$\rho = \frac{H_{max} - H(S)}{H_{max}} = 1 - \eta$$

# Information flow of a DMS

- ▶ Suppose that message  $s_i$  takes time  $t_i$  to be transmitted via some channel.
- ▶ Definition: the **information flow** of a DMS  $S$  is the average information transmitted per unit of time:

$$H_\tau(S) = \frac{H(S)}{\bar{t}} \leftarrow \text{average transmission time of messages}$$

where  $\bar{t}$  is the average duration of transmitting a message:

$$\bar{t} = \sum_i p_i t_i$$

- ▶ Measured in **bps** (bits per second)
- ▶ Important for data communication



# Distance between distributions

- ▶ How to measure how similar / how different are two distributions?

- ▶ must have the same number of messages
- ▶ example:  $p(s_1), \dots, p(s_n)$  and  $q(s_1), \dots, q(s_n)$

- ▶ **Definition:** the **Kullback–Leibler distance** of two distributions P and Q is

$$D_{KL} = D_{KL}(P||Q) = \sum_i p(s_i) \log\left(\frac{p(s_i)}{q(s_i)}\right)$$

$$\begin{matrix} (p_1, p_2, \dots, p_n) \\ (q_1, q_2, \dots, q_n) \end{matrix}$$

- ▶ It is a way to measure the distance (difference) between two distributions
- ▶ Also known as *relative entropy*, or the Kullback-Leibler *divergence*

$$\begin{aligned} S_1: & \begin{pmatrix} \Delta_1 & \Delta_2 & \Delta_3 \\ 0.5 & 0.2 & 0.3 \end{pmatrix} \\ S_2: & \begin{pmatrix} \Delta_1 & \Delta_2 & \Delta_3 \\ 0.51 & 0.18 & 0.31 \end{pmatrix} \\ S_3: & \begin{pmatrix} \Delta_1 & \Delta_2 & \Delta_3 \\ 0.52 & 0.19 & 0.3 \end{pmatrix} \end{aligned}$$
$$= p_1 \log \frac{p_1}{q_1} + p_2 \log \frac{p_2}{q_2} + \dots$$

# Properties of Kullback-Leibler distance

- ▶ Properties:
  - ▶  $D_{KL}(P||Q)$  is always  $\geq 0$ , and is equal to 0 only when P and Q are the same
  - ▶ the higher  $D_{KL}(P||Q)$  is, the more different the distributions are
  - ▶ it is **not commutative**:  $D_{KL}(P||Q) \neq D_{KL}(Q||P)$
- ▶ Example: at whiteboard
- ▶ Example usage: classification systems (cross-entropy loss)

# Extended DMS

- Definition: the **n-th order extension** of a DMS  $S$ ,  $S^n$  is a source which has as messages all the combinations of  $n$  messages of  $S$ :

$$\sigma_i = \underbrace{s_j s_k \dots s_l}_n$$

$$S : \begin{pmatrix} \Lambda_1 & \Lambda_2 & \Lambda_3 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \end{pmatrix}$$

$$S^2 : \begin{pmatrix} \Lambda_1 \Lambda_1 & \Lambda_1 \Lambda_2 & \Lambda_1 \Lambda_3 & \Lambda_2 \Lambda_1 & \Lambda_2 \Lambda_2 & \Lambda_2 \Lambda_3 & \Lambda_3 \Lambda_1 & \Lambda_3 \Lambda_2 & \Lambda_3 \Lambda_3 \\ \frac{1}{4} & \frac{1}{8} & \frac{1}{8} & \frac{1}{8} & \frac{1}{16} & \frac{1}{16} & \frac{1}{8} & \frac{1}{16} & \frac{1}{16} \end{pmatrix}$$

- If  $S$  has  $k$  messages,  $S^n$  has  $k^n$  messages
- Since  $S$  is DMS, probabilities multiply:

$$p(\sigma_i) = p(s_j) \cdot p(s_k) \cdot \dots \cdot p(s_l)$$

~

## Extended DMS - Example

► Examples:

$$S : \begin{pmatrix} s_1 & s_2 \\ \frac{1}{4} & \frac{3}{4} \end{pmatrix}$$

$$S^2 : \begin{pmatrix} \sigma_1 = s_1 s_1 & \sigma_2 = s_1 s_2 & \sigma_3 = s_2 s_1 & \sigma_4 = s_2 s_2 \\ \frac{1}{16} & \frac{3}{16} & \frac{3}{16} & \frac{9}{16} \end{pmatrix}$$

$$S^3 : \begin{pmatrix} s_1 s_1 s_1 & s_1 s_1 s_2 & s_1 s_2 s_1 & s_1 s_2 s_2 & s_2 s_1 s_1 & s_2 s_1 s_2 & s_2 s_2 s_1 & s_2 s_2 s_2 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

# Extended DMS - Another example

$$18 \cdot H(s) = 9 \cdot H(s^2) = 2.25 \cdot H(s^8)$$

- Long sequence of binary messages:

010011001110010100...

18 messages from  $\rightarrow S : \begin{pmatrix} 0 & 1 \\ 1/2 & 1/2 \end{pmatrix} H(s) = 1 \text{ bit}$

9 messages from  $\rightarrow S^2 : \begin{pmatrix} 00 & 01 & 10 & 11 \\ . & . & . & . \end{pmatrix}$

messages from  $\rightarrow S^8 : \begin{pmatrix} 00000000 & 00000001 & 00000010 & 00000011 & 00000100 & 00000101 & 00000110 & 00000111 & 00001000 & 00001001 & 00001010 & 00001011 & 00001100 & 00001101 & 00001110 & 00001111 & 00010000 & 00010001 & 00010010 & 00010011 & 00010100 & 00010101 & 00010110 & 00010111 & 00011000 & 00011001 & 00011010 & 00011011 & 00011100 & 00011101 & 00011110 & 00011111 & 00100000 & 00100001 & 00100010 & 00100011 & 00100100 & 00100101 & 00100110 & 00100111 & 00101000 & 00101001 & 00101010 & 00101011 & 00101100 & 00101101 & 00101110 & 00101111 & 00110000 & 00110001 & 00110010 & 00110011 & 00110100 & 00110101 & 00110110 & 00110111 & 00111000 & 00111001 & 00111010 & 00111011 & 00111100 & 00111101 & 00111110 & 00111111 & 01000000 & 01000001 & 01000010 & 01000011 & 01000100 & 01000101 & 01000110 & 01000111 & 01001000 & 01001001 & 01001010 & 01001011 & 01001100 & 01001101 & 01001110 & 01001111 & 01010000 & 01010001 & 01010010 & 01010011 & 01010100 & 01010101 & 01010110 & 01010111 & 01011000 & 01011001 & 01011010 & 01011011 & 01011100 & 01011101 & 01011110 & 01011111 & 01100000 & 01100001 & 01100010 & 01100011 & 01100100 & 01100101 & 01100110 & 01100111 & 01101000 & 01101001 & 01101010 & 01101011 & 01101100 & 01101101 & 01101110 & 01101111 & 01110000 & 01110001 & 01110010 & 01110011 & 01110100 & 01110101 & 01110110 & 01110111 & 01111000 & 01111001 & 01111010 & 01111011 & 01111100 & 01111101 & 01111110 & 01111111 & 10000000 & 10000001 & 10000010 & 10000011 & 10000100 & 10000101 & 10000110 & 10000111 & 10001000 & 10001001 & 10001010 & 10001011 & 10001100 & 10001101 & 10001110 & 10001111 & 10010000 & 10010001 & 10010010 & 10010011 & 10010100 & 10010101 & 10010110 & 10010111 & 10011000 & 10011001 & 10011010 & 10011011 & 10011100 & 10011101 & 10011110 & 10011111 & 10100000 & 10100001 & 10100010 & 10100011 & 10100100 & 10100101 & 10100110 & 10100111 & 10101000 & 10101001 & 10101010 & 10101011 & 10101100 & 10101101 & 10101110 & 10101111 & 10110000 & 10110001 & 10110010 & 10110011 & 10110100 & 10110101 & 10110110 & 10110111 & 10111000 & 10111001 & 10111010 & 10111011 & 10111100 & 10111101 & 10111110 & 10111111 & 11000000 & 11000001 & 11000010 & 11000011 & 11000100 & 11000101 & 11000110 & 11000111 & 11001000 & 11001001 & 11001010 & 11001011 & 11001100 & 11001101 & 11001110 & 11001111 & 11010000 & 11010001 & 11010010 & 11010011 & 11010100 & 11010101 & 11010110 & 11010111 & 11011000 & 11011001 & 11011010 & 11011011 & 11011100 & 11011101 & 11011110 & 11011111 & 11100000 & 11100001 & 11100010 & 11100011 & 11100100 & 11100101 & 11100110 & 11100111 & 11101000 & 11101001 & 11101010 & 11101011 & 11101100 & 11101101 & 11101110 & 11101111 & 11110000 & 11110001 & 11110010 & 11110011 & 11110100 & 11110101 & 11110110 & 11110111 & 11111000 & 11111001 & 11111010 & 11111011 & 11111100 & 11111101 & 11111110 & 11111111 \end{pmatrix}$

- Can be grouped in bits, half-bytes, bytes, 16-bit words, 32-bit long words, and so on
- Can be considered:
  - N messages from a binary source (with 1 bit), or
  - N/2 messages from a source with 4 messages (with 2 bits)...
  - etc

$$H(S^8) = 8 \cdot H(s) = 8 \text{ bits}$$

# Property of DMS

Proof:  $S = \begin{pmatrix} \Lambda_1 & \Lambda_2 & \dots & \Lambda_N \\ p_1 & p_2 & \dots & p_N \end{pmatrix}$   
 $S^n = \begin{pmatrix} \Lambda_1 & \Lambda_2 & \dots & \Lambda_{N^n} \\ p(\Lambda_1) & p(\Lambda_2) & \dots & p(\Lambda_{N^n}) \end{pmatrix}$

$$H(S^n) = - \sum_{i=1}^{N^n} p(\Lambda_i) \cdot \log(p(\Lambda_i))$$

$$= - \sum_{j=1}^N \underbrace{\sum_{k=1}^N \dots \sum_{l=1}^N}_{n \text{ terms}} (\underbrace{p_j \cdot p_k \cdot \dots \cdot p_l}_n) \cdot \log(\underbrace{p_j \cdot p_k \cdot \dots \cdot p_l}_{(\log p_j + \log p_k + \dots + \log p_l)})$$

- Theorem: The entropy of a  $n$ -th order extension is  $n$  times larger than the entropy of the original DMS

$$H(S^n) = nH(S)$$

- Interpretation: grouping messages from a long sequence in blocks of  $n$  does not change total information (e.g. groups of 8 bits = 1 byte)

$$= - \sum_{j=1}^N \sum_{k=1}^N \dots \sum_{l=1}^N (\underbrace{p_j \cdot p_k \cdot \dots \cdot p_l}_n) \cdot \log p_j + \dots + \sum_{j=1}^N \sum_{k=1}^N \dots \sum_{l=1}^N (\underbrace{p_j \cdot p_k \cdot \dots \cdot p_l}_n) \cdot \log p_k + \dots + \sum_{j=1}^N \sum_{k=1}^N \dots \sum_{l=1}^N (\underbrace{p_j \cdot p_k \cdot \dots \cdot p_l}_n) \cdot \log p_l =$$

$$= n \cdot H(S)$$

$$= - \sum_j p_j \log p_j - \sum_k p_k \log p_k - \sum_l p_l \log p_l - \dots - \sum_j p_j \log p_j - \sum_k p_k \log p_k - \sum_l p_l \log p_l - \dots - \sum_j p_j \log p_j - \sum_k p_k \log p_k - \sum_l p_l \log p_l - \dots$$

$H(S)$   $H(S)$   $H(S)$   $H(S)$   $H(S)$

# An example [memoryless is not enough]

- ▶ The distribution (frequencies) of letters in English:

letter	probability	letter	probability
A	.082	N	.067
B	.015	O	.075
C	.028	P	.019
D	.043	Q	.001
E →	.127	R	.060
F	.022	S	.063
G	.020	T	.091
H	.061	U	.028
I	.070	V	.010
J	.002	W	.023
K	.008	X	.001
L	.040	Y	.020
M	.024	Z	.001

— \_ o ? — —

- ▶ Text from a memoryless source with these probabilities:

→ OCRO HLI RGWR NMIELWIS EU LL NBNESEBYA TH EEI  
ALHENHTTPA OOBTTVA NAH BRL

(taken from *Elements of Information Theory*, Cover, Thomas)

- ▶ What's wrong? **Memoryless**

Handwritten notes illustrating the concept of memoryless sources:

- A bracket groups the following equations:
$$\begin{cases} l = \cancel{.95} \\ p = \\ s = \end{cases}$$
- Below the bracket, the sequence  $a \ c \ c \ e \ -$  is written.
- To the right, the sequence  $- \ l \ e \ e \ -$  is written.
- Below these, the sequence  $- \ l \ e \ \textcircled{-} \ -$  is written, with a question mark above the circled dash.
- Further down, the sequence  $- \ l \ e \ \textcircled{-} \ -$  is written again, with a question mark above the circled dash.

# Sources with memory

- **Definition:** A source has **memory of order  $m$**  if the probability of a message depends on the last  $m$  messages.
- The last  $m$  messages = the **state** of the source (notation  $S_i$ ).
- A source with  $n$  messages and memory  $m \Rightarrow$  has  $n^m$  states in all.
- For every state, messages can have a different set of probabilities.  
Notation:  $p(s_i|S_k) =$  "probability of  $s_i$  in state  $S_k$ ".
- Also known as Markov sources.

$$S : \begin{pmatrix} \Delta_1 & \Delta_2 \end{pmatrix}$$

memory = 2

How many states?

~~Last 2~~

$$\begin{aligned} \Delta_1 \Delta_1 &= S_1 \\ \Delta_1 \Delta_2 &= S_2 \\ \Delta_2 \Delta_1 &= S_3 \\ \Delta_2 \Delta_2 &= S_4 \end{aligned}$$

Last 2 messages = state

state old

$$P(S_3 | S_4) = P(\Delta_1 | S_4)$$

$\Delta_2 \Delta_2 \Delta_1$   
 $\underbrace{\hspace{1cm}}_{S_4} \quad S_3$

$\Delta_1 \Delta_2 \Delta_2$   
 $\underbrace{\hspace{1cm}}_{S_2} \quad S_4$

$$P(\Delta_2 | S_2) \Leftrightarrow P(S_4 | S_2)$$



# Example

- ▶ A source with  $n = 4$  messages and memory  $m = 1$

- ▶ if the source is in state  $S_1$ , choose next message with distribution

States

state  $S_1$ :  $\begin{pmatrix} s_1 & s_2 & s_3 & s_4 \\ 0.4 & 0.3 & 0.2 & 0.1 \end{pmatrix} \Rightarrow H(S_1) = -\sum p_k \cdot \log p_k$

- ▶ if last message was  $s_2$ , choose next message with distribution

state  $S_2$ :  $\begin{pmatrix} s_1 & s_2 & s_3 & s_4 \\ 0.33 & 0.37 & 0.15 & 0.15 \end{pmatrix}$

- ▶ if last message was  $s_3$ , choose next message with distribution

state  $S_3$ :  $\begin{pmatrix} s_1 & s_2 & s_3 & s_4 \\ 0.2 & 0.35 & 0.41 & 0.04 \end{pmatrix}$

- ▶ if last message was  $s_4$ , choose next message with distribution

state  $S_4$ :  $\begin{pmatrix} s_1 & s_2 & s_3 & s_4 \\ 0.1 & 0.2 & 0.3 & 0.4 \end{pmatrix}$

$$S = \begin{pmatrix} \Lambda_1 & \Lambda_2 & \Lambda_3 & \Lambda_4 \end{pmatrix}$$

$\Lambda_1$   
 $\Lambda_2$   
 $\Lambda_3$   
 $\Lambda_4$

"prob. of  $\Lambda_3$  if source is in state  $S_2$ "

$$P(\Lambda_3 | S_2)$$

"conditioned by"

$$P(S_1 | S_1) = 0.4$$

$$P(S_2 | S_1) = 0.3$$

# Transitions

- ▶ When a new message is provided, the source transitions to a new state:

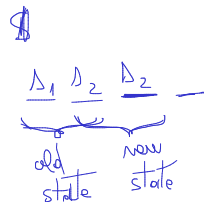
$$\dots \underbrace{s_i s_j s_k}_{\text{old state}} \underline{s_l}$$

old state

$$\dots s_i \underbrace{s_j s_k s_l}_{\text{new state}}$$

new state

- ▶ The message probabilities = the probabilities of transitions from some state  $S_u$  to another state  $S_v$



# Transition matrix

- ▶ The transition probabilities are organized in a transition matrix  $[T]$

$$[T] = \begin{matrix} \begin{matrix} \text{old state} & \text{new state} \end{matrix} & \begin{matrix} p_{11} & p_{12} & \dots & p_{1N} \\ p_{21} & p_{22} & \dots & p_{2N} \\ \dots & \dots & \dots & \dots \\ p_{N1} & p_{N2} & \dots & p_{NN} \end{matrix} \end{matrix}$$

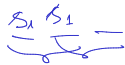
$$P_{kl} = P(S_l | S_k)$$

- ▶  $p_{ij}$  is the transition probability from state  $S_i$  to state  $S_j$
- ▶  $N$  is the total number of states

$$= P(S_j | S_i)$$

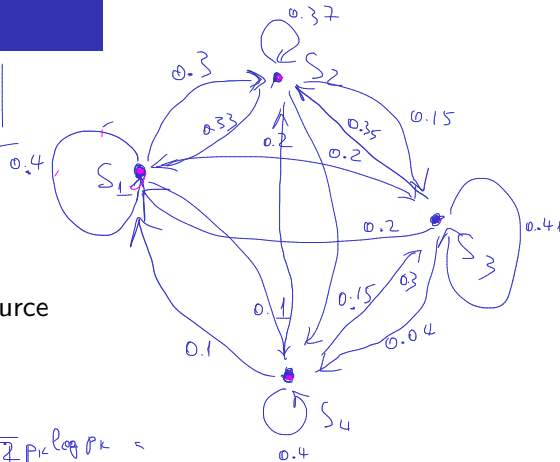
# Graphical representation

acce -  
 p 100  
 l 101  
 u 110  
 : 11



$S_1 S_2$

Sum of output transitions = 1



At whiteboard: draw states and transitions for previous example (source with  $n = 4$  messages and memory  $m = 1$ )

new

	$S_1$	$S_2$	$S_3$	$S_4$	
$S_1$	0.4	0.3	0.2	0.1	$\rightarrow H(S_1) = -\sum p_k \log p_k =$
$S_2$	0.33	0.37	0.15	0.15	$\rightarrow H(S_2) =$
$S_3$	0.2	0.35	0.41	0.04	$\rightarrow H(S_3) =$
$S_4$	0.1	0.2	0.3	0.4	$\rightarrow H(S_4) =$

4x4



# Entropy of sources with memory

- ▶ What entropy does <sup>the</sup> source with memory have?
- ▶ Each state  $S_k$  has a different distribution  $\rightarrow$  each state has a different entropy  $H(S_k)$

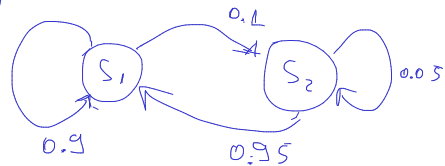
$$\underline{H(S_k)} = - \sum_i \underbrace{p(s_i|S_k)} \cdot \log(\underbrace{p(s_i|S_k)})$$

- ▶ Global entropy = average entropy

$$\underline{H(S)} = \sum_k \underbrace{p_k}_{\text{circled}} \underbrace{H(S_k)}$$

where  $p_k$  = probability that the source is in state  $S_k$

- ▶ (i.e. after a very long sequence of messages, the fraction of time when the source was in state  $S_k$ )



$$p_1 \approx 90\%$$

$$p_2 = 10\%$$

$$H(S_1)$$

$$H(S_2)$$

$$H(S) = 90\% \frac{H(S_1)}{10\%} + 10\% \frac{H(S_2)}{10\%}$$

# Ergodic sources

- ▶ How to find out the weights  $p_k$ ?
- ▶ They are known as the **stationary probabilities**
- ▶  $p_k$  = probability that the source is in state  $S_i$ , after running for a very long time
  - ▶ (i.e. after a very long sequence of messages, the fraction of time when the source was in state  $S_k$ )
- ▶ We need to answer the following question:  
If we know the state  $S_k$  at time  $n$ , what will be the state at time  $n + 1$ ?

# Ergodic sources

- ▶ Let  $p_i^{(n)}$  = the probability that source  $S$  is in state  $S_i$  at time  $n$ .
- ▶ In what state will it be at time  $n + 1$ ? (after one more message)
  - ▶ i.e. what are the probabilities of the states at time  $n + 1$ ?

- ▶ Just multiply with  $T$

$$\underbrace{[p_1^{(n)}, p_2^{(n)}, \dots, p_N^{(n)}]}_{\text{at time } n} \cdot [T] = \underbrace{[p_1^{(n+1)}, p_2^{(n+1)}, \dots, p_N^{(n+1)}]}_{\text{at time } n+1}$$

- ▶ After one more message:

$$\underbrace{[p_1^{(n)}, p_2^{(n)}, \dots, p_N^{(n)}] \cdot [T] \cdot [T]} = [p_1^{(n+2)}, p_2^{(n+2)}, \dots, p_N^{(n+2)}]$$

- ▶ For every new moment of time, one more multiplication with  $T$

$$\frac{P(S_1 \text{ after 1 step})}{P(S_2)} = P_1 \cdot P_{11} + P_2 \cdot P_{21} + P_3 \cdot P_{31} + P_4 \cdot P_{41}$$

~~$S_2$~~

$n : S_2$

$$n+1 : \begin{Bmatrix} S_1 & S_2 & S_3 & S_4 \\ 0.33 & 0.37 & 0.15 & 0.15 \end{Bmatrix}$$

$$\underbrace{[0 \ 1 \ 0 \ 0]}_{\text{time } n} \cdot \begin{bmatrix} T \\ T \\ T \\ T \end{bmatrix} = \underbrace{[0.33 \ 0.37 \ 0.15 \ 0.15]}_{\text{time } n+1}$$

$$[P_1 \ P_2 \ P_3 \ P_4] \cdot \begin{bmatrix} S_1 \\ S_2 \\ S_3 \\ S_4 \end{bmatrix}$$

# Ergodic sources

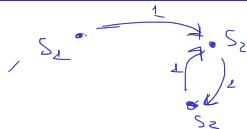
- In general, starting from time 0, after  $n$  messages the probabilities that the source is in a certain state are:

$$[p_1^{(0)}, p_2^{(0)}, \dots, p_N^{(0)}] \cdot [T]^n = [p_1^{(n)}, p_2^{(n)}, \dots, p_N^{(n)}]$$



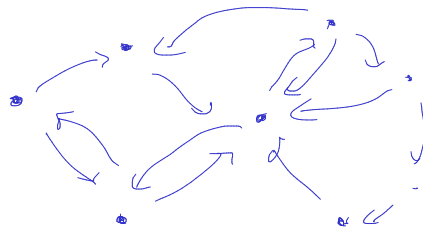


# Ergodicity



- ▶ A source is called ergodic if every state can be reached from every state, in a finite number of steps.
- ▶ Property of ergodic sources:
  - ▶ After many messages, the probabilities of the states *become stationary* (converge to some fixed values), irrespective of the initial probabilities (no matter what state the source started from initially)

$$\lim_{n \rightarrow \infty} [p_1^{(n)}, p_2^{(n)}, \dots, p_N^{(n)}] = [p_1, p_2, \dots, p_N]$$



# Finding the stationary probabilities

- ▶ How to find the value of the stationary probabilities?
- ▶ When  $n$  is very large, after  $n$  messages and after  $n+1$  messages the probabilities are the same:

$$\underbrace{[p_1, p_2, \dots, p_N]}_{\text{at time } n} \cdot [T] = \underbrace{[p_1, p_2, \dots, p_N]}_{\text{at time } n+1}$$

if  $n \rightarrow \infty$

system with  
 $N$  unknowns ( $p_k$ )  
and  $N$  equations

- ▶ This is an equation system in matrix form
- ▶ One line should be removed (linear combination), and replaced with:

$$p_1 + p_2 + \dots + p_N = 1$$

- ▶ Solve the resulting system of equations, find values of  $p_k$

# Entropy of ergodic sources with memory

- The entropy of an ergodic source with memory is

$$\underline{H(S)} = \sum_k \underbrace{p_k}_{\text{}} \underbrace{H(S_k)}_{\text{}} = - \sum_k p_k \cdot \underbrace{\sum_i p(s_i|S_k) \cdot \log(p(s_i|S_k))}_{\#(S_k)}$$

## Exercise

1. Consider a discrete source with memory, with the graphical representation given below. The states are defined as follows:

$S_1 : s_1 s_1$ ,  $S_2 : s_1 s_2$ ,  $S_3 : s_2 s_1$ ,  $S_4 : s_2 s_2$ .

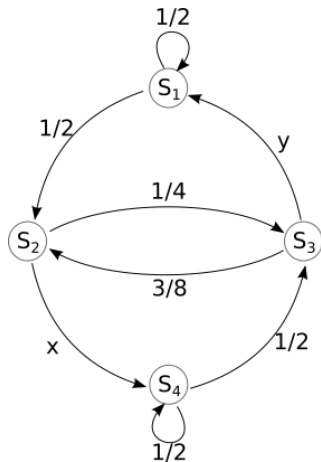


Figure 3: Graphical representation of the source

## Exercise (continued)

Questions:

- a. What are the values of  $x$  and  $y$ ?
- b. Write the transition matrix  $[T]$ ;
- c. Compute the entropy in state  $S_4$ ;
- d. Compute the global entropy of the source;
- e. What are the memory order,  $m$ , and the number of messages of the source,  $n$ ?
- f. If the source is initially in state  $S_2$ , in what states and with what probabilities will the source be after 2 messages?

# Example English text as sources with memory

(taken from *Elements of Information Theory*, Cover, Thomas)

- ▶ Memoryless source, equal probabilities:



XFOML RXKHRJFFJUJ ZLPWCFWKCYJ  
FFJEYVKCQSGXYD QPAAMKBZAACIBZLHJQD

- ▶ Memoryless source, probabilities of each letter as in English:



OCRO HLI RGWR NMIELWIS EU LL NBNESEBYA TH EEI  
ALHENHTTPA OOBTTVA NAH BRL

- ▶ Source with memory  $m = 1$ , frequency of pairs as in English:



QN IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY  
ACHIN D ILONASIVE TUCOOWE AT TEASONARE EUSO  
TIZIN ANDY TOBE SEACE CTISBE

A ~~A~~

$S : (A \dots Z)$

# Example English text as sources with memory

- ▶ Source with memory  $m = 2$ , frequency of triplets as in English:

IN NO IST LAT WHEY CRATICT FROURE BERS GROCID  
PONDENOME OF DEMONSTURES OF THE REPTAGIN IS  
REGOACTIONA OF CRE

- ▶ Source with memory  $m = 3$ , frequency of 4-plets as in English:

THEGENERATEDJOBPROVIDUALBETTERTRANDTHEDISPLAYED  
CODE, ABOVEVERY UPONDULTS WELL THE CODERST IN THESTICAL  
IT DO HOCK BOTHE MERG. (INSTATES CONS ERATION. NEVER  
ANY OF PUBLEAND TO THEORY. EVENTIAL CALLEGAND TO ELAST  
BENERATED IN WITH PIES AS IS WITH THE )

NE ?

EVE

## Example application

- ▶ Suppose we receive a text with random missing letters
- ▶ We need to fill the blanks with the appropriate letters
- ▶ How?
  - ▶ build a model: source with memory of some order
  - ▶ fill the missing letter with the most likely letter given by the model



## Chapter summary

- ▶ Information of a message:  $i(s_k) = -\log_2(p(s_k))$
- ▶ Entropy of a memoryless source:  
$$H(S) = \sum_k p_k i(s_k) = -\sum_k p_k \log_2(p_k)$$
- ▶ Properties of entropy:
  1.  $H(S) \geq 0$
  2. Is maximum when all messages have equal probability  
( $H_{\max}(S) = \log(n)$ )
  3. *Diversification* of the source always increases the entropy
- ▶ Sources with memory: definition, transitions
- ▶ Stationary probabilities of ergodic sources with memory:  
 $[p_1, p_2, \dots, p_N] \cdot [T] = [p_1, p_2, \dots, p_N], \sum_i p_i = 1.$
- ▶ Entropy of sources with memory:

$$H(S) = \sum_k p_k H(S_k) = -\sum_k p_k \sum_i p(s_i|S_k) \cdot \log(p(s_i|S_k))$$