

## Information Theory

## Chapter II: Source coding

# What does coding do?

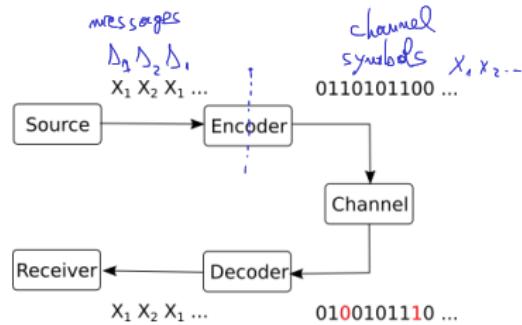


Figure 1: Communication system

## ► Why coding?

### 1. Source coding

- ▶ Convert source messages to channel symbols (for example 0,1)
- ▶ Minimize number of symbols needed
- ▶ (Adapt probabilities of symbols to maximize mutual information)

### 2. Error control

- ▶ Protection against channel errors / Adds new (redundant) symbols

## Source-channel separation theorem

Source-channel separation theorem (informal):

- ▶ It is possible to obtain the best reliable communication by performing the two tasks separately:
  - 1. Source coding: to minimize number of symbols needed
  - 2. Error control coding (channel coding): to provide protection against noise

## Source coding

- ▶ Assume we code for transmission over ideal channels with no noise
- ▶ Transmitted symbols are perfectly recovered at the receiver
- ▶ Main concerns:
  - ▶ minimize the number of symbols needed to represent the messages
  - ▶ make sure we can decode the messages
- ▶ Advantages:
  - ▶ Efficiency
  - ▶ Short communication times
  - ▶ Can decode easily

## Definitions

- Let  $S = \{s_1, s_2, \dots, s_N\}$  = an input discrete memoryless source
- Let  $X = \{x_1, x_2, \dots, x_M\}$  = the alphabet of the code = the symbols
  - Example: binary:  $\{0, 1\}$
- A code is a mapping from  $S$  to the set of all codewords:

$$C = \{c_1, c_2, \dots, c_N\}$$

Message	Codeword
$s_1$	$c_1 = x_1 x_2 x_1 \dots$
$s_2$	$c_2 = x_1 x_2 x_2 \dots$
...	....
$s_N$	$c_N = x_2 x_2 x_2 \dots$

$$s_1 \rightarrow c_1 = 01010$$

$$s_2 \rightarrow c_2 = 01111$$

$$s_3 \rightarrow c_3 = 101$$

$$s_4 \rightarrow c_4 = \dots$$

$$\vdots$$

$$s_N \rightarrow c_N = \dots$$

	Code A	Code B
$s_1$	001	1
$s_2$	100	10
$s_3$	11	100
$s_4$	01	1000

Which one  
is better?

- Codeword length  $l_i$  = the number of symbols in  $c_i$

$l_i$  = Length of codeword  $c_i$

# Encoding and decoding

$\Delta_1 \Delta_2 \Delta_1 \Delta_2 \Delta_2 \rightarrow 01010 \underbrace{01111}_{\Delta_2} \underbrace{01010}_{\Delta_1} \underbrace{01111}_{\Delta_2} \underbrace{01111}_{\Delta_2}$

- ▶ **Encoding:** given a sequence of messages, replace each message with its codeword
- ▶ **Decoding:** given a sequence of symbols, deduce the original sequence of messages
- ▶ Example: at blackboard

$01010 \underbrace{01111}_{\Delta_2} \underbrace{010100}_{\Delta_1} \underbrace{1111}_{\Delta_2} \underbrace{01111}_{\Delta_2} \rightarrow \Delta_1 \Delta_2 \Delta_1 \Delta_2 \Delta_2$

## Example: ASCII code

Message

Letter	ASCII Code	Binary	Letter	ASCII Code	Binary
a	097	01100001	A	065	01000001
b	098	01100010	B	066	01000010
c	099	01100011	C	067	01000011
d	100	01100100	D	068	01000100
e	101	01100101	E	069	01000101
f	102	01100110	F	070	01000110
g	103	01100111	G	071	01000111
h	104	01101000	H	072	01001000
i	105	01101001	I	073	01001001
j	106	01101010	J	074	01001010
k	107	01101011	K	075	01001011
l	108	01101100	L	076	01001100
m	109	01101101	M	077	01001101
n	110	01101110	N	078	01001110
o	111	01101111	O	079	01001111
p	112	01110000	P	080	01010000
q	113	01110001	Q	081	01010001
r	114	01110010	R	082	01010010
s	115	01110011	S	083	01010011
t	116	01110100	T	084	01010100
u	117	01110101	U	085	01010101
v	118	01110110	V	086	01010110
w	119	01110111	W	087	01010111
x	120	01111000	X	088	01011000
y	121	01111001	Y	089	01011001
z	122	01111010	Z	090	01011010

Figure 2: ASCII code (partial)

# Average code length

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

- ▶ How to measure representation efficiency of a code?
- ▶ **Average code length** = average of the codeword lengths:

$$\bar{l} = \bar{l} = \sum_i p(s_i) l_i$$

- ▶ The probability of a codeword = the probability of the corresponding message

- ▶ Smaller average length: code more efficient (better)

- ▶ How small can the average length be?

2) Can you recover messages from the binary seq.?

$$011 \xrightarrow{s_4} s_2 s_3$$

	Code A	Code B	Code C	Code D
$s_1$	001	1	0	0
$s_2$	100	10	0	01
$s_3$	11	100	1	1
$s_4$	01	1000	1	011

$$\bar{l}_A = \frac{0.4 \cdot 3}{0.4} + \frac{0.3 \cdot 3}{0.3} + \frac{0.2 \cdot 2}{0.2} + \frac{0.1 \cdot 2}{0.1} = 2.7 \text{ bits}$$

$$\boxed{\bar{l}_B} = \frac{0.4 \cdot 1}{0.4} + \frac{0.3 \cdot 2}{0.3} + \frac{0.2 \cdot 3}{0.2} + \frac{0.1 \cdot 4}{0.1} = 2 \text{ bits}$$

$$\bar{l}_C = 0.4 + 0.3 + 0.2 + 0.1 = 1 \text{ bits}$$

$$\widehat{01} \xrightarrow{s_3 \text{ or } s_2}$$

which one  
is better?

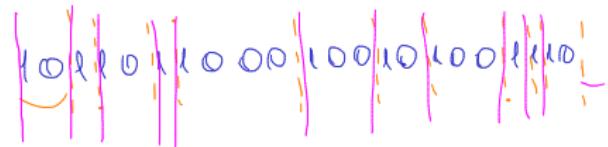
# Definitions

01011

A code can be:

- ▶ **non-singular**: all codewords are different
- ▶ **uniquely decodable**: for any received sequence of symbols, there is only one corresponding sequence of messages
  - ▶ i.e. no sequence of messages produces the same sequence of symbols
  - ▶ i.e. there is never a confusion at decoding
- ▶ **instantaneous** (also known as **prefix-free**): no codeword is prefix to another code
  - ▶ A *prefix* = a codeword which is the beginning of another codeword

Examples: at the blackboard

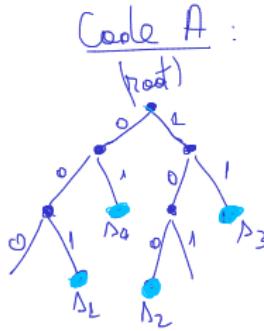


A : 0111100100  
S4 D3 S2 D2  
S5 D4 S2 D2

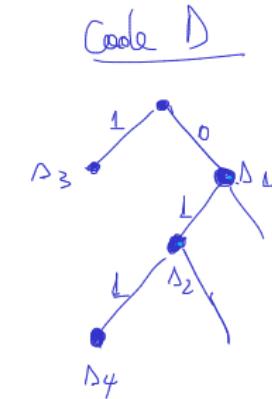
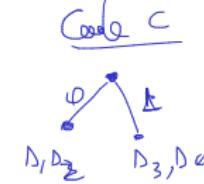
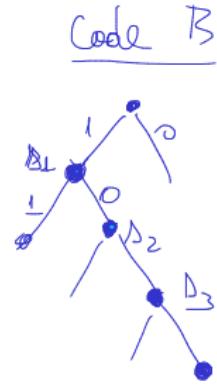
tree

# The graph of a code

nodes  
branches



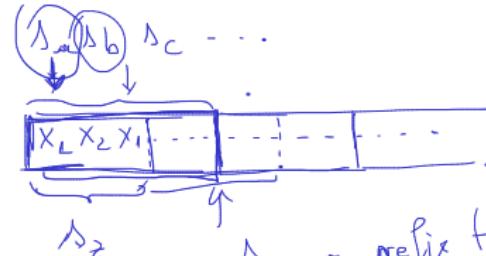
Example at blackboard



## Instantaneous codes are uniquely decodable

- ▶ Theorem:
  - ▶ An instantaneous code is uniquely decodable
- ▶ Proof:
  - ▶ There is exactly one codeword matching the beginning of the sequence
    - ▶ Suppose the true initial codeword is  $c$
    - ▶ There can't be a shorter codeword  $c'$ , since it would be prefix to  $c$
    - ▶ There can't be a longer codeword  $c''$ , since  $c$  would be prefix to it
  - ▶ Remove first codeword from sequence
  - ▶ By the same argument, there is exactly one codeword matching the new beginning, and so on ...
- ▶ Note: the converse is not necessary true; there exist uniquely decodable codes which are not instantaneous

Proof :



$A_z$  a prefix to  $A_a$   
 $A_a$  a prefix to  $A_z$

Code A	
$A_1$	001
$A_2$	100
$A_3$	11
$A_4$	01



# Uniquely decodable codes are non-singular

non-singular = different codewords  
singular = at least 2 identical codewords

$A_1$	0 0 1
$A_2$	0 0 1
:	

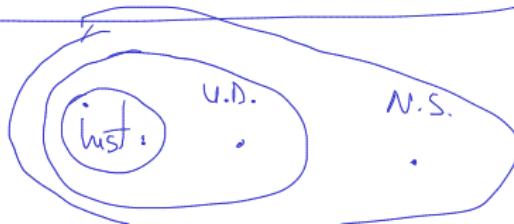
0 0 1  $\xrightarrow{A_1}$  ?  
 $\xrightarrow{A_2}$  ?

Singular  $\Rightarrow$  not unq. decodable

$$\begin{array}{ccc} A & \xrightarrow{\quad} & B \\ & \Downarrow & \\ & \text{not } B & \xrightarrow{\quad} \text{not } A \end{array}$$

- ▶ Theorem:
  - ▶ An uniquely decodable code is non-singular
- ▶ Proof:
  - ▶ If the code is singular, some codewords are not unique (different messages, same codeword)
  - ▶ Don't know which of those messages was there  $\Rightarrow$  not uniquely decodable
  - ▶ So if the code is uniquely-decodable, it must also be non-singular ( $A \rightarrow B \Leftrightarrow \overline{B} \rightarrow \overline{A}$ )
- ▶ Relation between code types:
  - ▶ Instantaneous  $\subset$  uniquely decodable  $\subset$  non-singular

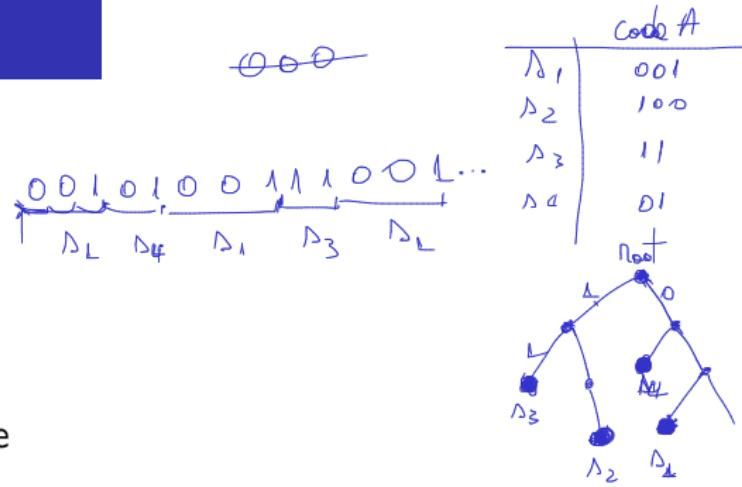
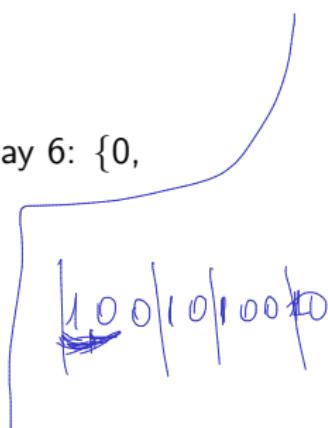
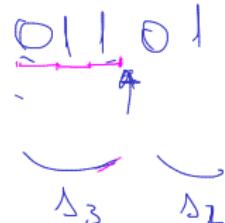
Venn diagrams:



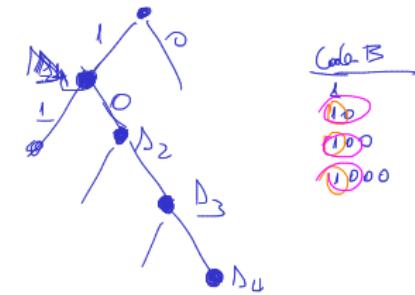
# Graph-based decoding of instantaneous codes

- ▶ How to decode an instantaneous code: graph-based decoding
  - ▶ Illustrate at whiteboard
- ▶ Advantage on instantaneous code over uniquely decodable: simple decoding
- ▶ Why the name *instantaneous*?
  - ▶ The codeword can be decoded as soon as it is fully received
  - ▶ Counter-example: Uniquely decodable, non-instantaneous, delay 6: {0, 01, 011, 1110}

Codeword:	$\Delta_1$	0
	$\Delta_2$	01
	$\Delta_3$	011
	$\Delta_4$	1110



Code B



# Existence of instantaneous codes

- When can an instantaneous code exist?

- Kraft inequality theorem:

- There exists an instantaneous code with  $D$  symbols and codeword lengths  $l_1, l_2, \dots, l_n$  if and only if the lengths satisfy the following inequality:

$$\sum_i D^{-l_i} \leq 1.$$

$D=2 \quad \{0,1\}$

- Proof: At blackboard
- Comments:

- If lengths do not satisfy this, no instantaneous code exists
- If the lengths of a code satisfy this, that code can be instantaneous or not (there exists an instantaneous code, but not necessarily that one)
- Kraft inequality means that the codewords lengths cannot be all very small

Sum of all shadows  $\leq 2^N \Rightarrow \sum_i 2^{N-l_i} \leq 2^N : 2^N \Rightarrow \sum_i 2^{-l_i} \leq 1$  Kraft q.e.d.

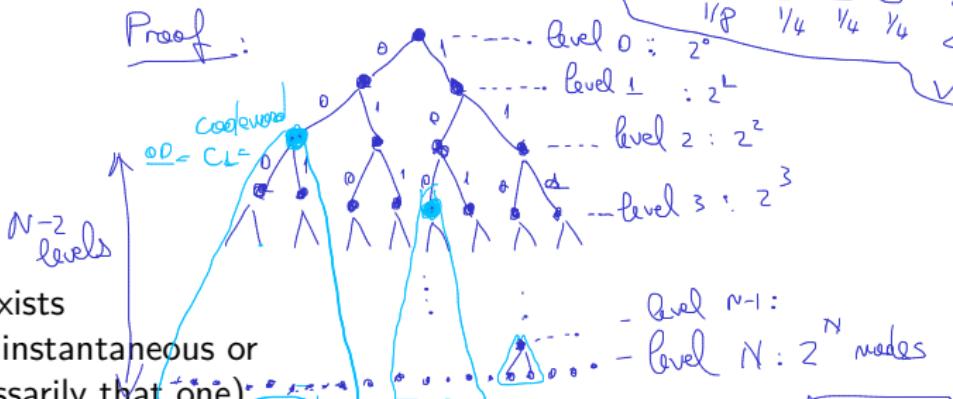
	$D_1$	$D_2$	$D_3$	$D_4$		$D_1$	$D_2$	$D_3$	$D_4$	
	0 1	1				0 1	1			
	0 0	0				0 0	1			
	0 1	1				1 1	2			
	1					1 0	2			

lengths

check formula

a)  $2^{-1} + 2^{-2} + 2^{-2} + 2^{-1} = \frac{9}{8} \leq 1$

b)  $2^{-3} + 2^{-2} + 2^{-2} + 2^{-2} = \frac{7}{8} \leq 1$



How many nodes in the last layer are shadowed by codeword  $c_i$ ?

$\sum_i 2^{-l_i} \leq 1$  Kraft q.e.d.

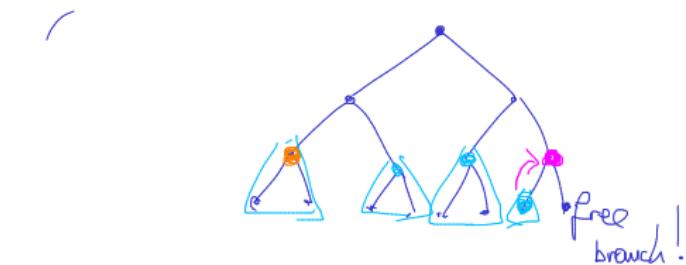
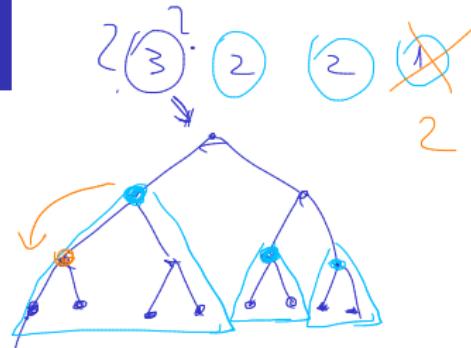
# Instantaneous codes with equality in Kraft

- ▶ From the proof  $\Rightarrow$  we have equality in the relation

$$\sum_i D^{-l_i} = 1$$

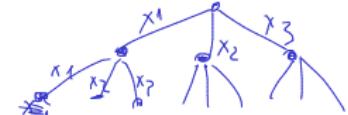
only if the lowest level is fully covered  $\Leftrightarrow$  no unused branches

- ▶ For an instantaneous code which satisfies Kraft with  $= 1$ , all the graph branches terminate with codewords (there are no unused branches)
  - ▶ This is most economical: codewords are as short as they can be



$\sim \{0, 1\}$

$\{x_1, x_2, x_3\}$



# Kraft inequality for uniquely decodable codes

- ▶ Instantaneous codes must obey Kraft inequality
- ▶ How about uniquely decodable codes?
- ▶ McMillan theorem (no proof given):
  - ▶ Any uniquely decodable code **also** satisfies the Kraft inequality:



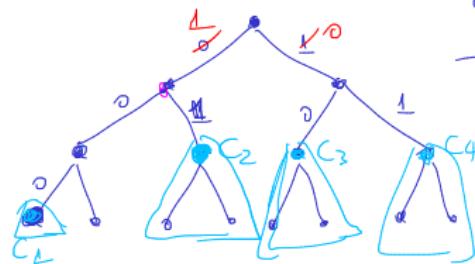
$$\sum_i D^{-l_i} \leq 1.$$

- ▶ Consequence:
  - ▶ For every uniquely decodable code, there exists an instantaneous code with the same lengths!
  - ▶ Even though the class of uniquely decodable codes is larger than that of instantaneous codes, it brings no benefit in codeword length
  - ▶ We can always use just instantaneous codes.

# Finding an instantaneous code for given lengths

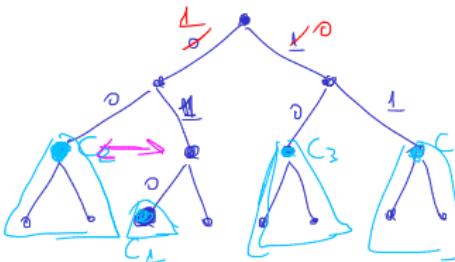
lengths: { 3, 2, 2, 2 }

$$\begin{aligned}c_1 &= \underline{\underline{0}} \underline{\underline{0}} \underline{\underline{0}} \\c_2 &= \underline{\underline{0}} \underline{\underline{1}} \\c_3 &= \underline{\underline{1}} \underline{\underline{0}} \\c_4 &= \underline{\underline{1}} \underline{\underline{1}}\end{aligned}$$



$$\begin{aligned}c_1 &= \underline{\underline{1}} \underline{\underline{0}} \\c_2 &= \underline{\underline{1}} \underline{\underline{1}} \\c_3 &= \underline{\underline{0}} \underline{\underline{0}} \\c_4 &= \underline{\underline{0}} \underline{\underline{1}}\end{aligned}$$

- ▶ How to find an instantaneous code with code lengths  $\{l_i\}$ 
  1. Check that lengths satisfy Kraft relation
  2. Draw graph
  3. Assign nodes in a certain order (e.g. descending probability)  
*(optional)*  
Always keep left!
- ▶ Easy, standard procedure
- ▶ Example: at blackboard



$$\begin{aligned}c_1 &= \underline{\underline{1}} \underline{\underline{0}} \\c_2 &= \underline{\underline{1}} \underline{\underline{1}} \\c_3 &= \\c_4 &=\end{aligned}$$

# Optimal codes

- We want to minimize the average length of a code:

$$\bar{l} = \bar{l} = \sum_i p(s_i)l_i \quad \text{Want minimum}$$

$$l_i = 1$$

$$\begin{aligned}c_1 &= - \\c_2 &= - \\c_3 &= - \\c_4 &= - \\c_5 &= -\end{aligned}$$

- But the lengths must obey the Kraft inequality (for uniquely decodable)
- So we reach the following **constrained optimization problem**:

Find  $l_i$  such that

$$\boxed{\begin{array}{l}\text{minimize } \sum_i p(s_i)l_i \\ \text{subject to } \sum_i D^{-l_i} \leq 1\end{array}}$$

## The method of Lagrange multipliers

- ▶ Method of Lagrange multipliers: standard mathematical tool
- ▶ To solve the following constrained optimization problem

$$\begin{aligned} & \text{minimize } f(x) \\ & \text{subject to } g(x) = 0 \end{aligned}$$

one must build a new function  $L(x, \lambda)$  (the **Lagrangean function**):

$$L(x, \lambda) = f(x) - \lambda g(x)$$

and the solution  $x$  is among the solutions of the system:

$$\left\{ \begin{array}{l} \frac{\partial L(x, \lambda)}{\partial x} = 0 \\ \frac{\partial L(x, \lambda)}{\partial \lambda} = 0 \end{array} \right.$$

- ▶ If there are multiple variables  $x_i$ , derivation is done for each one

# Solving for minimum average length of code

In our case : The unknowns =  $\{l_i\}$   
 $f(x) = \bar{l} = \sum_i p(s_i) \cdot l_i$

$$g(x) = \sum_i D^{-l_i} - L = 0$$

► In our case:

- The unknown  $x$  are  $l_i$
- The function is  $f(x) = \bar{l} = \sum_i p(s_i)l_i$
- The constraint is  $g(x) = \sum_i D^{-l_i} - L = 0$

► (Solve at blackboard)

► The optimal values are:

$$l_i = -\log(p(s_i))$$

► Intuition: using  $l_i = -\log(p(s_i))$  satisfies Kraft with equality, so the lengths cannot be any shorter, in general

$$(x^*)' = a \cdot \ln a \frac{\partial}{\partial l_1} \left\{ D^{-l_1} \right\} = -D^{-l_1} \cdot \ln D$$

continued  
↓ ↓

$$\left( x \cdot \ln D \right)' = -1$$

$$p(s_1) \cdot l_1 + p(s_2) \cdot l_2 + \dots + p(s_n) \cdot l_n + \lambda \left( \sum_i D^{-l_i} - 1 \right) = 0$$

$f(x)$        $g(x)$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial l_1} = p(s_1) + \lambda \cdot \left( D^{-l_1} \cdot \ln D \right) = 0 \\ \frac{\partial L}{\partial l_2} = p(s_2) + \lambda \cdot \left( D^{-l_2} \cdot \ln D \right) = 0 \\ \vdots \\ \frac{\partial L}{\partial l_N} = p(s_N) + \lambda \cdot \left( D^{-l_N} \cdot \ln D \right) = 0 \\ \frac{\partial L}{\partial \lambda} = - \left( \sum_i D^{-l_i} - 1 \right) = 0 \end{array} \right. \Rightarrow \sum_i D^{-l_i} = L$$

add all ||

$$\sum_i p(s_i) + \lambda \cdot \ln D \cdot \left( \sum_i D^{-l_i} \right) = 0$$

$$\Rightarrow 1 + \lambda \cdot \ln D = 0 \Rightarrow$$

# Optimal lengths

pt.

- ▶ The optimal codeword lengths are:

$$l_i = -\log(p(s_i))$$

$p(s_i)$  big  $\Rightarrow$  short codeword

- ▶ Higher probability  $\Rightarrow$  smaller codeword

- ▶ more efficient
- ▶ language examples: "da", "nu", "the", "le" ...

- ▶ Smaller probability  $\Rightarrow$  longer codeword

- ▶ it appears rarely  $\Rightarrow$  no problem

- ▶ Overall, we obtain the minimum average length

↓ continued

$$\begin{aligned} p(s_i) + D^{-l_1} &= 0 \Leftrightarrow p(s_i) = D^{-l_1} \quad | \log_D \\ &\Rightarrow \log_D p(s_i) = -l_1 \\ &\Rightarrow l_1 = -\log_D p(s_i) \\ &\vdots \\ &\text{some for all: } l_i = -\log_D p(s_i) \end{aligned}$$

$D = 2$

## Entropy = minimal codeword average length

- If the optimal values are:

$$l_i = l_i = -\log(p(s_i))$$

- Then the minimal average length is:

$$\min \bar{l} = \sum_i p(s_i)l_i = -\sum_i p(s_i) \log(p(s_i)) = H(S)$$

- The entropy of a source = the minimum average length necessary to encode the messages

- e.g. the minimum number of bits required to represent the data in binary form

## Meaning of entropy

- ▶ This tells us something about entropy
  - ▶ This is what entropy means in practice
  - ▶ Small entropy => can be written (encoded) with few bits
  - ▶ Large entropy => requires more bits for encoding
- ▶ This tells us something about the average length of codes
  - ▶ The average length of an uniquely decodable code must be at least as large as the source entropy

$$H(S) \leq \bar{l}$$

- ▶ One can never represent messages, on average, with a code having average length less than the entropy

## Analogy of entropy and codes



### ► Analogy: 1 liter of water

- ▶ 1 liter of water = the quantity of water that can fit in any bottle of size  $\geq 1$  liter, but not in any bottle  $< 1$  liter

$$Bottle \geq water$$

- ▶ Information of the source = the water
- ▶ The code used for representing the messages = the bottle that carries the water

$$\bar{I} \geq H(S)$$

/    //

## Efficiency and redundancy of a code

- Efficiency of a code ( $M$  = size of code alphabet):

$$\eta = \frac{H(S)}{\bar{l} \log_2 M}$$

in binary

$$\eta = \frac{H(S)}{\bar{l} \cdot \log_2 M}$$

\* Efficiency of a code :  
 $\frac{H(S)}{H_{\max}}$

\* Efficiency of a code  
 $\frac{H(S)}{\bar{l} \cdot \log_2 M}$

- ▶ usually  $M = 2$  so  $\eta = \frac{H(S)}{\bar{l}}$
- ▶ but if  $M > 2$  a factor of  $\log M$  is needed because  $H(S)$  in bits (binary)  
but  $\bar{l}$  not in bits ( $M$  symbols)

- Redundancy of a code:

$$\rho = 1 - \eta$$

- ▶ These measures indicate how close is the average length to the optimal value
- ▶ When  $\eta = 1$ : optimal code

# Optimal codes

- ▶ Problem:  $l_i = -\log(p(s_i))$  might not be an integer number       $2.32 \rightarrow 3$

- ▶ but the codeword lengths must be natural numbers

- ▶ An optimal code = a code that attains the minimum average length

$$\bar{l} = H(S)$$

- ▶ An optimal code can always be found for a source where all  $p(s_i)$  are powers of 2

- ▶ e.g.  $1/2, 1/4, 1/2^n$ , known as *dyadic distribution*

- ▶ the lengths  $l_i = -\log(p(s_i))$  are all natural numbers  $\Rightarrow$  can be attained

- ▶ the code with lengths  $l_i$  can be found with the graph-based procedure

$$H(S) = -\frac{1}{2} \log_2 \frac{1}{2} - 2 \cdot \frac{1}{4} \log_2 \frac{1}{4} = \frac{1}{2} + \frac{1}{2} = 1.5 \text{ b}$$

$$\bar{l} = \frac{1}{2} \cdot 1 + 2 \cdot \frac{1}{4} \cdot 2 = \frac{1}{2} + 1 = 1.5 \text{ b} \quad \left. \begin{array}{l} \\ \end{array} \right\} \Rightarrow y=1$$

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

$$\begin{array}{l} l_1 = -\log_2 \frac{1}{2} = 1 \\ l_2 = -\log_2 \frac{1}{4} = 2 \\ l_3 = -\log_2 \frac{1}{4} = 2 \end{array} \quad \left. \begin{array}{l} 0 \\ 10 \\ 11 \end{array} \right\}$$

## Non-optimal codes

- ▶ What if  $-\log(p(s_i))$  is not a natural number? i.e.  $p(s_i)$  is not a power of 2
- ▶ Shannon's solution: round to next largest natural number

$$l_i = \lceil -\log(p(s_i)) \rceil$$

i.e.  $-\log(p(s_i)) = 2.15 \Rightarrow l_i = 3$

$$\text{floor} = \lfloor 1.8 \rfloor = 1$$

$$\text{ceil} = \lceil 2.15 \rceil = 3$$

$$\lceil 3 \rceil = 3$$

$$\begin{aligned}\text{floor}(1.8) &= 1 \\ \text{ceil}(1.8) &= 2\end{aligned}$$

$$\mathbb{S} = \begin{pmatrix} \Delta_1 & \Delta_2 & \Delta_3 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix}$$

$$l_{1,2,3} = -\log_2 \frac{1}{3} = 1.58 \text{ b}$$

$$\sum l_i = 1 \Leftrightarrow 2^{-1.58} + 2^{-1.58} + 2^{-1.58} = 1$$

# Shannon coding

- ▶ Shannon coding:

1. Arrange probabilities in descending order
2. Use codeword lengths  $l_i = \lceil -\log(p(s_i)) \rceil$
3. Find any instantaneous code for these lengths\*
4. For every message  $s_i$ :
  - compute the sum of all the probabilities up to this message
  - multiply this value with  $2^{l_i}$   
not including the current row
  - floor the result and convert to binary

- ▶ The code obtained = a “Shannon code”
- ▶ Simple scheme, better algorithms are available
  - ▶ Example: compute lengths for  $S : (0.9, 0.1)$
- ▶ But still enough to prove fundamental results

$$c_1: 0 \cdot 2^0 = 0$$

$$c_2: 0.5 \cdot 2^1 = 2 \rightarrow 10$$

$$c_3: 0.8 \cdot 2^4 = 12.8 : 1100$$

$$c_4: 0.9 \cdot 2^4 = 14.4 : 1110$$

Shannon Coding

$$\# S = \begin{pmatrix} \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 \\ 0.5 & 0.1 & 0.3 & 0.1 \end{pmatrix}$$

	$p(s_i)$	$\lceil l_i \rceil$	Codewords
$\Delta_1$	0.5	1	0
$\Delta_3$	0.3	$\lceil 1.73 \rceil = 2$	10
$\Delta_2$	0.1	$\lceil 3.32 \rceil = 4$	1100
$\Delta_4$	0.1	$\lceil 3.32 \rceil = 4$	1101

	$p(s_i)$	$\lceil l_i \rceil$	Cumulative sum	
$\Delta_1$	0.5	1	0	0
$\Delta_3$	0.3	$\lceil 1.73 \rceil = 2$	0.5	10
$\Delta_2$	0.1	$\lceil 3.32 \rceil = 4$	0.8	1100
$\Delta_4$	0.1	$\lceil 3.32 \rceil = 4$	0.9	1110

## Average length of Shannon code

$$H(S) = 7.3$$

$$\bar{l} \in \{7.3, 8.3\}$$

Theorem:

- The average length of a Shannon code satisfies :

$$H(S) \leq \bar{l} < H(S) + 1$$



# Average length of Shannon code

Proof:

1. The first inequality is because  $H(S)$  is minimum length
2. The second inequality:

a. Use Shannon code:

$$l_i = \lceil -\log(p(s_i)) \rceil = -\log(p(s_i)) + \epsilon_i$$

where  $0 \leq \epsilon_i < 1$

b. Compute average length:

$$\bar{l} = \bar{l} = \sum_i p(s_i) l_i = H(S) + \underbrace{\sum_i p(s_i) \epsilon_i}_{<1}$$

c. Since  $\epsilon_i < 1 \Rightarrow \sum_i p(s_i) \epsilon_i < \sum_i p(s_i) = 1$

average

$$\lceil 2.32 \rceil = 3 = 2.32 + \underbrace{0.68}_{\epsilon < 1}$$

$$\sum_i p(s_i) \cdot (-\log(p(s_i)) + \epsilon_i) = \underbrace{-\sum_i p(s_i) \log(p(s_i))}_{H(S)} + \underbrace{\sum_i p(s_i) \cdot \epsilon_i}_{<1} \underbrace{\sum_i p(s_i)}_{<1}$$

$$\sum_i p(s_i) \cdot \epsilon_i < \sum_i p(s_i) = 1$$

$$\bar{l} = H(S) + (<1) \quad \square$$

$$\Rightarrow \boxed{\bar{l} < H(S) + 1}$$

## Average length of Shannon code

- ▶ Average length of Shannon code is **at most 1 bit longer** than the minimum possible value
  - ▶ That's quite efficient
  - ▶ There exist even better codes, in general
- ▶ Q: Can we get even closer to the minimum length?
- ▶ A: Yes, as close as we want!
  - ▶ In theory, at least . . . :)
  - ▶ See next slide.

## Shannon's first theorem

Shannon's first theorem (coding theorem for noiseless channels):

- ▶ It is possible to encode an infinitely long sequence of messages from a source  $S$  with an average length as close as desired to  $H(S)$ , but never below  $H(S)$

$$H(S) = 4.5$$

$$\overline{l} = 4.500000000 L$$

Key points:

- ▶ we can always obtain  $\overline{l} \rightarrow H(S)$
- ▶ for an infinitely long sequence

# Shannon's first theorem

Proof:

- ▶ Average length can never go below  $H(S)$  because this is minimum
- ▶ How can it get very close to  $H(S)$  (from above)?

- . 1. Use  **$n$ -th order extension  $S^n$**  of  $S$
- . 2. Use Shannon coding for  $S^n$ , so it satisfies *for  $n$  messages of orig.  $S$*

$$H(S^n) \leq \overline{I_{S^n}} < H(S^n) + 1 \quad | : n$$

- . 3. But  $H(S^n) = nH(S)$ , and average length **per message of  $S$**  is

$$\overline{I_S} = \frac{\overline{I_{S^n}}}{n}$$

because messages of  $S^n$  are just  $n$  messages of  $S$  glued together

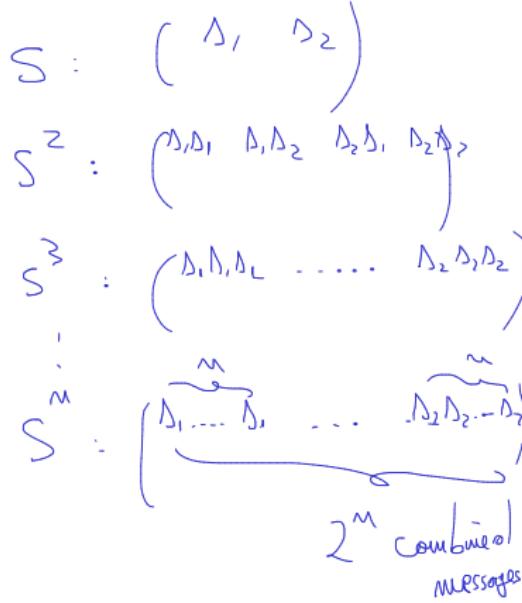
- . 4. So, dividing by  $n$ :

$$H(S) \leq \overline{I_S} < H(S) + \frac{1}{n}$$

- . 5. If extension order  $n \rightarrow \infty$ , then

$$\overline{I_S} \rightarrow H(S)$$

□



# Shannon's first theorem

- ▶ Analogy: how to buy things online without paying for delivery :)
- ▶ FanCourier taxes 15 lei per delivery
  - ▶ not efficient to buy something worth a few lei
- ▶ How to improve efficiency? Buy  $n$  things bundled together!
- ▶ The delivery cost **per unit** is now  $\frac{15}{n}$
- ▶ As  $n \rightarrow \infty$ , the delivery cost per unit  $\rightarrow 0$ 
  - ▶ What's 15 lei when you pay  $\infty$  lei...

$$S : \begin{pmatrix} \Delta_1 & \Delta_2 \\ 0.9 & 0.1 \end{pmatrix}$$

$$H(S) = -0.9 \log 0.9 - 0.1 \log 0.1 = 0.46$$

Shannon coding :  $\overline{l}$

			Cum. sum	codewords
S <sub>1</sub>	0.9	$\lceil 0.15 \rceil = 1$	0	0
S <sub>2</sub>	0.1	$\lceil 3.32 \rceil = 4$	0.9	1110

$$\overline{l} = 0.9 \cdot 1 + 0.1 \cdot 4 = 1.3$$

$$y = \frac{H(S)}{\overline{l}} = \frac{0.46}{1.3}$$

$$S^2 : \begin{pmatrix} \Delta_1 \Delta_1 & \Delta_1 \Delta_2 & \Delta_2 \Delta_1 & \Delta_2 \Delta_2 \\ 0.81 & 0.09 & 0.09 & 0.01 \end{pmatrix}$$

Shannon :

			codewords
0.81	1		
0.09	4		
0.09	4		
0.01	7		
			...

8

## Shannon's first theorem

$S$ : 256 messages

$$S^5: 256^5 \text{ messages} \approx (2^8)^5 = 2^{40} \approx 1000000000000 \text{ messages}$$

Comments:

- ▶ Shannon's first theorem shows that we can approach  $H(S)$  to any desired accuracy using extensions of large order of the source
  - ▶ This is not practical: the size of  $S^n$  gets too large for large  $n$
  - ▶ Other (better) algorithms than Shannon coding are used in practice to approach  $H(S)$

$$\begin{aligned}\overline{L} &= 0.81 \cdot 1 + 0.09 \cdot 4 + \\ &0.09 \cdot 4 + 0.01 \cdot 7 = \\ &= 1.6 \text{ bits for 2 messages of } S! \\ \Rightarrow 0.8 &\text{ for each message of } S !!\end{aligned}$$

# Coding with the wrong code

$$S_p : \begin{pmatrix} s_1 & \dots & s_m \\ p(s_1) & \dots & p(s_m) \end{pmatrix}$$

$$S_q : \begin{pmatrix} s_1 & \dots & s_m \\ q(s_1) & \dots & q(s_m) \end{pmatrix}$$

↓  
Code

- ▶ Consider a source with probabilities  $p(s_i)$
- ▶ We use a code designed for a different source:  $l_i = -\log(q(s_i))$
- ▶ The message probabilities are  $p(s_i)$  but the code is designed for  $q(s_i)$
- ▶ Examples:
  - ▶ design a code based on a sample data file (like in lab)
  - ▶ but we use it to encode various other files => probabilities might differ slightly
  - ▶ e.g. design a code based a Romanian text, but encode a text in English
- ▶ What happens?

Penalty

## Coding with the wrong code

- ▶ We lose some efficiency:
  - ▶ Codeword lengths  $\bar{l}_i$  are not optimal for our source  $\Rightarrow$  increased  $\bar{l}$
- ▶ If code were optimal, best average length = entropy  $H(S)$ :

*ignoring roundings:* 
$$\overline{l_{optimal}} = - \sum p(s_i) \underbrace{\log p(s_i)}_{\text{codeword length}}$$

- ▶ But the actual average length we obtain is:

$$\overline{l_{actual}} = \sum p(s_i) l_i = - \sum p(s_i) \underbrace{\log q(s_i)}_{+}$$

# The Kullback–Leibler distance

- ▶ Difference between average lengths is:

$$\overline{I_{actual}} - \overline{I_{optimal}} = \sum_i p(s_i) \log\left(\frac{p(s_i)}{q(s_i)}\right) = D_{KL}(p||q)$$

- ▶ The difference = **the Kullback–Leibler distance** between the two distributions
  - ▶ is always  $\geq 0 \Rightarrow$  improper code means increased  $\bar{I}$  (bad)
  - ▶ distributions more different  $\Rightarrow$  larger average length (worse)
- ▶ The KL distance between the distributions = the number of extra bits used because of a code optimized for a different distribution  $q(s_i)$  than the true distribution of our data  $p(s_i)$

# The Kullback–Leibler distance

Reminder: where is the Kullback–Leibler distance used

- ▶ Here: Using a code optimized for a different distribution:
  - ▶ Average length is increased with  $D_{KL}(p||q)$
- ▶ In chapter IV (Channels): Definition of mutual information:
  - ▶ Distance between  $p(x_i \cap y_j)$  and the distribution of two independent variables  $p(x_i) \cdot p(y_j)$

$$I(X, Y) = \sum_{i,j} p(x_i \cap y_j) \log\left(\frac{p(x_i \cap y_j)}{p(x_i)p(y_j)}\right)$$

Haven't  
done  
it  
yet

## Shannon-Fano coding (binary)

$$S = \begin{pmatrix} 1 & 2 & 3 & 4 \\ \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 \\ 0.5 & 0.1 & 0.3 & 0.1 \end{pmatrix}$$

Shannon-Fano (binary) coding procedure:

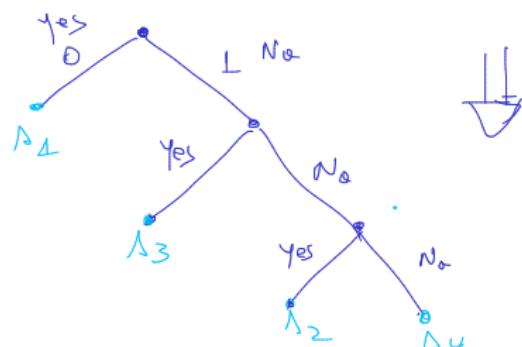
1. Sort the message probabilities in descending order
2. Split into two subgroups as nearly equal as possible
3. Assign first bit 0 to first group, first bit 1 to second group
4. Repeat on each subgroup
5. When reaching one single message => that is the codeword

Example: blackboard

Comments:

- Shannon-Fano coding does not always produce the shortest code lengths
- Connection: yes-no answers (example from first chapter)

$\Delta_1$	0.5	0
$\Delta_3$	0.3	1 0
$\Delta_2$	0.1	1 1 0
$\Delta_4$	0.1	1 1 1



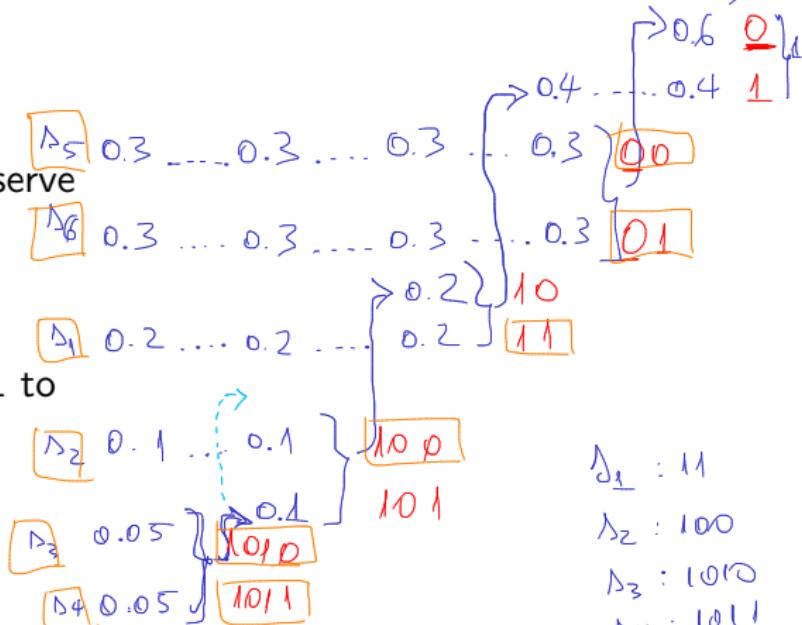
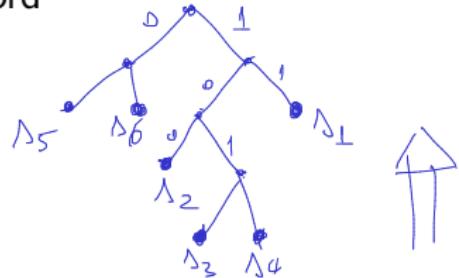
# Huffman coding (binary)

$$S: \begin{pmatrix} \Delta_1 & \Delta_2 & \Delta_3 & \Delta_4 & \Delta_5 & \Delta_6 \\ 0.2 & 0.1 & 0.05 & 0.05 & 0.3 & 0.3 \end{pmatrix}$$

## Huffman coding procedure (binary):

1. Sort the message probabilities in descending order
2. Join the last two probabilities, insert result into existing list, preserve descending order
3. Repeat until only two messages are remaining
4. Assign first bit 0 and 1 to the final two messages
5. Go back step by step: every time we had a sum, append 0 and 1 to the end of existing codeword

Example: blackboard

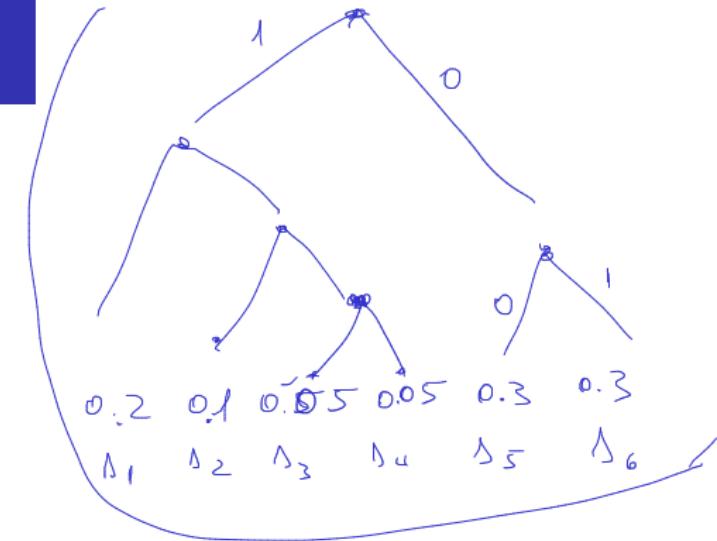


Δ₁ : 11  
Δ₂ : 100  
Δ₃ : 1010  
Δ₄ : 1011  
Δ₅ : 00  
Δ₆ : 01

# Properties of Huffman coding

Properties of Huffman coding:

- ▶ Produces a code with the smallest average length (better than Shannon-Fano)
- ▶ Assigning 0 and 1 can be done in any order => different codes, same lengths
- ▶ When inserting a sum into an existing list, may be equal to another value => options
  - ▶ we can insert above, below or in-between equal values
  - ▶ leads to codes with different *individual* lengths, but same average length
- ▶ Some better algorithms exist which do not assign a codeword to every single message (they code a ~~whole~~ sequence at once, not every message)



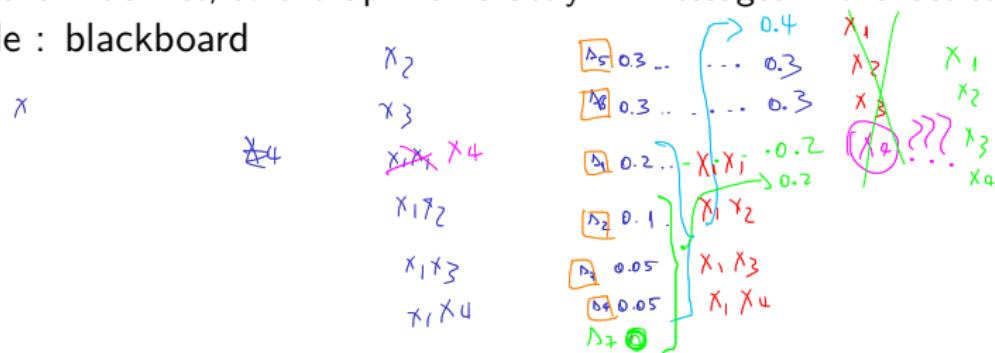
→ (but still, better ~~above~~)

# Huffman coding (M symbols)

$x_2 \ x_1 \ x_3 \ x_3 \ x_4$

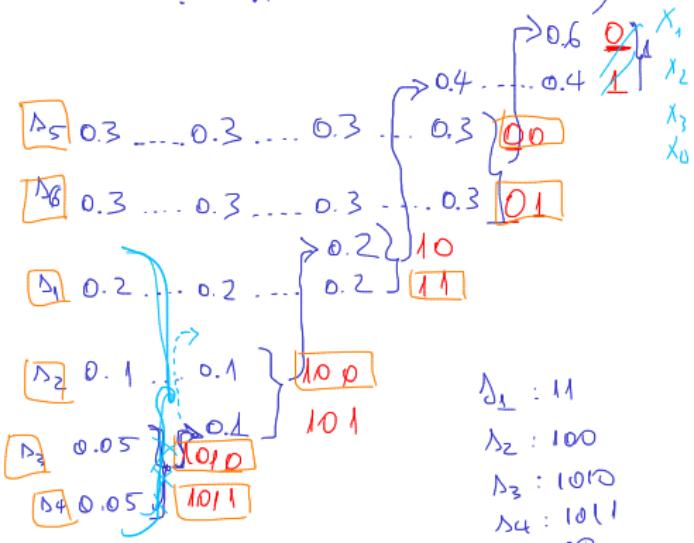
General Huffman coding procedure for codes with  $M$  symbols:

- ▶ Have  $M$  symbols  $\{x_1, x_2, \dots, x_M\} = \{x_1, x_2, x_3, x_4\}$
- ▶ Add together the last  $M$  symbols
- ▶ When assigning symbols, assign all  $M$  symbols
- ▶ **Important:** at the final step must have  $M$  remaining values
  - ▶ May be necessary to add virtual messages with probability 0 at the end of the initial list, to end up with exactly  $M$  messages in the last step
- ▶ Example : blackboard



Binary coding:  $\{0, 1\}$

$$S: \begin{pmatrix} s_1 & s_2 & s_3 & s_4 & s_5 & s_6 \\ 0.2 & 0.1 & 0.05 & 0.05 & 0.3 & 0.3 \end{pmatrix}$$



## Example: compare Huffman and Shannon-Fano

Example: compare binary Huffman and Shannon-Fano for:

$$p(s_i) = \{0.35, 0.17, 0.17, 0.16, 0.15\}$$

## Probability of symbols

$$P_0 = ? \quad P_1 = ?$$

in the resulting encoded sequence

- For every symbol  $x_i$  we can compute the average number of symbols  $x_i$  in a code

0 or 1

$$\bar{l}_{x_i} = \sum_i p(s_i) l_{x_i}(s_i)$$

$\Delta_1 \Delta_2 \Delta_1 \Delta_1 \Delta_4 \Delta_3 \dots$   
 $0 \ 1 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \dots$

Mess	$p(\Delta_i)$	Codewords
$\Delta_1$	0.5	0
$\Delta_3$	0.3	1 0
$\Delta_2$	0.1	1 1 0
$\Delta_4$	0.1	1 1 1

- $l_{x_i}(s_i)$  = number of symbols  $x_i$  in the codeword of  $s_i$
- e.g.: average number of 0's and 1's in a code

- Divide by average length  $\Rightarrow$  probability (frequency) of symbol  $x_i$

$$p(x_i) = \frac{\bar{l}_{x_i}}{\bar{l}}$$

$$P(0) = \frac{\bar{l}_0}{\bar{l}} = \frac{0.9}{1.7} = 52.9\%$$

$$P(1) = \frac{\bar{l}_1}{\bar{l}} = \frac{0.8}{1.7} = 47.1\%$$

$$P_0 = ?$$

$$P_1 = ?$$

$$\bar{l}_0 = 0.5 \cdot 1 + 0.3 \cdot 1 + 0.1 \cdot 1 =$$

$$= 0.9$$

$$\bar{l}_1 =$$

$$= 0.5 \cdot 0 + 0.3 \cdot 1 + 0.1 \cdot 2 +$$

$$0.1 \cdot 3 = 0.8$$

- These are the probabilities of the input symbols for the transmission channel

- they play an important role in Chapter IV (transmission channels)

$$P(0)$$

## Source coding as data compression

- ▶ Consider that the messages are already written in a binary code
  - ▶ Example: characters in ASCII code
- ▶ Source coding = remapping the original codewords to other codewords
  - ▶ The new codewords are shorter, on average
- ▶ This means data **compression**
  - ▶ Just like the example in lab session
- ▶ What does data compression remove?
  - ▶ Removes redundancy; unused bits, patterns, regularities etc.
  - ▶ If you can guess somehow the next bit in a sequence, it means the bit is not really necessary, so compression will remove it
  - ▶ The compressed sequence looks like random data: impossible to guess, no discernable patterns

110 X

## Discussion: data compression with coding

- ▶ Consider data compression with Shannon or Huffman coding, like we did in lab
  - ▶ What property do we *exploit* in order to obtain compression?
  - ▶ How does *compressible data* look like?
  - ▶ How does *incompressible data* look like?
  - ▶ What are the limitation of our data compression method?
  - ▶ How could it be improved?

## Other codes: arithmetic coding

- ▶ Other types of coding do exist (info only)
  - ▶ Arithmetic coding
  - ▶ Adaptive schemes
  - ▶ etc.

## Chapter summary

- ▶ Average length:  $\bar{l} = \sum_i p(s_i)l_i$
- ▶ Code types: instantaneous  $\subset$  uniquely decodable  $\subset$  non-singular
- ▶ All instantaneous or uniquely decodable code must obey Kraft:

$$\sum_i D^{-l_i} \leq 1$$

- ▶ Optimal codes:  $l_i = -\log(p(s_i))$ ,  $\overline{l_{min}} = H(S)$
- ▶ Shannon's first theorem: use  $n$ -th order extension of  $S$ ,  $S^n$ :

$$H(S) \leq \overline{l_S} < H(S) + \frac{1}{n}$$

- ▶ average length always larger, but as close as desired to  $H(S)$
- ▶ Coding techniques:
  - ▶ Shannon: ceil the optimal codeword lengths (round to upper)
  - ▶ Shannon-Fano: split in two groups approx. equal
  - ▶ Huffman: group last two. Is best of all.