

# Basics of coding theory

## 1.) NOISELESS CODING THEORY

→ Shannon entropies

→ Huffman coding

## 2.) Noisy coding theory

→ Error correcting codes

→ Block codes

## Noisless coding theory

Random variable  $X = \{x_0, \dots, x_{n-1}\}$

$p_0, \dots, p_{n-1}$

$$\sum_{i=0}^{n-1} p_i = 1$$

$\Sigma \rightarrow$  alphabet

$$\Sigma = \{0, 1\}$$

$x_0 \rightarrow$   
 $x_1 \rightarrow$   
 $x_2 \rightarrow$   
 $x_3 \rightarrow$

00
01
10
11

→ is this the most **efficient** way?

code  $C$

$\log_2 n$  bits to find codewords for  $n$  different messages

$$AVG(C) = \sum_{i=0}^{n-1} p_i \cdot |C_i|$$

↳ length of codeword for message  $x_i$ .

1.) given a probability distribution (random variable)  
what is the best achievable  $AVG(C)$ ?

2.) How to construct such a code?

---

1.) for a random variable  $X$  w.p.  $(p_0, \dots, p_{m-1})$

$$S(X) = - \sum_{i=0}^{n-1} p_i \log_2 p_i$$

Shannon entropy

I. Average length of any code  $C$  for r.v.  $X \geq S(X)$

II. Encoding multiple messages helps

III. As the number of messages encoded together approaches infinity  
 $S(X)$  is achievable.

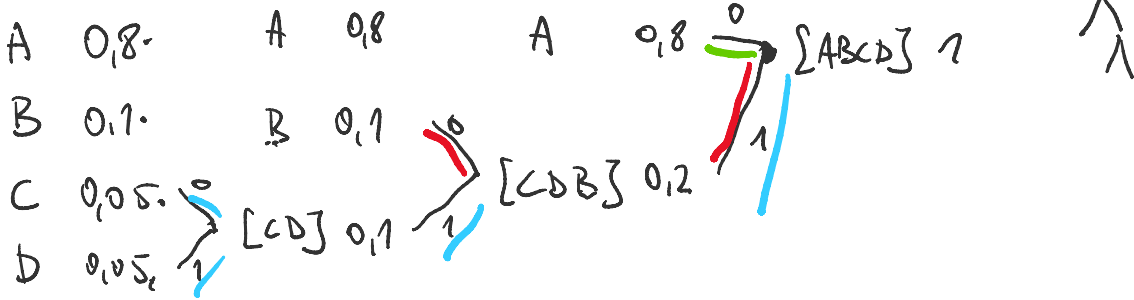
2.) Huffman coding

Alg.

INPUT: Probability distribution

OUTPUT: CODE

EX 1.2



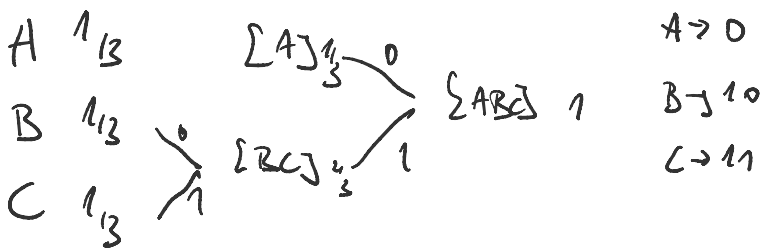
- A → 0.
- B → 10.
- C → 110.
- D → 111.

$$AVG(C) = 0.8 \cdot 1 + 0.1 \cdot 2 + 0.05 \cdot 3 + 0.05 \cdot 3$$

$$Vl = 1.3$$

$$S(X) = - (0.8 \cdot \log_2 0.8 + 0.1 \cdot \log_2 0.1 + 0.05 \log_2 0.05 + 0.05 \log_2 0.05)$$

$$= 1.02 \dots$$

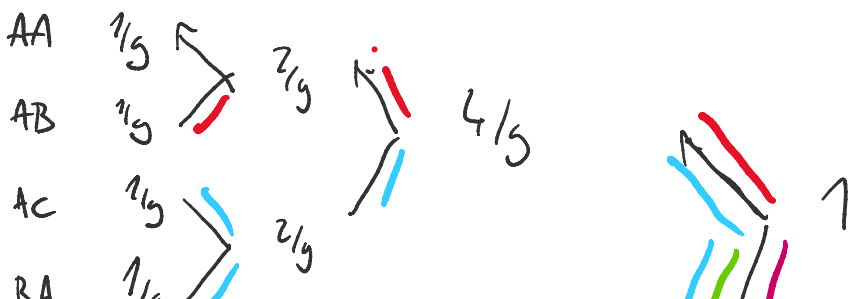


$$AVG(C) = \frac{1}{3} \cdot 1 + 2 \cdot (\frac{1}{3} \cdot 2)$$

$$Vl = \sqrt[3]{3} = 1.442 \dots$$

$$S(X) = -3 \cdot \frac{1}{3} \cdot \log_2 \frac{1}{3}$$

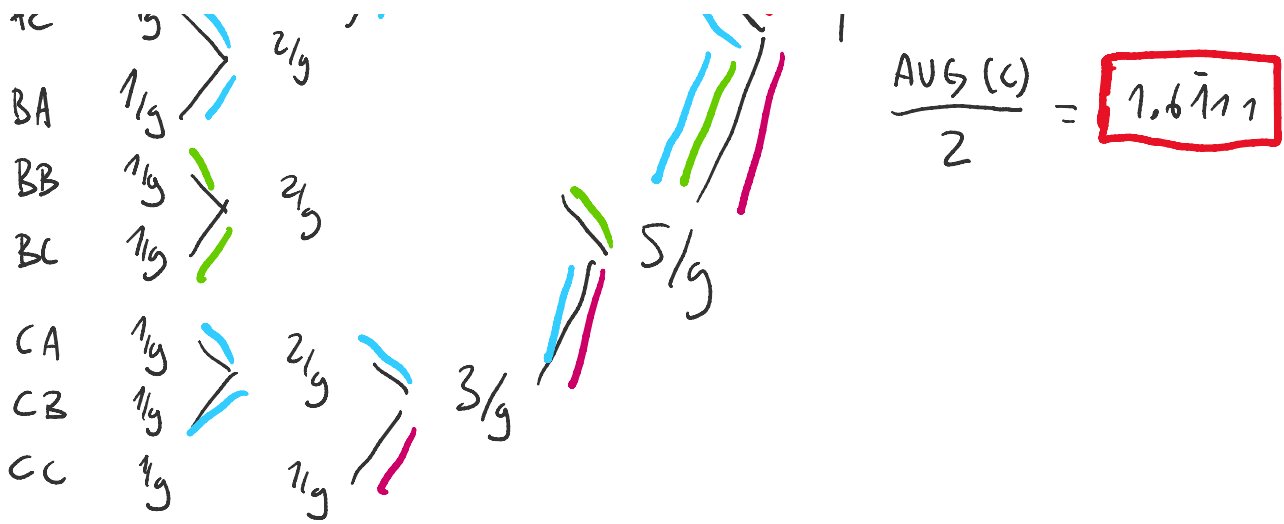
$$= -\log_2 \frac{1}{3} = 1.58$$



$$AVG(C) = 7 \cdot (\frac{1}{5} \cdot 3) + 2 \cdot (\frac{1}{5} \cdot 4)$$

$$= 3.222 \dots$$

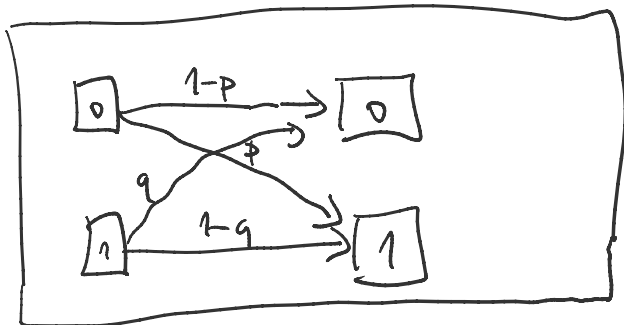
$$AVG(C) = 1.611 \dots$$



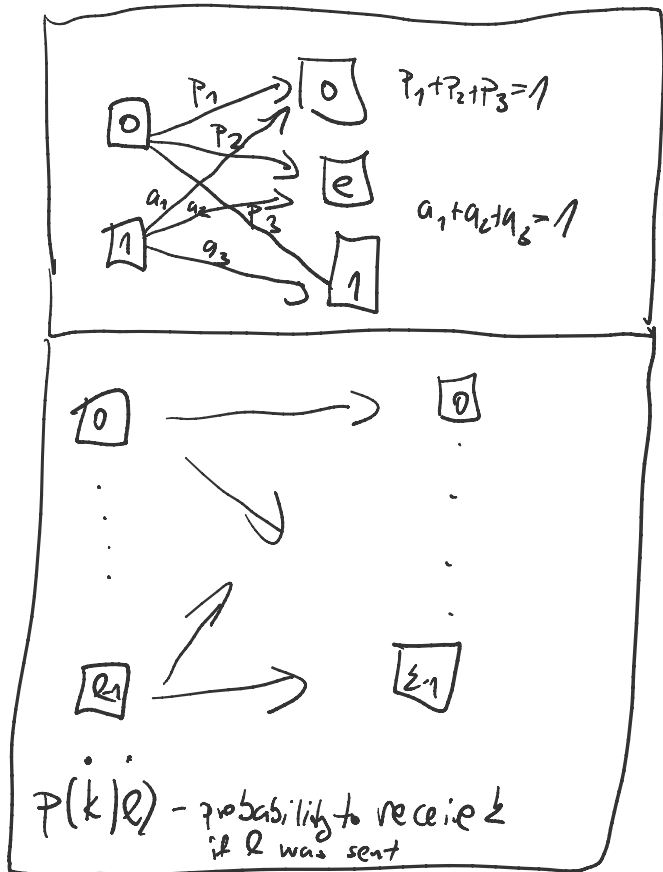
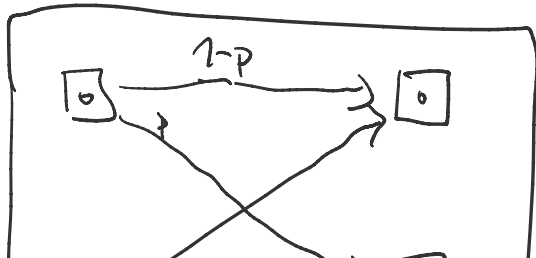
If you encode  $k$  symbols together

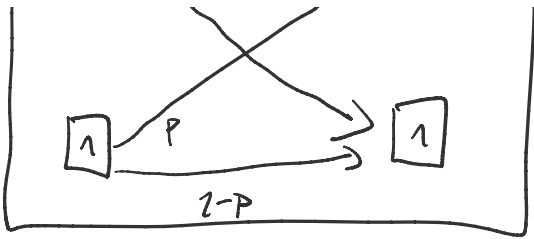
$$\text{then } \lim_{k \rightarrow \infty} \frac{AVG(\text{Huff}(k))}{k} = S(X)$$

## Noisy coding theory (Error correcting codes)



## Binary symmetric channel





$$p < 1/2$$

## Principle of maximum likelihood

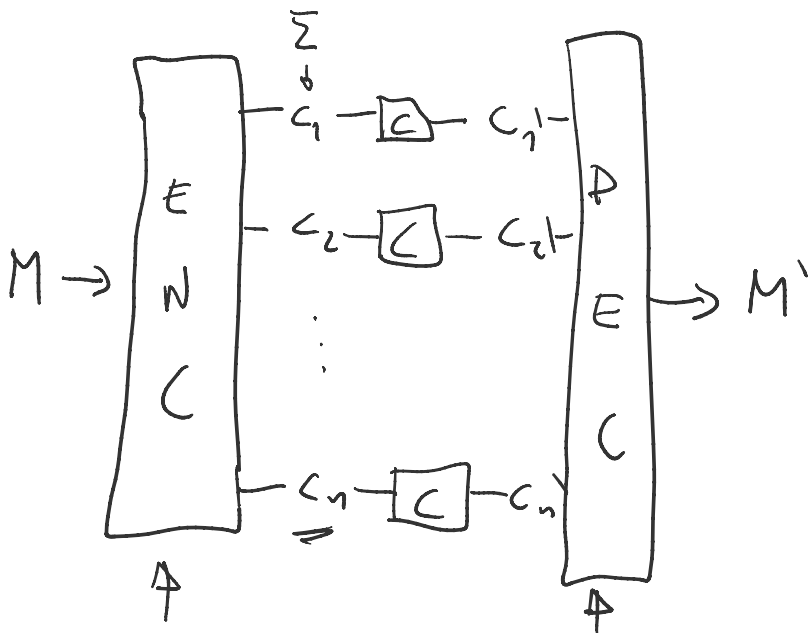
You receive 0 from a binary symmetric channel

How do you decode it?

$$0 \xrightarrow{1-p} 0$$

$$1 \xrightarrow{p} 0$$

$p < 1/2 \Rightarrow p < 1-p \Rightarrow$  it is more likely 0 was sent



	ENCODER	DECODER	
$M \in \{0, 1\}$	$0 \rightarrow 000$	$\#1 \geq 2$	decode as 1
$\Sigma \in \{0, 1\}$	$1 \rightarrow 111$ in, out	$\#1 < 2$	decode as 0

$$\Sigma \in \{0,1\} \quad 1 \rightarrow \underset{\substack{\uparrow \text{in} \\ \downarrow \text{out}}}{11} \quad \#1 < 2 \quad \text{decode as } 0$$

$$\Sigma \in \{0,1\} \quad \cdot P(000 | 001) = (1-p)(1-p) \cdot p \quad p < 1/2$$

$$P(111 | 001) = \overset{V1}{p} \cdot p \cdot (1-p)$$

Repetition code can achieve arbitrarily low probability of decoding

$n \in \mathbb{Z}, n \geq 1$	ENC	DEC
$\Sigma \in \{0,1\}$	$0 \rightarrow \overbrace{0 \dots 0}^{2k+1}$	$\#1 \geq k+1 \Rightarrow 1$
$\Sigma \in \{0,1\}$	$1 \rightarrow 1 \dots 1$	$\#1 < k+1 \Rightarrow 0$

000000...0  $\rightarrow$  less than  $k+1$  errors, the decoding is correct  
 $\rightarrow$  more than  $k$  errors, the decoding is incorrect

$$Pr(\text{correct decoding}) = \sum_{i=0}^k \binom{2k+1}{i} p^i (1-p)^{2k+1-i}$$

$$\lim_{k \rightarrow \infty} Pr(\text{correct decoding}) = 1$$

$$\frac{\# \text{Messages}}{\text{length code words}} = \frac{2}{2k+1} \underset{k \rightarrow \infty}{=} 0$$

code rate

2 homework exercise

#4  
#6

#4

#2

## Hamming distance

$$C_i \in \{0,1\}^n \quad C_i \in C \quad C \subseteq \{0,1\}^n$$

C - code

$C_i$  - codeword

$\text{Ham}(C_i, C_j)$  the number of positions in which  $C_i$  and  $C_j$  differ

**Ex 1.6**  $C = \{10001, 00110, 11010, 01101\}$

$$\text{Ham}(10001, 00110) = 4 \quad \text{Ham}(00110, 11010) = 4$$

$$\text{Ham}(10001, 11010) = 3 \quad \text{Ham}(00110, 01101) = 3$$

$$\text{Ham}(10001, 01101) = 3 \quad \text{Ham}(11010, 01101) = 4$$

$n = 5$	lengths of the codewords
$M = 4$	number of codewords
$d = 3$	minimum distance

Error detection  $\rightarrow$  if at most  $d-1$  errors happen they are detected  
(output is not a codeword)

Error correction  $\rightarrow$   $d=2t+1$  code can correct up to  $t$  errors

$$C = \{10001, 01101, 00110, 11010\}$$

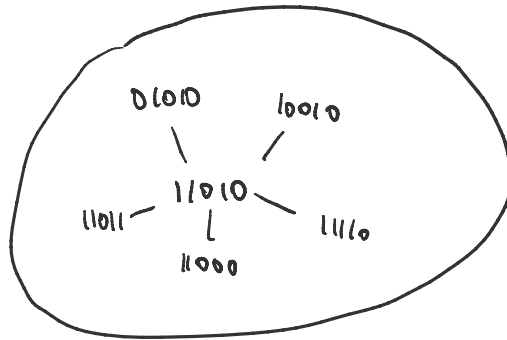
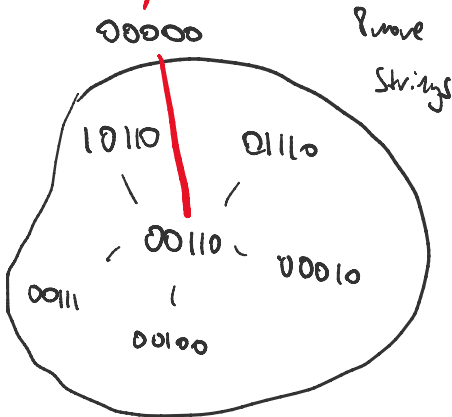
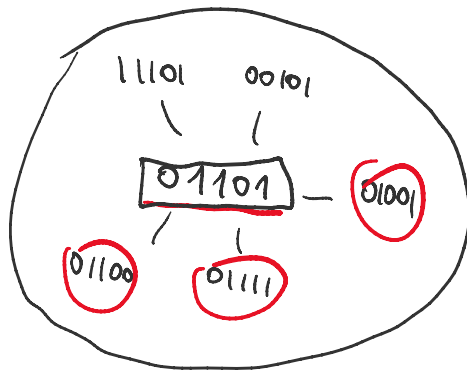
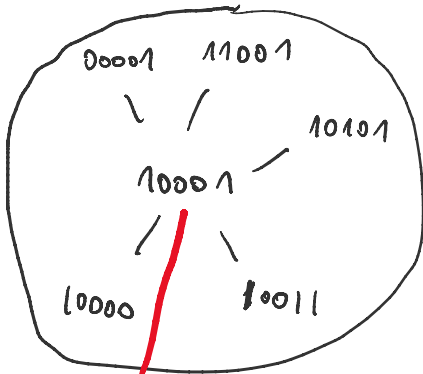
$$\text{Ham}(10101, 10001) = 1$$

MLP

$$\text{Ham}(01101, 10101) = 2 \Rightarrow \text{decode as } 10001$$

$$\text{Ham}(00110, 10101) = 3$$

$$\text{Ham}(11010, 10101) = 4$$



$A_q(n, M) =$  the largest  $d$  of a code with codewords of length  $n$  (over  $q$ -ary alphabet)

alphabet

size (usually  $q=2$ )

$\rightarrow n, M, d$



---

Code equivalence (slide 53)

$C_1$  from  $C_0$  by

- 1.) permutations of the positions of the code
- 2.) permutation of symbols appearing in a fixed position

$C_1$  is equivalent to  $C_0$

$$C_0 = \{ \overset{\curvearrowright}{10001}, \overset{\curvearrowright}{01101}, \overset{\curvearrowright}{00110}, \overset{\curvearrowright}{11010} \}$$

$$C_1 = \{ \underline{01001}, \underline{10101}, \underline{00110}, \underline{11010} \}$$

$$C_2 = \{ \underline{00001}, \underline{11101}, \underline{10110}, \underline{01010} \}$$

---

More efficient

Linear codes

Cyclic codes