# Basics of Coding Theory

**1.) NOISLESS CODING THEORT**

→ Shannon Entropies

→ Huffman Coding

**2.) Noisy Coding theory**

→ Error correcting codes

## Noiseless Coding theory

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

$P_0,...,P_{n-1}$

$\sum$ → alphabet $\sum = \{0,1\}$

#codewords

$0 →$ $00$

$1 →$ $01$    is this the most efficient?

$2 →$ $10$

$3 →$ $11$ → Code C

$\log_2 n$ bits to find codewords for $n$ different signals.

$$AVG(C) = \sum_i P_i \cdot |C_i|$$

→ length of codeword for ith signal.

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

2.) How to construct the best code?

1.) for a random variable $X$ w.p. $(P_0,...,P_{n-1})$

$$S(\lambda) = -\sum_i P_i \log_2 P_i$$    Shannon entropy

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

II. Encoding multiple messages together helps

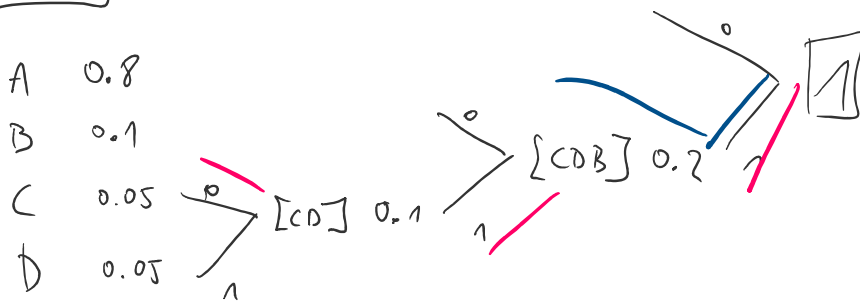III. As the number of messages encoded together approaches infinity $S(X)$ is achievable

2.) Huffman Coding

Alg.

INPUT: Probability distribution
OUTPUT: Optimal code

Ex 1.2

A   0.8
B   0.1
C   0.05
D   0.05

[CD] 0.1

[CDB] 0.2

1

A → 0
B → 10
C → 110
D → 111

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

$= 1.3$

$S(X) = -(0.8) \cdot \overline{\log_2(0.8)}$
$\qquad\quad - (0.1) \cdot \log_2(0.1)$
$\qquad\quad - 2 \cdot (0.05) \cdot \log_2(0.05)$

AA $(0.8)^2$
AB $(0.8) \cdot (0.1)$

$\vdots$

A  $1/3$
B  $1/3$
C  $1/3$

A → 0
B → 10
C → 11

$AVG(c) = 1/3 \cdot 1 + 2 \cdot (1/3 \cdot 2)$
$\qquad\qquad = \frac{5}{3} \approx 1.666$

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

$$= -\log_2 \frac{1}{3} = 1.5\ldots$$

AA
AB
AC
BA
BB
BC
CA
CB
CC

$\frac{1}{9} \rangle \frac{2}{9} \rangle \frac{4}{9}$
$\frac{1}{9}$
$\frac{1}{9} \rangle \frac{2}{9}$
$\frac{1}{9}$
$\frac{1}{9} \rangle \frac{2}{9} \rangle \frac{5}{9} \rangle 1$
$\frac{1}{9}$
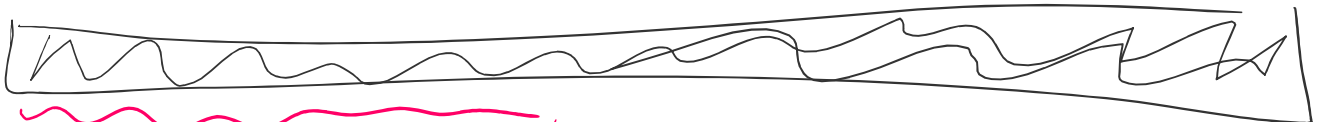$\frac{1}{9} \rangle \frac{3}{9}$
$\frac{1}{9}$
$\frac{1}{9}$

$$AVG(c) = 7 \cdot \left( \frac{1}{9} \cdot 3 \right) + 2 \left( \frac{1}{9} \cdot 4 \right)$$
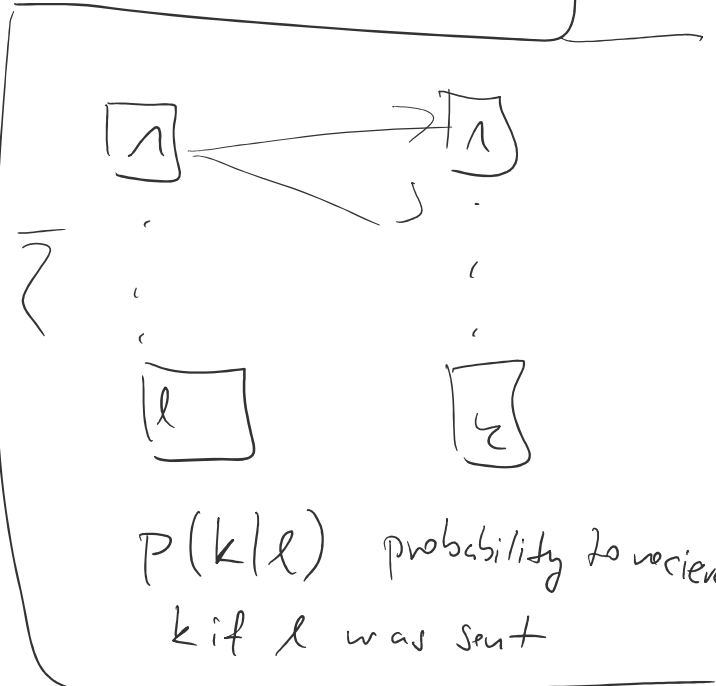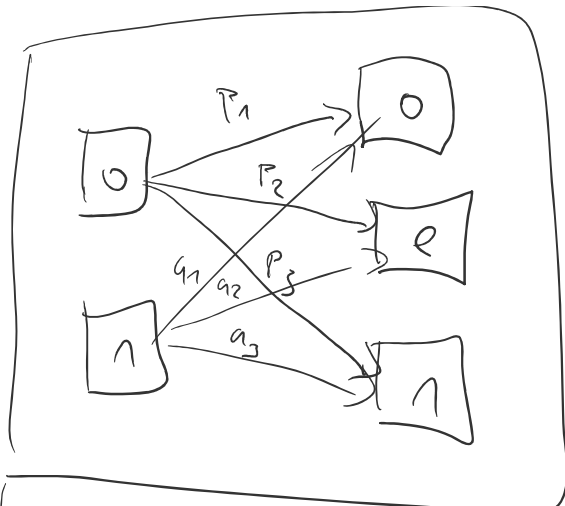
$$= 3.22222$$

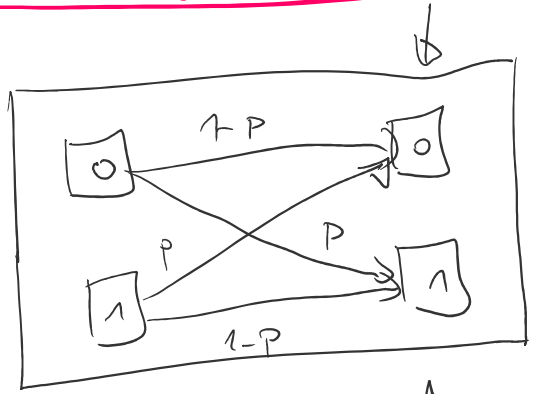$$\frac{AVG(c)}{2} = 1.6111$$

if you encode $k$ symbols together

Then $\quad \lim_{k \to \infty} \frac{AVG(C_{Huff})}{k} = S(X)$

Noisy Coding theory (ERROR CORRECTING CODES)

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

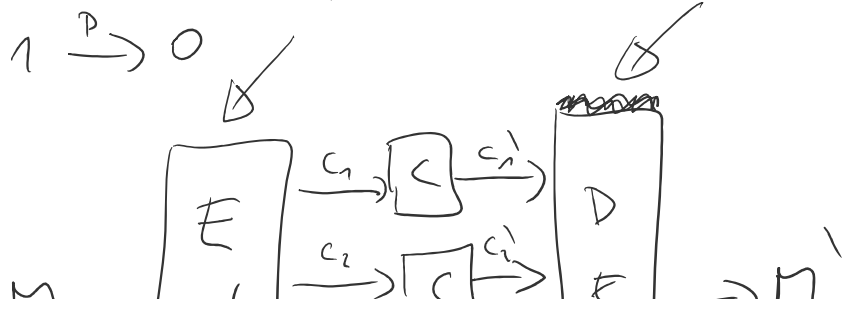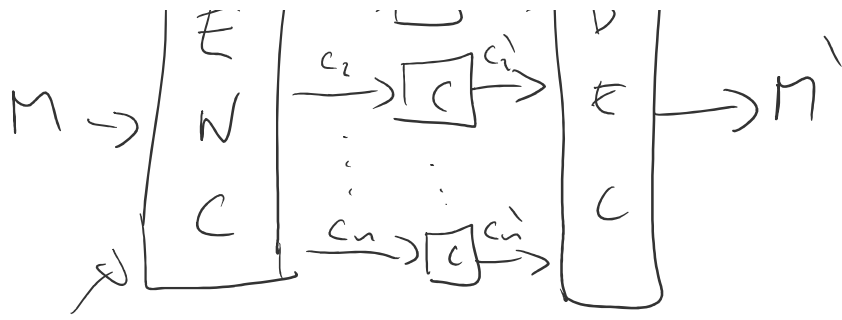## Binary symmetric channel



$P < \frac{1}{2}$

$P(k|\ell)$ probability to recieve $k$ if $\ell$ was sent

# PRINCIPLE OF MAXIMUN LIKELIHOOD

You receive 0 from a binary symmetric channel.
How do you interpret (decode) it?

$0 \xrightarrow{1-p} 0$    $p < \frac{1}{2} \Rightarrow$    $p < 1-p$

$1 \xrightarrow{p} 0$

$M \rightarrow$ | E N C | $\xrightarrow{c_2}$ | C | $\xrightarrow{c_i}$ ... $\xrightarrow{c_n}$ | C | $\xrightarrow{c_i'}$ | D E C | $\rightarrow M'$

---

$0 \rightarrow 000$

$1 \rightarrow 111$

Decoding rule     #1 $\geq 2$     decode as 1

                  #1 $< 2$     decode as 0

$\overset{\text{input}}{0}$    $\overset{\text{output}}{0}$

$Pr(000 | 001) = (1-p)(1-p)p \Leftarrow p < \frac{1}{2}$      $\boxed{01}$

$Pr(111 | 001) = p \cdot p \cdot (1-p)$

---

Repetition code can achieve arbitrary low probability of wrong decoding

$0 \rightarrow \overset{2k+1}{\overbrace{00 \cdots 0}}$

$1 \rightarrow 11 \cdots 1$

Decoding rule     #1 $\geq k+1 \Rightarrow 1$

$\boxed{\# < k+1 \Rightarrow 0}$

$Pr(\text{correct decoding})$

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

$\boxed{000}$  100  010  001

$\boxed{111}$  110  101  011

$= 1$

$\lim_{k \to \infty}$

---

$\dfrac{\# \text{Messages}}{\text{length of codewords}} = \dfrac{2}{k} \underset{k \to \infty}{=} 0$   Code rate

# Hamming Distance

$c_i$ — codewords

$C \subseteq \{0,1\}^n$   codes

$c_i \in C$

$\text{Ham}(c_i, c_j)$ the number of positions in which $c_i$ and $c_j$ __differ__.

EX 1.6   $\boxed{11111}$

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

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

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

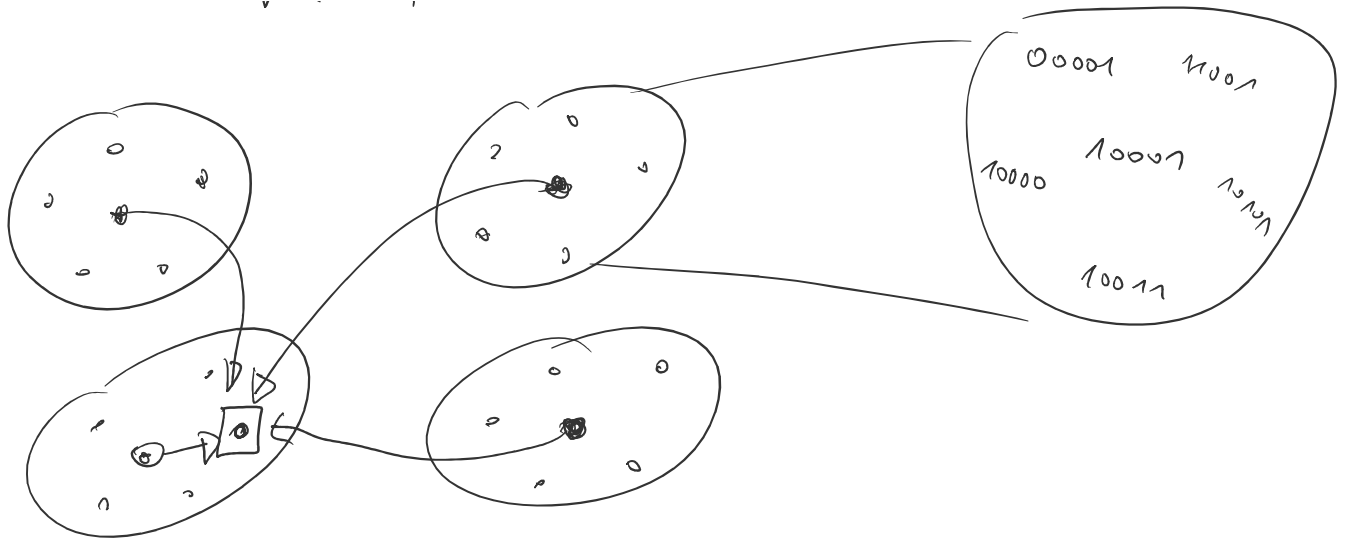$\text{Ham}(10001, 01101) = 3$   $\text{Ham}(11010, 01101) = 4$

$\boxed{\begin{array}{l} n = 5 \\ M = 4 \\ d = 3 \end{array}}$  (length of codewords)
(number of codewords)
(minimal distance)

Error detection $\longrightarrow$ Output of a channel is not a codeword ($d-1$ errors)

Error correction $\longrightarrow$ $d = 2t+1$ code can __correct__ up to $t$ errors

Pr of $e$ (specific errors)   $p^e (1-p)^{n-e}$

$\text{Pr}(e \text{ specific errors}) > \text{Pr}(e+1 \text{ specific errors})$

00001   11001
10000   10001   10101
       10011

10001 →error→ 00001
0
0
;

| 11111 |

---

n $\flat$

M $\phi$

d $\phi$    code alphabet size (usually q=2)

$A_q(n, M) =$ the largest d of a code with M codewords of length n (over q-ary alphabet)
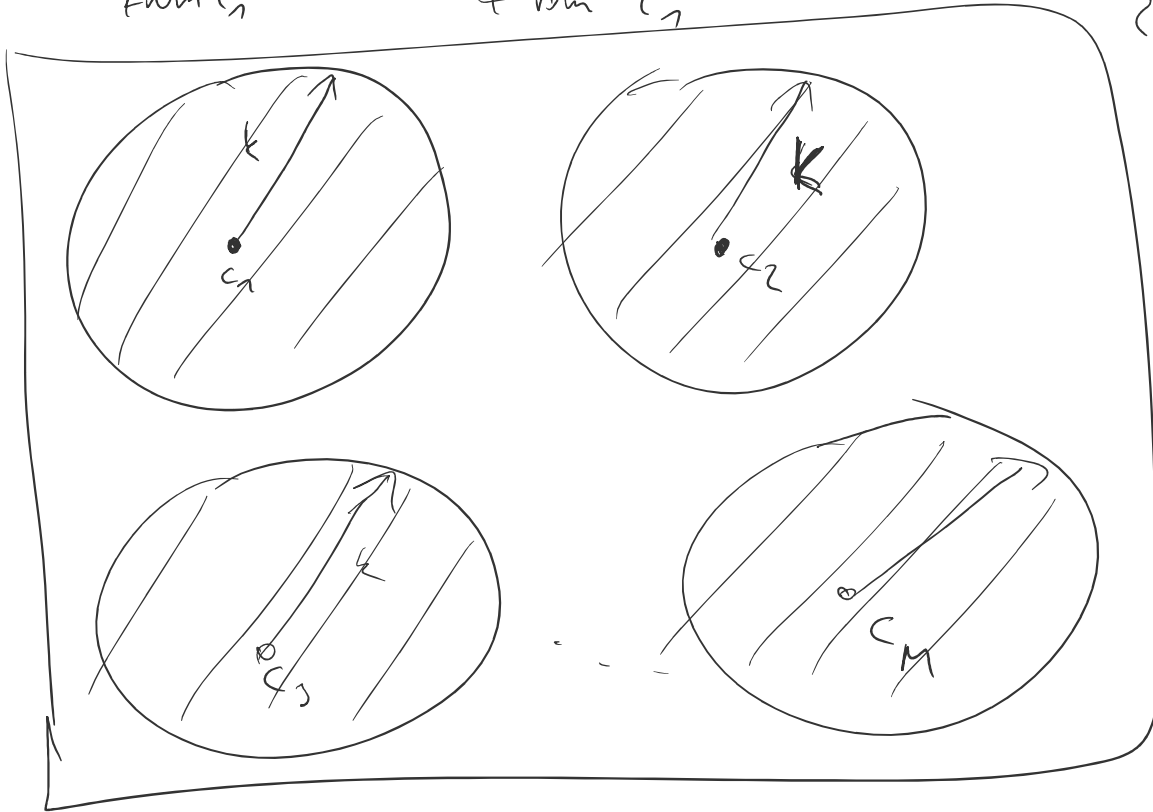
Sphere packing bound    $d = 2t + 1$

$$\boxed{M} \cdot \left[ \binom{n}{0} + \binom{n}{1} \cdot (q-1) + \binom{n}{2} \cdot (q-1)^2 + \cdots + \binom{n}{t} (q-1)^t \right] \leq q^n$$

Choose ↶

the number of codewords of distance 0

the number of codewords of distance 1

2

t

from $c_1$          from $c_2$          $\{0,...,q-1\}^n$



$\downarrow$

$\boxed{q^n}$ — total

number of

strings

$= 2 \cdot 1 + 1 \Rightarrow t = 1$          $q = 2$

$5, 4, 3$ — code
$\quad n \quad M \quad d$

$$4 \left\{ \binom{5}{0} + \binom{5}{1} \right\} \leq 2^5$$

$$4 \{ 1 + 5 \}$$

$$24 \qquad \leq 32$$

$0 \rightarrow 000$

$1 \rightarrow 111$