

Boris Kudryashov

ITMO University

December 15, 2016

- 1 Noiseless coding problem statement
- 2 Channel models
- 3 Mutual information. Average mutual information
- 4 Conditional average mutual information. Information rework theorem
- 5 Convexity of average mutual information
- 6 Information capacity and throughput
- 7 Fano inequality
- 8 Reverse coding theorem
- 9 Information capacity of memoryless channels
- 10 Symmetrical channels

Noiseless coding problem statement

- $X = \{0, 1\}. Y = X$
- Discrete channel with noise.
- Develop a code to eliminate errors.

Noiseless coding problem statement

- $X = \{0, 1\}. Y = X$
- Discrete channel with noise.
- Develop a code to eliminate errors.

Table: Example 1

Message	Codeword	Decisive area
0	000	$\{000, 001, 010, 100\}$
1	111	$\{011, 101, 110, 111\}$

Noiseless coding problem statement

Table: Example 2

Message	Codeword	Decisive area
00	00000	{00000,00001,00010,00100, 01000,10000,11000,10001}
01	10110	{10110,10111,10100,10010, 11110,00110,01110,00111}
10	01011	{01011,01010,01001,01111, 00011,11011,10011,11010}
11	11101	{11101,11100,11111,11001, 10101,01101,00101,01100}

Noiseless coding problem statement

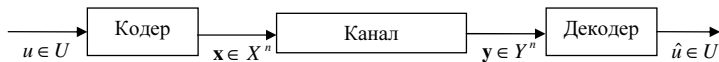


Figure: Communication system Scheme

Noiseless coding problem statement

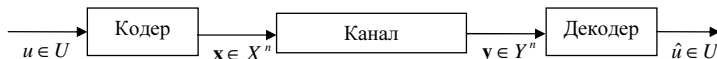


Figure: Communication system Scheme

- *Code of channel* over X is arbitrary set of sequences $A = \{\mathbf{x}_m\}$, $m = 1, \dots, M$, $A \in X^n$.
- These sequences are *codewords*.
- Their length n is *code length*.
- Number of sequences M is *code cardinality*. R , defined as:

$$R = \frac{\log M}{n} \quad (1)$$

is called *code rate* (bits per symbol).

- Event when $\hat{u} \neq u$ is *decoding error*.
- And it's probability is *error probability*

- *Channel model* is defined, if $\forall n$ and $\forall \mathbf{x} \in X^n$, $\mathbf{y} \in Y^n$ conditional probability $p(\mathbf{y}|\mathbf{x})$ is defined.

- *Channel model* is defined, if $\forall n$ and $\forall \mathbf{x} \in X^n$, $\mathbf{y} \in Y^n$ conditional probability $p(\mathbf{y}|\mathbf{x})$ is defined.
- Reminder: $\mathbf{x}_i^n = (x_i, \dots, x_n)$. Channel is called *stationary*, if $\forall j, n$ and $\forall \mathbf{x}_{j+1}^{j+n} \in X^n$, $\mathbf{y}_{j+1}^{j+n} \in Y^n$ conditional probabilities $p(\mathbf{y}_{j+1}^{j+n}|\mathbf{x}_{j+1}^{j+n})$ are defined by sequence characters and do not depend from index j .

- Channel is called *memoryless*, if $\forall j, n$ and $\forall \mathbf{x}_{j+1}^{j+n} \in X^n, \mathbf{y}_{j+1}^{j+n} \in Y^n$

$$p(\mathbf{y}_{j+1}^{j+n} | \mathbf{x}_{j+1}^{j+n}) = \prod_{i=j+1}^{j+n} p(y_i | x_i).$$

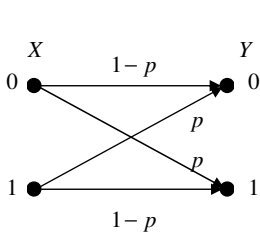
- Channel is called *memoryless*, if $\forall j, n$ and $\forall \mathbf{x}_{j+1}^{j+n} \in X^n, \mathbf{y}_{j+1}^{j+n} \in Y^n$

$$p(\mathbf{y}_{j+1}^{j+n} | \mathbf{x}_{j+1}^{j+n}) = \prod_{i=j+1}^{j+n} p(y_i | x_i).$$

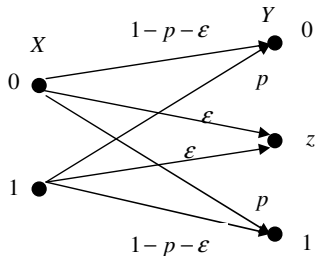
- Stationary channel without memory is called discrete stationary channel.

To describe a Discrete Stationary Channel it's enough to define conditional probabilities $\{p(y|x), x \in X, y \in Y\}$. Let $X = \{0, \dots, K-1\}$, $Y = \{0, \dots, L-1\}$. Let $p_{ij} = p(y=j|x=i)$, $i \in X, j \in Y$. Describe transition probabilities of channel p_{ij} in a *transition probability matrix*:

$$\begin{bmatrix} p_{00} & p_{01} & \cdots & p_{0,L-1} \\ p_{10} & p_{11} & \cdots & p_{1,L-1} \\ \vdots & \vdots & \ddots & \vdots \\ p_{K-1,0} & p_{K-1,1} & \cdots & p_{K-1,L-1} \end{bmatrix}.$$



а) ДСК



б) ДСКН

Figure: Discrete stationary channels examples

- Binary Symmetric Channel (BSC).
 $X = Y = \{0, 1\}$, $p_{10} = p_{01} = p$,
 $p_{00} = p_{11} = 1 - p$. Transition probability matrix:

$$P = \begin{bmatrix} 1 - p & p \\ p & 1 - p \end{bmatrix}.$$

- Binary Symmetric Channel (BSC).
 $X = Y = \{0, 1\}$, $p_{10} = p_{01} = p$,
 $p_{00} = p_{11} = 1 - p$. Transition probability matrix:

$$P = \begin{bmatrix} 1 - p & p \\ p & 1 - p \end{bmatrix}.$$

- Binary Symmetric Channel with Erasure (BSCE).

$$P = \begin{bmatrix} 1 - p - \varepsilon & \varepsilon & p \\ p & \varepsilon & 1 - p - \varepsilon \end{bmatrix}.$$

$X = 0, 1$, $Y = 0, 1, z$, where z is a special erasure symbol.

- For a given $XY = \{(x, y), p(x, y)\}$ of ensembles X and Y calculate the information about $x \in X$ by $y \in Y$.

- For a given $XY = \{(x, y), p(x, y)\}$ of ensembles X and Y calculate the information about $x \in X$ by $y \in Y$.
- Mutual information:

$$I(x; y) = I(x) - I(x|y). \quad (2)$$

- *Average mutual information* of X and Y is

$$I(X; Y) = \mathbf{M}[I(x; y)].$$

- *Average mutual information* of X and Y is

$$I(X; Y) = \mathbf{M}[I(x; y)].$$

- Dependence between average mutual information and joint probability distribution:

$$I(X; Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(y|x)}{p(y)}. \quad (3)$$

Properties of mutual information:

1 Symmetricity: $I(x; y) = I(y; x)$.

Properties of mutual information:

- 1 Symmetricity: $I(x; y) = I(y; x)$.
- 2 If x and y are independent, $I(x, y) = 0$.

Properties of mutual information:

- 1 Symmetricity: $I(x; y) = I(y; x)$.
- 2 If x and y are independent, $I(x, y) = 0$.
- 3 Symmetricity $I(X; Y) = I(Y; X)$.

Properties of mutual information:

- 1 Symmetricity: $I(x; y) = I(y; x)$.
- 2 If x and y are independent, $I(x, y) = 0$.
- 3 Symmetricity $I(X; Y) = I(Y; X)$.
- 4 Nonnegativity: $I(X; Y) \geq 0$.

Properties of mutual information:

- 1 Symmetricity: $I(x; y) = I(y; x)$.
- 2 If x and y are independent, $I(x, y) = 0$.
- 3 Symmetricity $I(X; Y) = I(Y; X)$.
- 4 Nonnegativity: $I(X; Y) \geq 0$.
- 5 Identity $I(X; Y) = 0$ holds iff ensembles X and Y are independent.

Properties of mutual information:

$$\begin{aligned} 6 \quad I(X; Y) &= H(X) - H(X|Y) = \\ &H(Y) - H(Y|X) = H(X) + H(Y) - H(XY). \end{aligned}$$

Properties of mutual information:

- 6 $I(X; Y) = H(X) - H(X|Y) =$
 $H(Y) - H(Y|X) = H(X) + H(Y) - H(XY).$
- 7 $I(X; Y) \leq \min \{H(X), H(Y)\}.$

Properties of mutual information:

- 6 $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = H(X) + H(Y) - H(XY).$
- 7 $I(X; Y) \leq \min \{H(X), H(Y)\}.$
- 8 $I(X; Y) \leq \min \{\log |X|, \log |Y|\}.$

Properties of mutual information:

$$6 \quad I(X; Y) = H(X) - H(X|Y) = \\ H(Y) - H(Y|X) = H(X) + H(Y) - H(XY).$$

$$7 \quad I(X; Y) \leq \min \{H(X), H(Y)\}.$$

$$8 \quad I(X; Y) \leq \min \{\log |X|, \log |Y|\}.$$

9 Mutual information $I(X; Y)$ is a convex \cap function of probability distribution $p(x)$.

Properties of mutual information:

$$6 \quad I(X; Y) = H(X) - H(X|Y) = \\ H(Y) - H(Y|X) = H(X) + H(Y) - H(XY).$$

$$7 \quad I(X; Y) \leq \min \{H(X), H(Y)\}.$$

$$8 \quad I(X; Y) \leq \min \{\log |X|, \log |Y|\}.$$

9 Mutual information $I(X; Y)$ is a convex \cap function of probability distribution $p(x)$.

10 Mutual information $I(X; Y)$ is a convex \cup function of conditional probabilities $p(y|x)$.

Conditional average mutual information.

- Consider $XYZ = \{(x, y, z), p(x, y, z)\}$. Fix $z \in Z$ and consider conditional probability distribution: $p(x, y|z) = \frac{p(x, y, z)}{p(z)}$.
- Average mutual information between X and Y :
$$I(X; Y|z) = \sum_{x \in X} \sum_{y \in Y} p(x, y|z) \log \frac{p(y|x, z)}{p(y|z)}.$$

Conditional average mutual information.

- Conditional average mutual information between X and Y :

$$I(X; Y|Z) = \mathbf{M} [I(X; Y|z)] = \sum_{x \in X} \sum_{y \in Y} \sum_{z \in Z} p(x, y, z) \log \frac{p(y|x, z)}{p(y|z)}$$

- Additional properties:

$$I(X; Y|Z) = H(Y|Z) - H(Y|XZ).$$

$$I(X; YZ) = I(X; Y) + I(X; Z|Y)$$

$$I(X; YZ) = I(X; Z) + I(X; Y|Z)$$

Conditional average mutual information.

A special case of information processing system,
which has 3 probability ensembles:

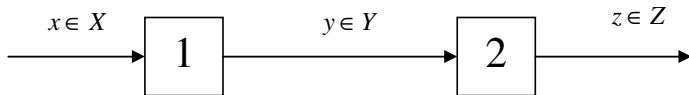


Figure: Information processing system

Conditional average mutual information.

Theorem

Let X, Y, Z be probability ensembles, which are formed by the information processing system at the previous slide. Then holds:

$$I(X; Y) \geq I(X; Z), \quad (4)$$

$$I(Y; Z) \geq I(X; Z). \quad (5)$$

Conditional average mutual information.

proof. Use properties of conditional average mutual information:

$$I(X; YZ) = I(X; Y) + I(X; Z|Y), \quad (6)$$

$$I(X; YZ) = I(X; Z) + I(X; Y|Z). \quad (7)$$

X and Z are independent. If Y is known, $I(X; Z|Y) = 0$. By equating the right sides of (6) and (7), we get

$$I(X; Y) = I(X; Z) + I(X; Y|Z).$$

Since the second term is non-negative, we obtain the inequality (4). Similarly we can prove (5).

Convexity of average mutual information

- Let $\mathbf{p} = (p_0, \dots, p_{K-1})$ be probabilities of input symbols $X = \{0, \dots, K-1\}$. Let use $I(\mathbf{p})$ instead of $I(X; Y)$ to emphasize that we are interested in the dependence between mutual information and input symbols distribution.
- Consider $Z = \{1, 2\}$ such that $p_z(1) = \alpha$, $p_z(2) = 1 - \alpha$
- Consider XYZ , where tuples (x, y, z) are created as follows:
 - (1) z is chosen according to $p(z)$.
 - (2) if $z = 1$, p_1 is used to choose x , otherwise, p_2 is used.
 - (3) After that, according to $p(y|x)$, y element is generated.

Convexity of average mutual information

- From the convexity definition: $\forall \mathbf{p}_1, \mathbf{p}_2, \alpha \in [0, 1]$ holds

$$I(\alpha \mathbf{p}_1 + (1 - \alpha) \mathbf{p}_2) \geq \alpha I(\mathbf{p}_1) + (1 - \alpha) I(\mathbf{p}_2). \quad (8)$$

- According to previous definitions:

$$\begin{aligned} I(X; Y|z = 1) &= I(\mathbf{p}_1); \\ I(X; Y|z = 2) &= I(\mathbf{p}_2); \\ I(X; Y|Z) &= \alpha I(\mathbf{p}_1) + (1 - \alpha) I(\mathbf{p}_2); \\ I(X; Y) &= I(\alpha \mathbf{p}_1 + (1 - \alpha) \mathbf{p}_2). \end{aligned}$$

- Inequation (8) is reduced to:

$$I(X; Y) \geq I(X; Y|Z). \quad (9)$$

Convexity of average mutual information

- consider mutual information $I(Y; XZ)$:

$$I(Y; XZ) = I(Y; X) + I(Y; Z|X); \quad (10)$$

$$I(Y; XZ) = I(Y; Z) + I(Y; X|Z). \quad (11)$$

- As long as Z and Y are independent,
 $I(Y; Z|X) = 0$
- By equating the right sides, we get (9).

Convexity of average mutual information

- Consider mutual information as function of conditional distribution $p(y|x)$.
- $\forall P_1, P_2, \alpha \in [0, 1]$ holds:

$$I(\alpha P_1 + (1 - \alpha)P_2) \leq \alpha I(P_1) + (1 - \alpha)I(P_2). \quad (12)$$

- Consider $Z = \{1, 2\}$. Consider XYZ , where tuples (x, y, z) are created as follows:
 - (1) $x \in X$ is chosen according to $p(x)$.
 - (2) z is chosen according to $p(z)$.
 - (3) Transition probability matrix P is chosen:
 $P = P_1(\text{if } z = 1) \text{ or } P = P_2(\text{if } z = 2)$
 - (4) After that, according to x and P , y element is generated.

Convexity of average mutual information

- According to previous definitions:

$$I(X; Y|z = 1) = I(P_1);$$

$$I(X; Y|z = 2) = I(P_2);$$

$$I(X; Y|Z) = \alpha I(P_1) + (1 - \alpha) I(P_2);$$

$$I(X; Y) = I(\alpha P_1 + (1 - \alpha) P_2).$$

- (12) is now reduced to

$$I(X; Y) \leq I(X; Y|Z). \quad (13)$$

Convexity of average mutual information

- Rewrite mutual information in two ways:

$$I(X; YZ) = I(X; Y) + I(X; Z|Y); \quad (14)$$

$$I(X; YZ) = I(X; Z) + I(X; Y|Z). \quad (15)$$

- By equating the right sides (14) and (15), we get (13) and (12). Thus, we proved convexity of mutual information as a function of conditional distributions.

Information capacity and throughput

- When using codewords of length n , average amount of information, received by decoder will be $I(X^n; Y^n)$ bit. This corresponds to information rate:

$$\frac{1}{n} I(X^n; Y^n) \text{ bit/channel symbol.}$$

- C_0 is called the Information Capacity of channel.

$$C_0 = \sup_n \max_{\{p(\mathbf{x})\}} \frac{1}{n} I(X^n; Y^n) \quad (16)$$

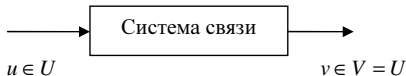


Figure: Information transfer system

- Messages are elements of $U = \{u\} = \{0, \dots, M - 1\}$.
- Receiver gets estimates of messages.
- Estimates are denoted $V = \{v\}$.
- U and V are bijective.
- If $u \neq v$ there is a decoding error.

- $UV = \{(u, v), p(u, v)\}$ and $p(u, v)$ are known.
- Error probability P_e is

$$P_e = \sum_u \sum_{v \neq u} p(u, v). \quad (17)$$

- Probability of correct decoding:

$$P_c = 1 - P_e = \sum_u \sum_{v=u} p(u, v). \quad (18)$$

Theorem

(Fano inequality)

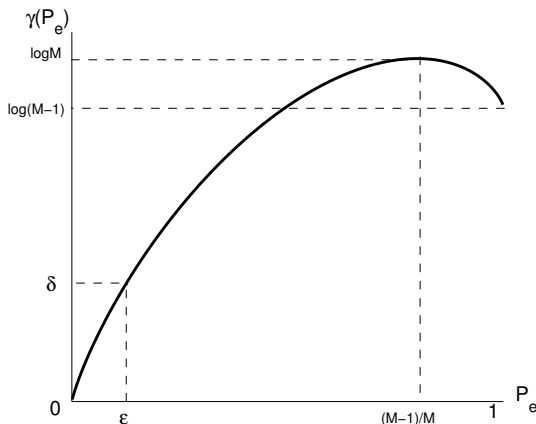
$$H(U|V) \leq \eta(P_e) + P_e \log(M - 1), \quad (19)$$

where $\eta(\cdot)$ denotes entropy of binary ensemble.

Fano inequality

Consider right side of Fano Inequality.

$$\gamma(P_e) = \eta(P_e) + P_e \log(M - 1) \quad (20)$$



Proof of Fano Inequality 2.

- Use (17) and (18), rewrite terms of (19):

$$H(U|V) = - \sum_u \sum_{v \neq u} p(u, v) \log p(u|v) - \sum_u \sum_{v=u} p(u, v) \log p(u|v), \quad (21)$$

$$\eta(P_e) = - \sum_u \sum_{v \neq u} p(u, v) \log P_e - \sum_u \sum_{v=u} p(u, v) \log P_c, \quad (22)$$

$$P_e \log(M-1) = \sum_u \sum_{v \neq u} p(u, v) \log(M-1). \quad (23)$$

- Consider Δ . For (19), we need to prove $\Delta \leq 0$.

$$\Delta = H(U|V) - \eta(P_e) - P_e \log(M - 1).$$

- Subtract from (21) corresponding parts of (22) and (23).

$$\Delta = \sum_u \sum_{v \neq u} p(u, v) \log \frac{P_e}{p(u|v)(M-1)} + \sum_u \sum_{v=u} p(u, v) \log \frac{P_c}{p(u|v)}.$$

- Use $\log x \leq (x - 1) \log e$

$$\begin{aligned} \Delta \leq (\log e) & \left[\sum_u \sum_{v \neq u} p(u, v) \frac{P_e}{p(u|v)(M-1)} - \sum_u \sum_{v \neq u} p(u, v) + \right. \\ & \left. + \sum_u \sum_{v=u} p(u, v) \frac{P_c}{p(u|v)} - \sum_u \sum_{v=u} p(u, v) \right]. \end{aligned}$$

- Use $p(u, v) = p(v)p(u|v)$ and (17) и (18).

$$\Delta \leq \log e \times \left[\frac{P_e}{M-1} \sum_u \sum_{v \neq u} p(v) - P_e + P_c \sum_u \sum_{v=u} p(v) - P_c \right]. \quad (24)$$

- Note, that

$$\sum_u \sum_{v \neq u} p(v) = (M-1) \sum_v p(v) = (M-1). \quad (25)$$

- Moreover

$$\sum_u \sum_{v=u} p(v) = \sum_u p(u) = 1. \quad (26)$$

- Substitute (25) and (26) in (24) and get $\Delta \leq 0$.

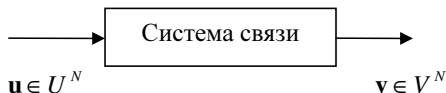


Figure: Система передачи информации

- Let input be a sequence of messages
 $\mathbf{u} = (u_1, \dots, u_N)$
- Let output be a sequence of decisions
 $\mathbf{v} = (v_1, \dots, v_N)$
- $u_i, v_i \in U = V = \{0, \dots, M - 1\}$, $i = 1, \dots, N$
- Let error probability in i -th message be
 $P_{ei} = P(u_i \neq v_i)$

- Let Average error probability of sequence of length N be

$$\bar{P}_e = \frac{1}{N} \sum_{i=1}^N P_{ei}.$$

Theorem

For sequences $(\mathbf{u}, \mathbf{v}) \in U^N V^N$, which consist of elements of M , holds

$$\frac{1}{N} H(U^N | V^N) \leq \eta(\bar{P}_e) + \bar{P}_e \log(M - 1) \quad (27)$$

- Use properties of Conditional Entropy

$$H(U^N|V^N) = \sum_{i=1}^N H(U_i|U_1 \dots U_{i-1} V^N) \leq \sum_{i=1}^N H(U_i|V_i).$$

- Divide both sides on N and use Fano Inequality

$$\frac{1}{N} H(U^N|V^N) \leq \frac{1}{N} \sum_{i=1}^N \eta(P_{ei}) + \frac{1}{N} \sum_{i=1}^N P_{ei} \log(M-1).$$

- As long as entropy is convex \cap function, we get (27) from the last inequality.

Theorem

Reverse coding theorem. For Discrete Memoryless Channel with information capacity C_0 , $\forall \delta > 0 \exists \varepsilon > 0$, such that \forall code with code rate $R > C_0 + \delta$ average error probability satisfies the inequality:

$$\bar{P}_e \geq \varepsilon.$$

Reverse coding theorem

Proof of Reverse Coding Theorem

- $R = \log |C|/n$
- Let $\mathbf{v} \in V^N$ be decoded sequences.

$$\begin{aligned} nR &= \log |C| = \\ &\stackrel{(a)}{=} H(X^n) \stackrel{(b)}{\leq} H(U^N) = \\ &= H(U^N) - H(U^n|V^N) + H(U^N|V^N) = \\ &\stackrel{(c)}{=} I(U^N; V^N) + H(U^N|V^N) \leq \\ &\stackrel{(d)}{=} I(X^n; Y^n) + H(U^N|V^N) \leq \\ &\stackrel{(e)}{\leq} nC_0 + n\gamma(\bar{P}_e). \end{aligned}$$

- $\gamma(\bar{P}_e) \geq R - C_0 > \delta.$

Information capacity of m-l channels

- Conditional probabilities $p(\mathbf{y}|\mathbf{x})$:

$$p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^n p(y_i|x_i). \quad (28)$$

- Information capacity of channel is:

$$C_0 = \sup_n \max_{\{p(\mathbf{x})\}} \frac{1}{n} I(X^n; Y^n). \quad (29)$$

Information capacity of m-l channels

Theorem

Information capacity of discrete memoryless channel can be calculated as:

$$C_0 = \max_{\{p(x)\}} I(X; Y). \quad (30)$$

Information capacity of m-l channels

Proof of theorem (5)

- Mutual information between input and output:

$$I(X^n; Y^n) = H(Y^n) - H(Y^n|X^n). \quad (31)$$

- Use (28)

$$\begin{aligned} H(Y^n|X^n) &= \mathbf{M}[-\log p(\mathbf{y}|\mathbf{x})] = \\ &= \mathbf{M}\left[-\log \prod_{i=1}^n p(y_i|x_i)\right] = \\ &= \sum_{i=1}^n \mathbf{M}[-\log p(y_i|x_i)] = \\ &= \sum_{i=1}^n H(Y_i|X_i). \end{aligned}$$

Information capacity of m-l channels

- Use properties of Entropy:

$$H(Y^n) \leq \sum_{i=1}^n H(Y_i), \quad (32)$$

- Take into account: (31)

$$I(X^n; Y^n) \leq \sum_{i=1}^n [H(Y_i) - H(Y_i|X_i)] = \sum_{i=1}^n I(X_i; Y_i). \quad (33)$$

Information capacity of m-l channels

- Input distribution is:

$$p(\mathbf{y}) = \sum_{\mathbf{x} \in X^n} p(\mathbf{x}) p(\mathbf{y} | \mathbf{x}).$$

- Assume that input characters are independent:

$$\begin{aligned} p(\mathbf{y}) &= \sum_{\mathbf{x} \in X^n} \prod_{i=1}^n p(x_i) \prod_{i=1}^n p(y_i | x_i) = \\ &= \sum_{\mathbf{x} \in X^n} \prod_{i=1}^n p(x_i) p(y_i | x_i) = \\ &= \sum_{x_1 \in X} \sum_{x_2 \in X} \cdots \sum_{x_n \in X} p(x_1) p(y_1 | x_1) \cdot p(x_2) p(y_2 | x_2) \cdot \dots \\ &\quad \dots \cdot p(x_n) p(y_n | x_n). \end{aligned}$$

Information capacity of m-l channels



$$p(\mathbf{y}) = \prod_{i=1}^n \sum_{x_i \in X} p(x_i) p(y_i | x_i) = \prod_{i=1}^n p(y_i),$$

- Substitute (33) into (29):

$$C_0 = \sup_n \max_{\{p(\mathbf{x})\}} \frac{1}{n} \sum_{i=1}^n I(X_i; Y_i).$$

Information capacity of m-l channels

- Search for maximum independently for each term:

$$C_0 = \sup_n \frac{1}{n} \sum_{i=1}^n \max_{\{p(x_i)\}} I(X_i; Y_i).$$

- As long as we have memoryless channel,

$$C_0 = \sup_n \max_{\{p(x)\}} I(X; Y) = \max_{\{p(x)\}} I(X; Y).$$

Information capacity of m-l channels

Information capacity of m-l channels

Information capacity of m-l channels