

Delivering Information Reliably

I-Hsiang Wang

Department of Electrical Engineering
National Taiwan University

ihwang@ntu.edu.tw

October 19, 2023

The information processing task motivating the study of this lecture:

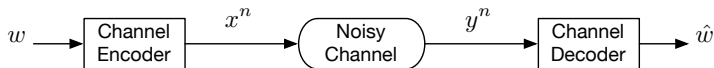
For a length- k i.i.d. $\text{Ber}(\frac{1}{2})$ bit sequence, given a noisy channel with input alphabet \mathcal{X} and output alphabet \mathcal{Y} , design a good encoding scheme to represent the k -bit sequence using n symbols in \mathcal{X} , and a good decoding scheme to reconstruct the k -bit sequence reliably from the observed n symbols in \mathcal{Y} at the output of the channel.

Note: how the encoding and decoding scheme work are known *a priori*.

Fundamental Questions:

- What is the maximum possible ratio $\frac{k}{n}$ (transmission rate) ?
- How to achieve that fundamental limit?

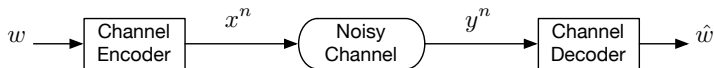
The channel coding problem (Shannon's abstraction)



Meta Description

- 1 Message:** Random message $W \sim \text{Unif}\{1, \dots, 2^k\}$.
- 2 Channel:** Consist of an input alphabet \mathcal{X} , an output alphabet \mathcal{Y} , and a family of conditional laws $\{P_{Y_k|X^k, Y^{k-1}} \mid k \in \mathbb{N}\}$ determining the stochastic relationship between the output symbol Y_k and the input symbol X_k along with all past signals (X^{k-1}, Y^{k-1}) .
- 3 Encoder:** Encode the message w by a length- n **codeword** $x^n \in \mathcal{X}^n$.
- 4 Decoder:** Reconstruct message \hat{w} from the channel output y^n .
- 5 Efficiency:** Maximize the **code rate** $R := \frac{k}{n}$ bits/channel use, given a certain decoding criterion.

Decoding criterion: small error probability



A key performance measure: Error Probability $P_e^{(n)} := \Pr\{W \neq \hat{W}\}$.

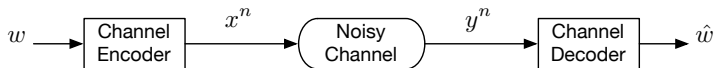
Question: Is it possible to get zero error probability?

Answer: Probably not, unless the channel noise has some special structure.

Following the development of lossless source coding, Shannon turned the attention to answering the following question:

Is it possible to have a sequence of encoder/decoder pairs such that $P_e^{(n)} \rightarrow 0$ as $n \rightarrow \infty$? If so, what is the largest possible code rate R where vanishing error probability is possible?

Encoding and decoding




Decoder takes the observation Y^n generated from the noisy channel and make a decision

$$\hat{W} = \text{dec}(Y^n) \in \mathcal{W} \equiv \{1, \dots, 2^k\}.$$

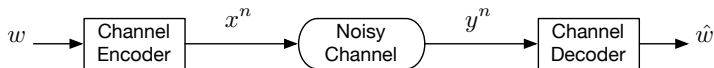
This is just a 2^k -ary detection (hypothesis testing) problem, with

$$\mathcal{H}_w : Y^n \sim P_w \equiv P_{Y^n|X^n}(\cdot|x^n(w)), \quad w = 1, \dots, 2^k.$$

Looks pretty simple – just use maximum likelihood. 

But the analysis of the error probability is tricky.

Encoding and decoding (cont'd)



The even more tricky part is the encoder.

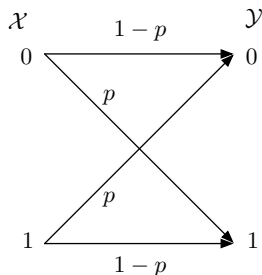
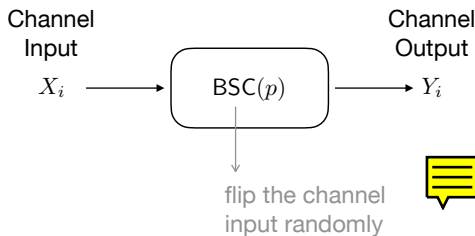
Encoder places the 2^k *codewords* in \mathcal{X}^n so that after passing through the noisy channel $P_{Y^n|X^n}$, the distributions of the observations

$$\left\{ P_w(\cdot) \equiv P_{Y^n|X^n}(\cdot | x^n(w)) \mid w = 1, \dots, 2^k \right\}$$

are **well-separated**.

To make these challenges concrete and see that there are ways to resolve them, let us begin with an example – the **binary symmetric channel (BSC)**, one of the simplest yet non-trivial noisy channel.

Motivating example: the binary symmetric channel



- Binary input/output: $x_i, y_i \in \{0, 1\}$, $i = 1, 2, \dots$
- Channel flips the input bit i.i.d. with probability $p \in (0, 1/2)$:

$$Y_i = X_i \oplus Z_i, \quad Z_i \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p), \quad i = 1, 2, \dots$$

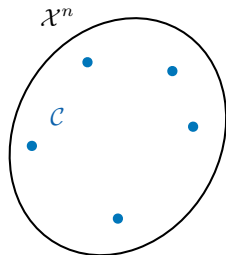
Encoding

$$\begin{array}{ccc} W \equiv [B_1 & B_2 & \dots & B_k] & \longrightarrow & \boxed{\text{enc}} & \longrightarrow & [X_1 & X_2 & \dots & \dots & X_n] \\ \text{message} & & & & & & & & \text{codeword} \\ B \sim \text{Unif}(\{0, 1\}^k) & & & & & & & & \mathbf{X} \in \{0, 1\}^n \end{array}$$

- Goal: introduce redundancy to combat channel noises (for BSC, bit flips).
- Codebook: the range of the encoding function enc :

$$\mathcal{C} = \{ \text{enc}(w) \equiv \mathbf{x}(w) \mid w \in \mathcal{W} \equiv \{1, \dots, 2^k\} \}.$$

- eg. Repetition coding $\mathcal{C} = \{[0 \dots 0], [1 \dots 1]\}$.



Pictorially, the encoding function should place the codewords far apart so that they do not get confused with one another.

Decoding

$$\begin{array}{ccc} [Y_1 & Y_2 & \dots & \dots & Y_n] & \longrightarrow & \boxed{\text{dec}} & \longrightarrow & [\hat{B}_1 & \hat{B}_2 & \dots & \hat{B}_k] \equiv \hat{W} \\ \text{received signal} & & & & & & & & \text{decoded message} \\ Y \in \{0, 1\}^n & & & & & & & & \hat{W} \in \mathcal{W} \end{array}$$

- Performance metric: probability of error $P_e^{(n)} = \Pr\{\hat{W} \neq W\}$.
- Since $W \sim \text{Unif}(\mathcal{W})$ (uniform prior), *maximum likelihood* decoding is the optimal: $\hat{w}_{\text{ML}}(\mathbf{y}) = \arg \max_{w \in \mathcal{W}} P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}(w))$.
- For BSC, likelihood function is simple and determined by the *Hamming distance* between \mathbf{y} and the candidate codeword $\mathbf{x}(w)$:

$$P_{\mathbf{Y}|\mathbf{X}}(\mathbf{y}|\mathbf{x}(w)) = p^{\text{d}_H(\mathbf{y}, \mathbf{x}(w))} (1-p)^{n-\text{d}_H(\mathbf{y}, \mathbf{x}(w))} = (1-p)^n \left(\frac{p}{1-p}\right)^{\text{d}_H(\mathbf{y}, \mathbf{x}(w))}.$$

Since $0 < p < 1/2$, we have

$$\hat{w}_{\text{ML}}(\mathbf{y}) = \arg \min_{w \in \mathcal{W}} \text{d}_H(\mathbf{y}, \mathbf{x}(w))$$

Error probability analysis

Under ML decoding, the error probability is given as

$$P_{e,ML}^{(n)} = \frac{1}{2^k} \sum_{w \in \mathcal{W}} \underbrace{\Pr \{ \exists \tilde{w} \neq w \text{ such that } d_H(\mathbf{Y}, \mathbf{x}(\tilde{w})) \leq d_H(\mathbf{Y}, \mathbf{x}(w)) \mid W = w \}}_{P_{e,ML}^{(n)}(w)}$$

The event that $\mathbf{x}(\tilde{w})$ is closer to \mathbf{Y} than the actual $\mathbf{x}(w)$ is just the event that BSC flips more than half of the bits at which $\mathbf{x}(w)$ and $\mathbf{x}(\tilde{w})$ differ.

Hence by the union bound, we may get to

$$P_{e,ML}^{(n)}(w) \leq \sum_{\substack{\tilde{w} \in \mathcal{W} \\ \tilde{w} \neq w}} \sum_{j=\lceil \frac{d_H(\mathbf{x}(\tilde{w}), \mathbf{x}(w))}{2} \rceil}^{d_H(\mathbf{x}(\tilde{w}), \mathbf{x}(w))} \binom{d_H(\mathbf{x}(\tilde{w}), \mathbf{x}(w))}{j} p^j (1-p)^{d_H(\mathbf{x}(\tilde{w}), \mathbf{x}(w))-j}.$$

It largely depends on how the encoding function populates the codewords in the codebook and hard to analyze, let alone the asymptote as $n \rightarrow \infty$.

Reflections

- Recall: for lossless source coding, we leveraged the *concentration of probability* (AEP) when $n \rightarrow \infty$ to achieve data compression.

Reason: distribution of the random source is known and “well-behaved”.

- But for the channel coding problem, the encoding function enc makes the distribution of coded symbols non-i.i.d. in general, and it becomes hard to control the statistical behavior of the codewords.

$$\underset{\text{i.i.d. over time}}{\mathbf{B}} \longrightarrow \boxed{\text{enc}} \longrightarrow \underset{\text{not i.i.d. over time}}{\mathbf{x}(\mathbf{B})} = \begin{bmatrix} x_1(\mathbf{B}) & x_2(\mathbf{B}) & \dots & x_n(\mathbf{B}) \end{bmatrix}$$

Work-around:

- 1 “Random” encoding
- 2 “Typicality” decoding

Typicality decoding

Observation: bit-flips are i.i.d. $\text{Ber}(p)$, and hence the number of bit-flips follows $\text{Binom}(n, p)$ and concentrates at np when $n \rightarrow \infty$. So, it is “typical” that there are $\approx np$ bit-flips, that is, $\forall \varepsilon > 0$,

$$\Pr\left\{\overbrace{w(\mathbf{Z})}^{\text{\# of bit-flips}} \leq n(p + \varepsilon)\right\} \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

Typicality decoding:

$\hat{w}_T =$ the unique $w \in \mathcal{W}$ such that $d_H(\mathbf{y}, \mathbf{x}(w)) \leq n(p + \varepsilon)$.

Obviously, $P_{e, \text{ML}}^{(n)}(w) \leq P_{e, T}^{(n)}(w)$. Furthermore, by the union bound,

$$\begin{aligned} P_{e, T}^{(n)}(w) &\leq \underbrace{\Pr\left\{\overbrace{w(\mathbf{Z})}^{d_H(\mathbf{x}(w) \oplus \mathbf{Z}, \mathbf{x}(w))} > n(p + \varepsilon)\right\}}_{\rightarrow 0 \text{ as } n \rightarrow \infty} \\ &\quad + \sum_{\tilde{w} \in \mathcal{W}, \tilde{w} \neq w} \Pr\{d_H(\mathbf{x}(w) \oplus \mathbf{Z}, \mathbf{x}(\tilde{w})) \leq n(p + \varepsilon)\} \end{aligned} \quad (*)$$

Random encoding

$$\Pr\{d_H(\mathbf{x}(w) \oplus \mathbf{Z}, \mathbf{x}(\tilde{w})) \leq n(p + \varepsilon)\}$$

This term remains difficult to analyze as it depends on the structure of the codebook (more specifically, pairwise Hamming distance between codewords).

To overcome the difficulty, Shannon came up with a *random coding* idea:

*To prove the existences of codebook a codeb with $P_e^{(n)} \leq \epsilon$, it suffices to show that over a **set** of codebooks, the **averaged** $P_e^{(n)} \leq \epsilon$.*

One may consider the codebook $\mathcal{C} = \{\mathbf{x}(1), \dots, \mathbf{x}(2^k)\}$ as a $2^k \times n$ matrix

$$\mathbf{C} = \begin{bmatrix} x_1(1) & x_2(1) & \cdots & x_n(1) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(2^k) & x_2(2^k) & \cdots & x_n(2^k) \end{bmatrix}$$

Random codebook simply means that the codebook matrix is random and follows a certain distribution: $\mathbf{C} \sim P_{\mathbf{C}} \in \mathcal{P}(\{0, 1\}^{2^k \times n})$.

With the random codebook matrix \mathbf{C} , the goal now turns to proving

$$\mathbb{E}_{\mathbf{C} \sim \mathbf{P}_{\mathbf{C}}, W \sim \text{Unif}(\mathcal{W})} \left[\mathbf{P}_{\mathbf{e}, \mathbf{T}}^{(n)}(W; \mathbf{C}) \right] \leq \epsilon.$$

Clearly, \mathbf{C} is chosen to be independent of W . Hence, it suffices to show that

$$\forall w \in \mathcal{W}, \mathbb{E}_{\mathbf{C} \sim \mathbf{P}_{\mathbf{C}}} \left[\mathbf{P}_{\mathbf{e}, \mathbf{T}}^{(n)}(w; \mathbf{C}) \right] \leq \epsilon.$$

It boils down to show that $\mathbb{E}_{\mathbf{C} \sim \mathbf{P}_{\mathbf{C}}} [(*)] \leq \epsilon/2$, that is,

$$\sum_{\tilde{w} \in \mathcal{W}, \tilde{w} \neq w} \mathbf{P}_{\mathbf{C}, \mathbf{Z}} \left\{ \underbrace{\mathbf{d}_H(\mathbf{X}(w) \oplus \mathbf{Z}, \mathbf{X}(\tilde{w}))}_{\mathbf{w}(\mathbf{X}(w) \oplus \mathbf{Z} \oplus \mathbf{X}(\tilde{w}))} \leq n(p + \varepsilon) \right\} \leq \epsilon/2.$$

Let's pick a distribution on \mathbf{C} so that the analysis becomes simple:

$$X_i(w) \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(1/2) \quad \forall i = 1, 2, \dots, n, \quad \forall w = 1, 2, \dots, 2^k,$$

that is, the entries of \mathbf{C} are i.i.d. $\text{Ber}(1/2)$. Then, each codeword (row of \mathbf{C}) is a length- n i.i.d. $\text{Ber}(1/2)$ sequence.

Since $\mathbf{X}(w)$, $\mathbf{X}(\tilde{w})$, \mathbf{Z} are independent, $\mathbf{X}(w) \oplus \mathbf{Z} \oplus \mathbf{X}(\tilde{w})$ is also a length- n i.i.d. $\text{Ber}(1/2)$ sequence.

As a result, $\forall \tilde{w} \neq w$,

$$\begin{aligned} & P_{\mathbf{C}, \mathbf{Z}} \{ \mathbf{w}(\mathbf{X}(w) \oplus \mathbf{Z} \oplus \mathbf{X}(\tilde{w})) \leq n(p + \varepsilon) \} \\ &= \Pr \{ \text{Binom}(n, \tfrac{1}{2}) \leq n(p + \varepsilon) \} \leq 2^{-n d_b(p + \varepsilon \| \frac{1}{2})} \end{aligned} \quad (\text{Chernoff Bound})$$

Wrapping-up: achievability

Putting everything together, we show that \exists a codebook \mathcal{C} of size 2^k such that

$$P_{e,ML}^{(n)} \leq \epsilon/2 + (2^k - 1)2^{-nd_b(p+\varepsilon \parallel \frac{1}{2})}.$$

A sufficient condition for the second term to vanish is $k < nd_b(p + \varepsilon \parallel \frac{1}{2})$.

Choosing $\varepsilon > 0$ judiciously close to 0, we can show that $\forall \delta > 0$ and $\forall \epsilon \in (0, 1)$, there exists a codebook \mathcal{C} of size 2^k such that

$$k > n \left(d_b(p \parallel \tfrac{1}{2}) - \delta \right) \quad \text{and} \quad P_{e,ML}^{(n)} \leq \epsilon.$$

Hence,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} k^*(n, \epsilon) \geq d_b(p \parallel \tfrac{1}{2}),$$

establishing the achievability part of the coding theorem for BSC(p).

Optimality

Is it true that $\lim_{n \rightarrow \infty} \frac{1}{n} k^*(n, \epsilon) = d_b(p \parallel \frac{1}{2}) = D(\text{Ber}(p) \parallel \text{Ber}(\frac{1}{2}))$?

The answer is **yes**. The divergence reminds us of binary hypothesis testing, and the proof below is built on the idea of *testing between two channel laws*.

For a given encoding/decoding pair (enc, dec) with $P_e^{(n)} \leq \epsilon$, consider a binary hypothesis testing problem

$$\begin{aligned}\mathcal{H}_0 : \quad & W \rightarrow \boxed{\text{enc}} \rightarrow \mathbf{X} \rightarrow \boxed{\text{BSC}(p)^{\otimes n}} \rightarrow \mathbf{Y} \rightarrow \boxed{\text{dec}} \rightarrow \hat{W} \\ \mathcal{H}_1 : \quad & W \rightarrow \boxed{\text{enc}} \rightarrow \mathbf{X} \rightarrow \boxed{\text{BSC}(1/2)^{\otimes n}} \rightarrow \mathbf{Y} \rightarrow \boxed{\text{dec}} \rightarrow \hat{W}\end{aligned}$$

with observation tuple $(W, \mathbf{X}, \mathbf{Y}, \hat{W})$ following P_0 and P_1 respectively.

Key observation 1: the second channel $\text{BSC}(1/2)$ completely breaks the dependency between \mathbf{Y} and \mathbf{X} , and hence the decoding performance is the same as pure random guess, that is, $P_1\{W = \hat{W}\} = \frac{1}{2^k}$.

For the above HT problem, one can use a test that accept \mathcal{H}_0 iff $W = \hat{W}$, and this test has

$$\begin{aligned}\pi_{1|0}^{(n)} &= \mathbb{P}_0\{W \neq \hat{W}\} \equiv \mathbb{P}_e^{(n)} \leq \epsilon \\ \pi_{0|1}^{(n)} &= \mathbb{P}_1\{W = \hat{W}\} = \frac{1}{2^k}.\end{aligned}$$

With the language of binary HT, we immediately have $\frac{1}{2^k} \geq \varpi_{0|1}^*(n, \epsilon)$.

Key observation 2: the above binary HT problem is equivalent to

$$\begin{aligned}\mathcal{H}_0 : \quad Z_i &\stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p), \quad i = 1, 2, \dots, n \\ \mathcal{H}_1 : \quad Z_i &\stackrel{\text{i.i.d.}}{\sim} \text{Ber}\left(\frac{1}{2}\right), \quad i = 1, 2, \dots, n.\end{aligned}$$

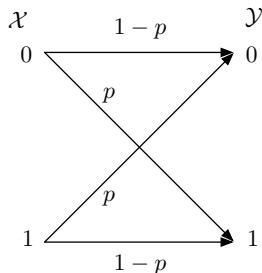
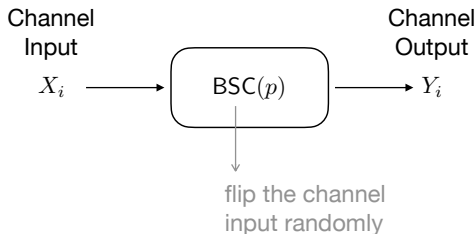
Hence, by the Chernoff-Stein Lemma, $\lim_{n \rightarrow \infty} -\frac{1}{n} \log \varpi_{0|1}^*(n, \epsilon) = d_b(p \parallel \frac{1}{2})$.

Combining everything together, we show that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} k^*(n, \epsilon) \leq \lim_{n \rightarrow \infty} -\frac{1}{n} \log \varpi_{0|1}^*(n, \epsilon) = d_b(p \parallel \frac{1}{2}),$$

and the optimal rate for communication over the BSC is characterized.

Summary



For the binary symmetric channel with flip probability p ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} k^*(n, \epsilon) = d_b(p \| \tfrac{1}{2}) = 1 - H_b(p) \quad \forall \epsilon \in (0, 1).$$

where $k^*(n, \epsilon)$ is the smallest k such that there exists channel codes of codebook size 2^k and error probability $\leq \epsilon$.

How to extend the result to general channels?

Outline

In this lecture, we will show that (for discrete memoryless channels) the fundamental limit is the **channel capacity**, a quantity that can be computed by maximizing the “**mutual information**” between the input and the output of the channel, when we want to reconstruct the bits with **vanishing error probability**.

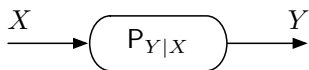
- 1 We begin with an intuitive motivation to introduce **mutual information**.
 - Motivation: single use of a channel.
 - Data processing inequality, chain rule, convexity.
- 2 Next we prove Shannon's noisy channel coding theorem.
 - Use Fano's inequality to prove the converse part.
 - Use typicality arguments to prove the achievability part.
 - Source-channel separation.

1 Mutual Information

- Definitions
- Properties

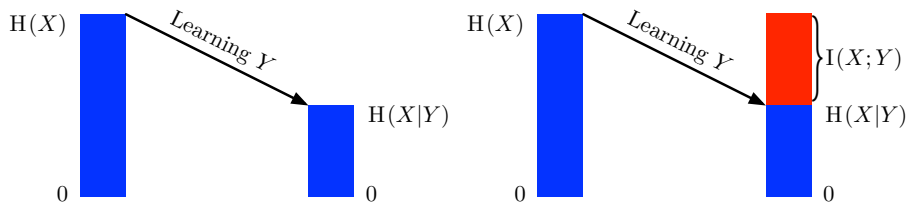
2 Noisy Channel Coding Theorem

Consider a single use of the channel $P_{Y|X}$:



$H(X)$ quantifies the amount of uncertainty of the input X .

$H(X|Y)$ quantifies the amount of uncertainty of X given the output Y .



Question: How much information does Y tell about X through the channel?

Answer: $H(X) - H(X|Y)$.

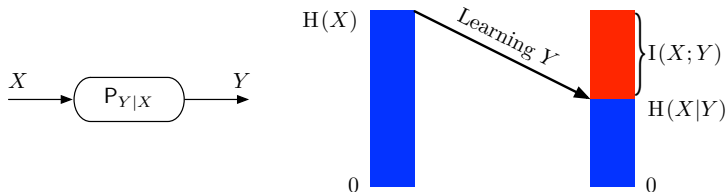
1 Mutual Information

- Definitions

- Properties

2 Noisy Channel Coding Theorem

Mutual information



Definition 1 (Mutual Information)

For a pair of jointly distributed (X, Y) , the mutual information between X and Y is defined as

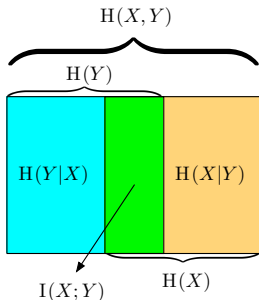
$$I(X; Y) := H(X) - H(X|Y).$$

Interpretation: Infer information about the channel input X from output Y .

Proposition 1 (An Identity)

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

pf: By chain rule: $H(X|Y) = H(X, Y) - H(Y)$. \square



Note: Mutual information is symmetric, that is,

$$I(X; Y) = I(Y; X).$$

The mutual information between X and itself is equal to its entropy:

$$I(X; X) = H(X) \text{ since } H(X|X) = 0.$$

Hence, the entropy is also called “self information” in some literatures.

Mutual information and information divergence

Proposition 2

For $(X, Y) \sim P_{X,Y} = P_X P_{Y|X} = P_Y P_{X|Y}$,

$$I(X; Y) = E_{(X,Y) \sim P_{X,Y}} \left[\log \frac{P_{X,Y}(X, Y)}{P_X(X) P_Y(Y)} \right] = D(P_{X,Y} \| P_X \times P_Y),$$

where $P_X \times P_Y$ is a product distribution of the two marginals P_X and P_Y .

pf: A simple corollary of the identity in Proposition 1. □

Corollary 1 (Extremal Values of Mutual Information)

- 1** $I(X; Y) \geq 0$, with equality iff X, Y are independent.
- 2** $I(X; Y) \leq H(X)$, with equality iff X is a deterministic function of Y .

pf: The proof of the first one is due to the fact that conditioning reduces entropy. The proof of the second one is due to $H(X|Y) \geq 0$.

Interpretation: the mutual information between X and Y , $I(X; Y)$ can also be viewed as a measure of the **dependency** between X and Y .

It is the divergence of the actual distribution $P_{X,Y}$ from the independent distribution $P_X \times P_Y$.

- If X is determined by Y (highly dependent), $I(X; Y)$ is maximized.
- If X is independent of Y (no dependency), $I(X; Y) = 0$.

Exercise 1

Prove the following identity:

$$I(X; Y) = D(P_{Y|X} \| P_Y | P_X) = D(P_{Y|X} \| Q_Y | P_X) - D(P_Y \| Q_Y),$$

$\forall Q_Y$ such that $D(P_Y \| Q_Y) < \infty$. Furthermore,

$$I(X; Y) = \min_{Q_Y: D(P_Y \| Q_Y) < \infty} D(P_{Y|X} \| Q_Y | P_X).$$

Example 1

Compute $I(X_1; X_2)$ for $X_1, X_2 \in \{0, 1\}$ with joint PMF

(x_1, x_2)	$(0, 0)$	$(0, 1)$	$(1, 0)$	$(1, 1)$
$P(x_1, x_2)$	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{6}$

sol: From the previous examples, we have

$$\begin{aligned} H(X_1, X_2) &= \log 3 + \frac{1}{3}, \quad H(X_1) = H(X_2) = 1, \\ H(X_1|X_2) &= H(X_2|X_1) = \log 3 - \frac{2}{3}. \end{aligned}$$

Hence, $I(X_1; X_2) = H(X_1) - H(X_1|X_2) = \frac{5}{3} - \log 3$.

Conditional mutual information

Definition 2 (Conditional Mutual Information)

For a tuple of jointly distributed r.v.'s (X, Y, Z) , the mutual information between X and Y given Z is defined as

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

Similar to the previous identity (Proposition 1), we have

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

Similar to Proposition ??, we have

- 1 $I(X; Y|Z) \geq 0$, with equality iff X, Y are independent given Z , that is, $X - Z - Y$ forms a Markov chain.
- 2 $I(X; Y|Z) \leq H(X|Z)$, with equality iff X is a (deterministic) function of Y, Z .

1 Mutual Information

- Definitions

- Properties

2 Noisy Channel Coding Theorem

Chain rule for mutual information

Theorem 1 (Chain Rule for Mutual Information)

$$I(X; Y^n) = \sum_{i=1}^n I(X; Y_i | Y^{i-1}).$$

pf: Proved by definition and the chain rule for entropy. □

Exercise 2

Show that

$$I(X; Z) \leq I(X; Y, Z) \quad \text{and} \quad I(X; Y|Z) \leq I(X; Y, Z).$$

Data processing inequality

Theorem 2 (Data Processing Inequality)

For a Markov chain $X - Y - Z$, that is, $P_{X,Y,Z} = P_X P_{Y|X} P_{Z|Y}$, we have

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

Interpretation: $X - Y - Z$ says that the information of X that Z can provide is contained in Y . Hence, the amount of information of X that can be inferred by Z is not greater than that can be inferred by Y .

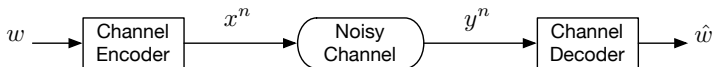
pf: Since $X - Y - Z$, we have $I(X; Z|Y) = 0$. Hence,

$$I(X; Y, Z) = I(X; Y) + I(X; Z|Y) = I(X; Y) \quad (\because I(X; Z|Y) = 0)$$

$$I(X; Y, Z) = I(X; Z) + I(X; Y|Z) \quad (\text{Chain Rule})$$

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





Markov chains are common in communication systems. For example, in channel coding (without feedback), the message W , the channel input X^n , the channel output Y^n , and the decoded message \hat{W} form a Markov chain $W - X^n - Y^n - \hat{W}$ (see the figure above).

Data processing inequality turns out to be crucial in obtaining **impossibility** results in information theory.

Exercise 3 (Functions of R.V.)

For $Z = g(Y)$ being a deterministic function of Y , show that

$$H(Y) \geq H(Z) \quad \text{and} \quad I(X; Y) \geq I(X; Z).$$

Exercise 4

Show that

$$X_1 - X_2 - X_3 - X_4 \implies I(X_1; X_4) \leq I(X_2; X_3).$$

Example 2

Consider two random variables $X_1, X_2 \in \{0, 1\}$ with the same joint PMF as that in Example 1. Let $X_3 = X_2 \oplus Z$, where $Z \sim \text{Ber}(p)$ and Z is independent of (X_1, X_2) .

- 1 Compute $I(X_1; X_3)$ and $I(X_1; X_2|X_3)$.
- 2 Show that $X_1 - X_2 - X_3$ forms a Markov chain.
- 3 Verify the data processing inequality $I(X_1; X_2) \geq I(X_1; X_3)$.

sol:

(x_1, x_2, x_3)	$(0, 0, 0)$	$(0, 0, 1)$	$(0, 1, 0)$	$(0, 1, 1)$
$P(x_1, x_2, x_3)$	$\frac{1}{6}(1-p)$	$\frac{1}{6}p$	$\frac{1}{3}p$	$\frac{1}{3}(1-p)$
(x_1, x_2, x_3)	$(1, 0, 0)$	$(1, 0, 1)$	$(1, 1, 0)$	$(1, 1, 1)$
$P(x_1, x_2, x_3)$	$\frac{1}{3}(1-p)$	$\frac{1}{3}p$	$\frac{1}{6}p$	$\frac{1}{6}(1-p)$

Then it is straightforward to compute mutual informations and verify the Markov chain $X_1 - X_2 - X_3$.

Conditioning reduces mutual information?

Sometimes yes, sometimes no.

Proposition 3 (Conditioning May Decrease Mutual Information)

For a Markov chain $X - Y - Z$, we have $I(X; Y) \geq I(X; Y|Z)$.

pf: The same argument as that of the data processing inequality. □

Example 3 (Conditioning May Increase Mutual Information)

Let X and Y be i.i.d. $\text{Ber}(\frac{1}{2})$ random variables, and $Z = X \oplus Y$. Evaluate $I(X; Y|Z)$ and show that $I(X; Y|Z) > I(X; Y)$.

sol:
$$\begin{aligned} I(X; Y|Z) &= H(X|Z) - H(X|Y, Z) = H(X|Z) - H(X|Y, X \oplus Y) \\ &= H(X|Z) - H(X|Y, X) = H(X|Z) = H(X) = 1. \end{aligned}$$

(note that X and Z are independent)

On the other hand, $I(X; Y) = 0$. Hence, $1 = I(X; Y|Z) > I(X; Y) = 0$.

Convexity and concavity of mutual information

Theorem 3

Let $(X, Y) \sim P_{X,Y} = P_X P_{Y|X}$.

- 1 With $P_{Y|X}$ fixed, $I(X; Y)$ is a **concave** function of P_X .
- 2 With P_X fixed, $I(X; Y)$ is a **convex** function of $P_{Y|X}$.

pf: For the first property, note that

$$I(X; Y) = H(Y) - H(Y|X) = H(P_Y) - \underbrace{\sum_{x \in \mathcal{X}} P_X(x) H(Y|X=x)}_{\text{linear in } P_X}.$$

It suffices to show that $H(P_Y)$ is concave in P_X .

P_Y is a linear function of P_X and $H(P_Y)$ is a concave function of P_Y . As a result, $H(P_Y)$ is concave in P_X , which completes the proof of concavity.

For the convexity, let us consider two conditional distributions $P_{Y|X}^{(0)}, P_{Y|X}^{(1)}$, and random variables X, Y, Θ such that $P_{\Theta, X, Y} = P_X P_{\Theta} P_{Y|X, \Theta}$, $\Theta \sim \text{Ber}(\lambda)$, and $P_{Y|X, \Theta}(y | x, \theta) = P_{Y|X}^{(\theta)}(y | x)$, $\forall y \in \mathcal{Y}, x \in \mathcal{X}, \theta \in \{0, 1\}$.

Hence, $\Theta \perp\!\!\!\perp X$ and $P_{Y|X} = (1 - \lambda)P_{Y|X}^{(0)} + \lambda P_{Y|X}^{(1)}$.

By the basic properties of mutual information and entropy,

$$I(X; Y) \leq I(X; Y, \Theta) = I(X; Y | \Theta) + I(X; \Theta).$$

Note that $\Theta \perp\!\!\!\perp X \implies I(X; \Theta) = 0$. Also note that

- $I(X; Y) = I(X; Y) |_{(X, Y) \sim P_X \{(1-\lambda)P_{Y|X}^{(0)} + \lambda P_{Y|X}^{(1)}\}}$.
- $I(X; Y | \Theta) = (1 - \lambda)I(X; Y | \Theta = 0) + \lambda I(X; Y | \Theta = 1)$ and

$$I(X; Y | \Theta = \theta) = I(X; Y) |_{(X, Y) \sim P_X P_{Y|X}^{(\theta)}}, \theta \in \{0, 1\}.$$

The proof is complete. □

Summary: mutual Information

- $I(X; Y) := H(X) - H(X|Y)$ measures the amount of information of X contained in Y (and vice versa). $I(X; Y) = D(P_{X,Y} \| P_X \times P_Y)$ also tells the level of dependency between X and Y .
- $I(X; Y|Z) := H(X|Z) - H(X|Y, Z)$ measures the amount of information of X in Y given Z .
- Nonnegative; Concave function of P_X with fixed $P_{Y|X}$; Convex function of $P_{Y|X}$ with fixed P_X .
- Chain rule: $I(X; Y^n) = \sum_{i=1}^n I(X; Y_i | Y^{i-1})$.
- Conditioning may reduce MI: $X - Y - Z \implies I(X; Y|Z) \leq I(X; Y)$.
- Data processing decreases MI: $X - Y - Z \implies I(X; Y) \geq I(X; Z)$.

Next: noisy channel coding theorem.

1 Mutual Information

- Definitions
- Properties

2 Noisy Channel Coding Theorem