# Delivering Information Reliably

I-Hsiang Wang

Department of Electrical Engineering
National Taiwan University

ihwang@ntu.edu.tw

October 26, 2023

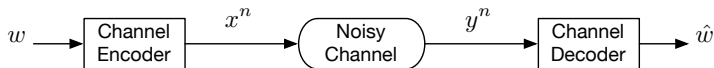The information processing task motivating the study of this lecture:

*For a length-$k$ i.i.d. $\mathrm{Ber}(\frac{1}{2})$ 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) ?
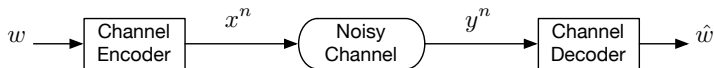- How to achieve that fundamental limit?

# The channel coding problem (Shannon's abstraction)



## Meta Description

1. **Message**: Random message $W \sim \mathrm{Unif}\{1, \ldots, 2^k\}$.

2. **Channel**: Consist of an input alphabet $\mathcal{X}$, an output alphabet $\mathcal{Y}$, and a family of conditional laws $\{P_{Y_i|X^i, Y^{i-1}} \mid i \in \mathbb{N}\}$ determining the stochastic relationship between the output symbol $Y_i$ and the input symbol $X_i$ along with all past signals $(X^{i-1}, Y^{i-1})$.

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 $\mathrm{R} := \frac{k}{n}$ bits/channel use, given a certain decoding criterion.

# Decoding criterion: small error probability



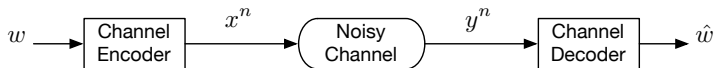A key performance measure: Error Probability $P_e^{(n)} := \Pr\{W \neq \hat{W}\}$.

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 $P_e^{(n)} \to 0$ as $n \to \infty$? If so, what is the largest possible code rate $R$ where vanishing error probability is possible?*

# Encoding and decoding



$$w \longrightarrow \boxed{\begin{array}{c}\text{Channel}\\\text{Encoder}\end{array}} \xrightarrow{x^n} \left(\begin{array}{c}\text{Noisy}\\\text{Channel}\end{array}\right) \xrightarrow{y^n} \boxed{\begin{array}{c}\text{Channel}\\\text{Decoder}\end{array}} \longrightarrow \hat{w}$$

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

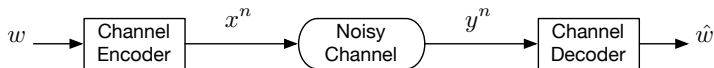$$\hat{W} = \text{dec}(Y^n) \in \mathcal{W} \equiv \{1, \ldots, 2^k\}.$$

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

$$\mathcal{H}_w : Y^n \sim \mathsf{P}_w \equiv \mathsf{P}_{Y^n|X^n}(\cdot|x^n(w)), \ w = 1, \ldots, 2^k.$$

Looks pretty simple – just use maximum likelihood.

But the analysis of the error probability is tricky.

# Encoding and decoding (cont'd)



$$w \longrightarrow \boxed{\begin{array}{c}\text{Channel}\\\text{Encoder}\end{array}} \xrightarrow{x^n} \left(\begin{array}{c}\text{Noisy}\\\text{Channel}\end{array}\right) \xrightarrow{y^n} \boxed{\begin{array}{c}\text{Channel}\\\text{Decoder}\end{array}} \longrightarrow \hat{w}$$

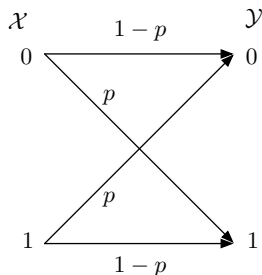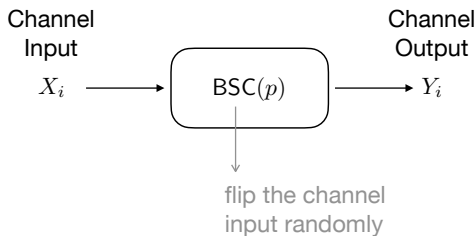The even more tricky part is the encoder.

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

$$\left\{ \mathsf{P}_w(\cdot) \equiv \mathsf{P}_{Y^n|X^n}(\cdot|x^n(w)) \,\Big|\, w = 1, \ldots, 2^k \right\}$$

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

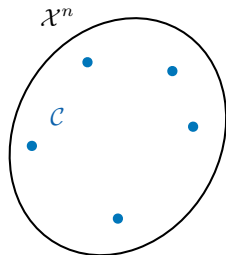$$Y_i = X_i \oplus Z_i, \ Z_i \overset{\text{i.i.d.}}{\sim} \text{Ber}(p), \ i = 1, 2, \ldots.$$

**Encoding**

$$W \equiv \begin{bmatrix} B_1 & B_2 & \ldots & B_k \end{bmatrix} \longrightarrow \boxed{\text{enc}} \longrightarrow \begin{bmatrix} X_1 & X_2 & \ldots & \ldots & X_n \end{bmatrix}$$

message                  codeword

$\boldsymbol{B} \sim \text{Unif}\big(\{0,1\}^k\big)$          $\boldsymbol{X} \in \{0,1\}^n$

- Goal: introduce redundancy to combat channel noises (for BSC, bit flips).

- Codebook: the range of the encoding function $\text{enc}$:

  $$\mathcal{C} = \big\{ \text{enc}(w) \equiv \boldsymbol{x}(w) \,\big|\, w \in \mathcal{W} \equiv \{1, \ldots, 2^k\}. \big\}.$$

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



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

## Decoding

$$\begin{bmatrix} Y_1 & Y_2 & \ldots & \ldots & Y_n \end{bmatrix} \longrightarrow \boxed{\text{dec}} \longrightarrow \begin{bmatrix} \hat{B}_1 & \hat{B}_2 & \ldots & \hat{B}_k \end{bmatrix} \equiv \hat{W}$$

received signal $\qquad\qquad\qquad\qquad$ decoded message

$\boldsymbol{Y} \in \{0,1\}^n \qquad\qquad\qquad\qquad\qquad \hat{W} \in \mathcal{W}$

- Performance metic: probability of error $\mathsf{P}_e^{(n)} = \Pr\{\hat{W} \neq W\}$.

- Since $W \sim \mathrm{Unif}(\mathcal{W})$ (uniform prior), *maximum likelihood* decoding is the optimal: $\hat{w}_{\mathrm{ML}}(\boldsymbol{y}) = \arg\max_{w \in \mathcal{W}} \mathsf{P}_{\boldsymbol{Y}|\boldsymbol{X}}(\boldsymbol{y}|\boldsymbol{x}(w))$.

- For BSC, likelihood function is simple and determined by the *Hamming distance* between $\boldsymbol{y}$ and the candidate codeword $\boldsymbol{x}(w)$:

$$\mathsf{P}_{\boldsymbol{Y}|\boldsymbol{X}}(\boldsymbol{y}|\boldsymbol{x}(w)) = p^{\mathrm{d_H}(\boldsymbol{y},\boldsymbol{x}(w))}(1-p)^{n-\mathrm{d_H}(\boldsymbol{y},\boldsymbol{x}(w))} = (1-p)^n \left(\frac{p}{1-p}\right)^{\mathrm{d_H}(\boldsymbol{y},\boldsymbol{x}(w))}.$$

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

$$\hat{w}_{\mathrm{ML}}(\boldsymbol{y}) = \arg\min_{w \in \mathcal{W}} \mathrm{d_H}(\boldsymbol{y}, \boldsymbol{x}(w))$$

## Error probability analysis

Under ML decoding, the error probability is given as

$$\mathsf{P}_{e,\mathrm{ML}}^{(n)} = \frac{1}{2^k} \sum_{w \in \mathcal{W}} \underbrace{\mathsf{Pr}\left\{\exists\, \tilde{w} \neq w \text{ such that } \mathrm{d_H}(\boldsymbol{Y}, \boldsymbol{x}(\tilde{w})) \leq \mathrm{d_H}(\boldsymbol{Y}, \boldsymbol{x}(w)) \,\middle|\, W = w\right\}}_{\mathsf{P}_{e,\mathrm{ML}}^{(n)}(w)}$$

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

Hence by the union bound, we may get to

$$\mathsf{P}_{e,\mathrm{ML}}^{(n)}(w) \leq \sum_{\substack{\tilde{w} \in \mathcal{W} \\ \tilde{w} \neq w}} \sum_{j=\left\lceil \frac{\mathrm{d_H}(\boldsymbol{x}(\tilde{w}), \boldsymbol{x}(w))}{2} \right\rceil}^{\mathrm{d_H}(\boldsymbol{x}(\tilde{w}), \boldsymbol{x}(w))} \binom{\mathrm{d_H}(\boldsymbol{x}(\tilde{w}), \boldsymbol{x}(w))}{j} p^j (1-p)^{\mathrm{d_H}(\boldsymbol{x}(\tilde{w}), \boldsymbol{x}(w))-j}.$$

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

## Reflections

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

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

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

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

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 $n \to \infty$. So, it is "typical" that there are $\approx np$ bit-flips, that is, $\forall\, \varepsilon > 0$,

$$\Pr\Big\{ \overbrace{\text{w}(\boldsymbol{Z})}^{\text{\# of bit-flips}} \leq n(p+\varepsilon) \Big\} \to 1 \quad \text{as } n \to \infty.$$

Typicality decoding:

$$\hat{w}_{\text{T}} = \text{the unique } w \in \mathcal{W} \text{ such that } d_{\text{H}}(\boldsymbol{y}, \boldsymbol{x}(w)) \leq n(p+\varepsilon).$$

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

$$\mathsf{P}_{\text{e,T}}^{(n)}(w) \leq \underbrace{\Pr\Big\{ \overbrace{\text{w}(\boldsymbol{Z})}^{d_{\text{H}}(\boldsymbol{x}(w) \oplus \boldsymbol{Z},\, \boldsymbol{x}(w))} > n(p+\varepsilon) \Big\}}_{\to 0 \text{ as } n \to \infty}$$
$$+ \sum_{\tilde{w} \in \mathcal{W}, \tilde{w} \neq w} \Pr\big\{ d_{\text{H}}(\boldsymbol{x}(w) \oplus \boldsymbol{Z},\, \boldsymbol{x}(\tilde{w})) \leq n(p+\varepsilon) \big\} \qquad (*)$$

**Random encoding**

$$\Pr\big\{\mathrm{d}_{\mathrm{H}}(\boldsymbol{x}(w) \oplus \boldsymbol{Z}, \boldsymbol{x}(\tilde{w})) \le n(p + \varepsilon)\big\}$$

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 with* $\mathsf{P}_{\mathrm{e}}^{(n)} \le \epsilon$*, it suffices to show that over a set of codebooks, the averaged* $\mathsf{P}_{\mathrm{e}}^{(n)} \le \epsilon$*.*

One may consider the codebook $\mathcal{C} = \{\boldsymbol{x}(1), \ldots, \boldsymbol{x}(2^k)\}$ 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: $\mathbf{C} \sim \mathsf{P}_{\mathbf{C}} \in \mathcal{P}(\{0,1\}^{2^k \times n})$.

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

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

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

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

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

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

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

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

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

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

$$\begin{aligned}
&\mathsf{P}_{\mathbf{C}, \boldsymbol{Z}}\big\{ \mathrm{w}(\boldsymbol{X}(w) \oplus \boldsymbol{Z} \oplus \boldsymbol{X}(\tilde{w})) \leq n(p + \varepsilon) \big\} \\
&= \mathsf{Pr}\big\{ \text{Binom}(n, \tfrac{1}{2}) \leq n(p + \varepsilon) \big\} \leq 2^{-n\mathsf{d_b}\left(p+\varepsilon \,\middle\|\, \frac{1}{2}\right)} \qquad \text{(Chernoff Bound)}
\end{aligned}$$

## Wrapping-up: achievability

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

$$\mathsf{P}_{e,\mathrm{ML}}^{(n)} \leq \epsilon/2 + (2^k - 1)2^{-n\mathsf{d}_\mathsf{b}\left(p+\varepsilon \,\big\|\, \frac{1}{2}\right)}.$$

A sufficient condition for the second term to vanish is $k < n\mathsf{d}_\mathsf{b}\left(p + \varepsilon \,\big\|\, \frac{1}{2}\right)$.

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

$$k > n\left(\mathsf{d}_\mathsf{b}\left(p \,\big\|\, \frac{1}{2}\right) - \delta\right) \quad \text{and} \quad \mathsf{P}_{e,\mathrm{ML}}^{(n)} \leq \epsilon.$$

Hence,

$$\liminf_{n\to\infty} \tfrac{1}{n} k^*(n, \epsilon) \geq \mathsf{d}_\mathsf{b}\left(p \,\big\|\, \tfrac{1}{2}\right),$$

establishing the achievability part of the coding theorem for $\mathrm{BSC}(p)$.

## Optimality

Is it true that $\lim_{n \to \infty} \frac{1}{n} k^*(n, \epsilon) = \mathsf{d}_\mathsf{b}\left(p \,\middle\|\, \frac{1}{2}\right) = \mathrm{D}\left(\mathrm{Ber}(p) \,\middle\|\, \mathrm{Ber}(\frac{1}{2})\right)$?

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 $(\mathrm{enc}, \mathrm{dec})$ with $\mathsf{P}_\mathsf{e}^{(n)} \le \epsilon$, consider a binary hypothesis testing problem

$$\mathcal{H}_0: \quad W \to \boxed{\mathrm{enc}} \to \boldsymbol{X} \to \boxed{\mathsf{BSC}(p)^{\otimes n}} \to \boldsymbol{Y} \to \boxed{\mathrm{dec}} \to \hat{W}$$

$$\mathcal{H}_1: \quad W \to \boxed{\mathrm{enc}} \to \boldsymbol{X} \to \boxed{\mathsf{BSC}(1/2)^{\otimes n}} \to \boldsymbol{Y} \to \boxed{\mathrm{dec}} \to \hat{W}$$

with observation tuple $(W, \boldsymbol{X}, \boldsymbol{Y}, \hat{W})$ following $\mathsf{P}_0$ and $\mathsf{P}_1$ respectively.

**Key observation 1**: the second channel $\mathsf{BSC}(1/2)$ completely breaks the dependency between $\boldsymbol{Y}$ and $\boldsymbol{X}$, and hence the decoding performance is the same as pure random guess, that is, $\mathsf{P}_1\{W = \hat{W}\} = \frac{1}{2^k}$.

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

$$\pi_{1|0}^{(n)} = \mathsf{P}_0\{W \neq \hat{W}\} \equiv \mathsf{P}_{\mathsf{e}}^{(n)} \leq \epsilon$$

$$\pi_{0|1}^{(n)} = \mathsf{P}_1\{W = \hat{W}\} = \tfrac{1}{2^k}.$$

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

$$\mathcal{H}_0 : \quad Z_i \overset{\text{i.i.d.}}{\sim} \mathrm{Ber}(p), \ i = 1, 2, \ldots, n$$
$$\mathcal{H}_1 : \quad Z_i \overset{\text{i.i.d.}}{\sim} \mathrm{Ber}(\tfrac{1}{2}), \ i = 1, 2, \ldots, n.$$

Hence, by the Chernoff-Stein Lemma, $\displaystyle\lim_{n \to \infty} -\tfrac{1}{n} \log \varpi_{0|1}^*(n, \epsilon) = \mathsf{d}_{\mathsf{b}}\big(p \,\big\|\, \tfrac{1}{2}\big)$.

Combining everything together, we show that

$$\limsup_{n \to \infty} \tfrac{1}{n} k^*(n, \epsilon) \leq \lim_{n \to \infty} -\tfrac{1}{n} \log \varpi_{0|1}^*(n, \epsilon) = \mathsf{d}_{\mathsf{b}}\big(p \,\big\|\, \tfrac{1}{2}\big),$$

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

## Summary



For the binary symmetric channel with flip probability $p$,

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

where $k^*(n,\epsilon)$ is the largest $k$ 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.
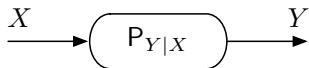   - Source-channel separation.

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

# Mutual information

---

### Definition 1 (Mutual Information)

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

$$\mathrm{I}(X;Y) := \mathrm{H}(X) - \mathrm{H}(X|Y).$$

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

**Proposition 1 (An Identity)**

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

**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 $(X, Y) \sim \mathsf{P}_{X,Y} = \mathsf{P}_X \mathsf{P}_{Y|X} = \mathsf{P}_Y \mathsf{P}_{X|Y}$,*

$$\mathrm{I}(X;Y) = \mathsf{E}_{(X,Y) \sim \mathsf{P}_{X,Y}} \left[ \log \frac{\mathsf{P}_{X,Y}(X,Y)}{\mathsf{P}_X(X)\mathsf{P}_Y(Y)} \right] = \mathrm{D}(\mathsf{P}_{X,Y} \| \mathsf{P}_X \times \mathsf{P}_Y),$$

*where $\mathsf{P}_X \times \mathsf{P}_Y$ is a product distribution of the two marginals $\mathsf{P}_X$ and $\mathsf{P}_Y$.*

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

## Corollary 1 (Extremal Values of Mutual Information)

**1** $\mathrm{I}(X;Y) \geq 0$, *with equality iff $X, Y$ are independent.*

**2** $\mathrm{I}(X;Y) \leq \mathrm{H}(X)$, *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 $\mathrm{H}(X|Y) \geq 0$.

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

It is the divergence of the actual distribution $\mathsf{P}_{X,Y}$ from the independent distribution $\mathsf{P}_X \times \mathsf{P}_Y$.

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

### Exercise 1

Prove the following identity:

$$\mathrm{I}(X;Y) = \mathrm{D}\left(\mathsf{P}_{Y|X}\big\|\mathsf{P}_Y\big|\mathsf{P}_X\right) = \mathrm{D}\left(\mathsf{P}_{Y|X}\big\|\mathsf{Q}_Y\big|\mathsf{P}_X\right) - \mathrm{D}\left(\mathsf{P}_Y\big\|\mathsf{Q}_Y\right),$$

$\forall\, \mathsf{Q}_Y$ such that $\mathrm{D}\left(\mathsf{P}_Y\|\mathsf{Q}_Y\right) < \infty$. Furthermore,

$$\mathrm{I}(X;Y) = \min_{\mathsf{Q}_Y\,:\,\mathrm{D}(\mathsf{P}_Y\|\mathsf{Q}_Y)<\infty} \mathrm{D}\left(\mathsf{P}_{Y|X}\big\|\mathsf{Q}_Y\big|\mathsf{P}_X\right).$$

### Example 1

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

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

**sol**: From the previous examples, we have

$$\mathrm{H}(X_1, X_2) = \log 3 + \tfrac{1}{3}, \ \mathrm{H}(X_1) = \mathrm{H}(X_2) = 1,$$
$$\mathrm{H}(X_1|X_2) = \mathrm{H}(X_2|X_1) = \log 3 - \tfrac{2}{3}.$$

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

# 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

$$\mathrm{I}(X;Y|Z) := \mathrm{H}(X|Z) - \mathrm{H}(X|Y,Z).$$

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

$$\mathrm{I}(X;Y|Z) = \mathrm{H}(X|Z) - \mathrm{H}(X|Y,Z) = \mathrm{H}(Y|Z) - \mathrm{H}(Y|X,Z)$$
$$= \mathrm{H}(X|Z) + \mathrm{H}(Y|Z) - \mathrm{H}(X,Y|Z).$$

Similar to Corollary 1, we have

1. $\mathrm{I}(X;Y|Z) \geq 0$, with equality iff $X, Y$ are independent given $Z$, that is, $X - Z - Y$ forms a Markov chain.

2. $\mathrm{I}(X;Y|Z) \leq \mathrm{H}(X|Z)$, with equality iff $X$ is a (deterministic) function of $Y, Z$.

# Chain rule for mutual information

**Theorem 1 (Chain Rule for Mutual Information)**

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

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

**Exercise 2**

Show that

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

# 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*

$$\mathrm{I}(X;Y) \geq \mathrm{I}(X;Z).$$

**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 $\mathrm{I}(X;Z|Y) = 0$. Hence,

$$\mathrm{I}(X;Y,Z) = \mathrm{I}(X;Y) + \mathrm{I}(X;Z|Y) = \mathrm{I}(X;Y) \qquad (\because \mathrm{I}(X;Z|Y) = 0)$$
$$\mathrm{I}(X;Y,Z) = \mathrm{I}(X;Z) + \mathrm{I}(X;Y|Z) \qquad \qquad (\text{Chain Rule})$$
$$\implies \mathrm{I}(X;Y) = \mathrm{I}(X;Z) + \mathrm{I}(X;Y|Z) \geq \mathrm{I}(X;Z). \qquad \qquad \square$$

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

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

## Exercise 4

Show that

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

### 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 $Z \sim \mathrm{Ber}(p)$ and $Z$ is independent of $(X_1, X_2)$.

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

**sol**:

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

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

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 $\mathrm{I}(X;Y) \geq \mathrm{I}(X;Y|Z)$.*

**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. $\mathrm{Ber}(\frac{1}{2})$ random variables, and $Z = X \oplus Y$. Evaluate $\mathrm{I}(X;Y|Z)$ and show that $\mathrm{I}(X;Y|Z) > \mathrm{I}(X;Y)$.

**sol**: $\mathrm{I}(X;Y|Z) = \mathrm{H}(X|Z) - \mathrm{H}(X|Y,Z) = \mathrm{H}(X|Z) - \mathrm{H}(X|Y, X \oplus Y)$
$$= \mathrm{H}(X|Z) - \mathrm{H}(X|Y,X) = \mathrm{H}(X|Z) = \mathrm{H}(X) = 1.$$

(note that $X$ and $Z$ are independent)

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

# 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 $X, Y, \Theta$ such that $P_{\Theta,X,Y} = P_X P_\Theta P_{Y|X,\Theta}$, $\Theta \sim \mathrm{Ber}(\lambda)$, and $P_{Y|X,\Theta}(y \mid x, \theta) = P_{Y|X}^{(\theta)}(y \mid x)$, $\forall y \in \mathcal{Y}$, $x \in \mathcal{X}$, $\theta \in \{0, 1\}$.

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,

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

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

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

- $\mathrm{I}(X;Y|\Theta) = (1-\lambda)\mathrm{I}(X;Y|\Theta = 0) + \lambda\mathrm{I}(X;Y|\Theta = 1)$ and

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

The proof is complete. $\qquad\square$

# 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.

# Discrete memoryless channel (DMC)

We focus on a main class of channels, *discrete memoryless channels* (DMC), where the perturbation experienced in the channel is memoryless.

**Definition 3 (Discrete Memoryless Channel)**

A discrete channel $\left(\mathcal{X}, \left\{ P_{Y_i|X^i, Y^{i-1}} \,\middle|\, i \in \mathbb{N} \right\}, \mathcal{Y}\right)$ is *memoryless* if $\forall\, i \in \mathbb{N}$,

$$P_{Y_i|X^i, Y^{i-1}}(y_i|x^i, y^{i-1}) = P_{Y|X}(y_i|x_i).$$

In other words, $Y_i - X_i - \left(X^{i-1}, Y^{i-1}\right)$. Here the conditional PMF $P_{Y|X}$ is called the *channel law* or *channel transition function*.

**Note**: The condition in the above definition is **NOT** sufficient to specify

$$P_{Y^n|X^n}(y^n|x^n),$$

the stochastic relationship between the channel input (codeword) $x^n$ and the channel output $y^n$.

Why? Because

$$P(y^n|x^n) = \frac{P(x^n, y^n)}{P(x^n)}$$

$$P(x^n, y^n) = \prod_{i=1}^{n} P(x_i, y_i|x^{i-1}, y^{i-1})$$

$$= \prod_{i=1}^{n} P(y_i|x^i, y^{i-1}) P(x_i|x^{i-1}, y^{i-1})$$

Hence, we need to further specify $\{P(x_i|x^{i-1}, y^{i-1}) \,|\, i \in \mathbb{N}\}$, which cannot be obtained from $P(x^n)$ alone.

**Interpretation**: $\{P(x_i|x^{i-1}, y^{i-1}) \,|\, i \in \mathbb{N}\}$ is induced by the encoding function, which implies that the encoder can potentially make use of the past channel output, i.e., *feedback*.

Let us first focus on the channel coding problem without feedback.

# DMC without feedback



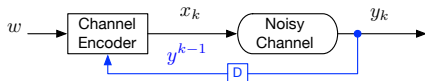No Feedback

With Feedback

Suppose in the system, the encoder has no knowledge about the realization of the channel output, then, $P(x_i | x^{i-1}, y^{i-1}) = P(x_i | x^{i-1})$ for all $i \in \mathbb{N}$.
$\implies$ It is said that the channel has no feedback.

In this case, specifying $\left\{ P(y_i | x^i, y^{i-1}) \,\middle|\, i \in \mathbb{N} \right\}$ suffices to specify $P(y^n | x^n)$.

**Proposition 4 (DMC without Feedback)**

*For a DMC $(\mathcal{X}, P_{Y|X}, \mathcal{Y})$ without feedback,*

$$P(y^n | x^n) = \prod_{i=1}^{n} P_{Y|X}(y_i | x_i).$$

# Channel coding without feedback: problem setup



$$w \longrightarrow \boxed{\begin{array}{c}\text{Channel} \\ \text{Encoder}\end{array}} \xrightarrow{x^n} \left(\begin{array}{c}\text{Noisy} \\ \text{Channel}\end{array}\right) \xrightarrow{y^n} \boxed{\begin{array}{c}\text{Channel} \\ \text{Decoder}\end{array}} \longrightarrow \hat{w}$$

**1** An $(n, k)$ channel code consists of
- an encoding function (encoder) $\mathrm{enc}_n : \mathcal{W} \equiv \{1, 2, ..., 2^k\} \to \mathcal{X}^n$ that maps each message $w$ to a length $n$ codeword $x^n$.
- a decoding function (decoder) $\mathrm{dec}_n : \mathcal{Y}^n \to \mathcal{W}$ that maps a channel output sequence $y^n$ to a reconstructed message $\hat{w}$

**2** An $(n, k)$ code with $\mathsf{P}_{\mathsf{e}}^{(n)} := \Pr\left\{W \neq \hat{W}\right\} \leq \epsilon$ is called an $(n, k, \epsilon)$ code.

**3** Let $k^*(n, \epsilon)$ denote the largest $k$ such that there exists an $(n, k, \epsilon)$ code. The $\epsilon$-capacity $\mathrm{C}(\epsilon) := \lim_{n \to \infty} \frac{1}{n} k^*(n, \epsilon)$ if the limit exists.

Meanwhile, the Shannon capacity $\mathrm{C} := \lim_{\epsilon \downarrow 0} \mathrm{C}(\epsilon)$.

# Channel coding theorem for DMC

**Theorem 4 (Channel Coding Theorem for DMC without Feedback)**

*The $\epsilon$-capacity $\mathrm{C}(\epsilon)$ of the DMC $\left(\mathcal{X}, \mathrm{P}_{Y|X}, \mathcal{Y}\right)$ without feedback is characterized as follows:*

$$\mathrm{C}(\epsilon) = \mathrm{C}^{\mathrm{I}} := \max_{\mathrm{P}_X \in \mathcal{P}(\mathcal{X})} \mathrm{I}(X;Y) \quad \forall \, \epsilon \in (0,1). \tag{1}$$

*Here $\mathcal{P}(\mathcal{X})$ denotes the set of probability distributions over the alphabet $\mathcal{X}$, that is, $\mathcal{P}(\mathcal{X}) \equiv \mathcal{P}_{|\mathcal{X}|}$ is a probability simplex in a $|\mathcal{X}|$-dimensional space.*

**Remarks**:

- Capacity formula (1) is intuitive: $\mathrm{I}(X;Y)$ is just the amount of information about the channel input $X$ that one can infer from the output $Y$.

- The maximization over $\mathrm{P}_X \in \mathcal{P}(\mathcal{X})$ stands for choosing the best possible input distribution so that the amount of information transfer is maximized.

- As a corollary, the Shannon capacity $\mathrm{C} = \max_{\mathrm{P}_X} \mathrm{I}(X;Y)$ as well.

# Shannon capacity and $\epsilon$-capacity: some remarks

In the previous formulation, definition of the Shannon capacity relies on the existence of the $\epsilon$-capacity, that is, the limit $\lim_{n \to \infty} \frac{1}{n} k^*(n, \epsilon)$.

In fact, capacity can be defined alternatively *without relying on the existence of the limit*. For example, in Chapter 19 of Polyanskiy and Wu, $\epsilon$-capacity is defined as the limit inferior $\liminf_{n \to \infty} \frac{1}{n} k^*(n, \epsilon)$ and hence always exist.

Meanwhile, historically, Shannon capacity was investigated first, which can be defined equivalently as follows:

*A rate* $\mathrm{R}$ *is said to be $\epsilon$-achievable (achievable) if there exist a sequence of* $(n, \lceil n\mathrm{R} \rceil)$ *codes such that* $\mathsf{P}_{\mathrm{e}}^{(n)} \leq \epsilon$ *when $n$ is large enough* ($\mathsf{P}_{\mathrm{e}}^{(n)} \to 0$ *as $n \to \infty$).*

*The $\epsilon$-capacity is defined as* $\mathrm{C}(\epsilon) := \sup \{\mathrm{R} \,|\, \mathrm{R} : \epsilon\text{-achievable}\}$.

*The Shannon capacity is defined as* $\mathrm{C} := \sup \{\mathrm{R} \,|\, \mathrm{R} : achievable\}$.

A weaker version of channel coding theorem (Theorem 4) with $C(\epsilon)$ replaced by $C$ was proved (converse: Fano 1952; achievability: Feinstein 1954) first. Later, the converse of the stronger version was proved by Wolfowitz in 1957.

Note: the converse of the stronger version (termed *strong converse* in the literature) equivalently says that if $R > C^I$, $P_e^{(n)} \to 1$, while the weaker version (*weak converse*) only guarantees $P_e^{(n)} \not\to 0$.

In this lecture, we will prove the weak converse using Fano's inequality, and demonstrate how to extend the strong converse argument of BSC capacity only when time permits.

# Roadmap

**1** **Capacity computation**:
Give some simple examples of DMCs to show how to compute capacity.

**2** **Weak converse**:
Prove that $\forall\,R > C^{I}$, it is impossible to have vanishing error probability, that is, $P_e^{(n)} \nrightarrow 0$ as $n \to \infty$.

**3** **Achievability**:
Prove that $\forall\,R < C^{I}$, there exist a sequence of encoding/decoding schemes such that the error probability vanishes, that is, $P_e^{(n)} \to 0$ as $n \to \infty$, based on a probabilistic argument called random coding.
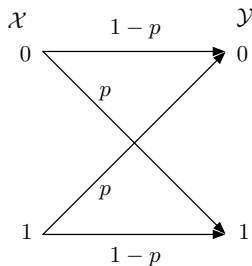
# Binary symmetric channel

A binary symmetric channel (BSC) consists of

- Binary input/output $\mathcal{X} = \mathcal{Y} = \{0, 1\}$.

- Channel law $\mathsf{P}_{Y|X} = \begin{bmatrix} 1-p & p \\ p & 1-p \end{bmatrix}$.

The capacity of BSC is $\mathrm{C}_{\mathrm{BSC}} = 1 - \mathsf{H}_{\mathsf{b}}(p)$.

$\mathcal{X}$ $\quad 1-p \quad$ $\mathcal{Y}$

$0 \longrightarrow 0$

$p$

$p$

$1 \longrightarrow 1$

$1-p$

To compute BSC capacity, observe $\mathrm{I}(X;Y) = \mathrm{H}(Y) - \mathrm{H}(Y|X)$, and

- $\mathrm{H}(Y|X=0) = \mathrm{H}(Y|X=1) = \mathsf{H}_{\mathsf{b}}(p) \implies \mathrm{H}(Y|X) = \mathsf{H}_{\mathsf{b}}(p)$.
- $\mathrm{H}(Y) \leq \log 2 = 1$, with equality iff $Y$ is uniform.

**Question**: Is it possible to choose a $\mathsf{P}_X$ such that $Y$ is uniform?

**Ans**: Yes – choose $X$ to be uniform $\implies \mathrm{C} = \max_{\mathsf{P}_X} \mathrm{I}(X;Y) = 1 - \mathsf{H}_{\mathsf{b}}(p)$.

# Binary erasure channel

A binary erasure channel (BEC) consists of

- Binary input $\mathcal{X} = \{0, 1\}$ and output with erasure $\mathcal{Y} = \{0, 1, *\}$.

- Channel law $\mathsf{P}_{Y|X} = \begin{bmatrix} 1-p & p & 0 \\ 0 & p & 1-p \end{bmatrix}$.

The capacity of BEC is $\mathrm{C}_{\mathrm{BEC}} = 1 - p$.



Suppose we begin with $\mathrm{I}(X;Y) = \mathrm{H}(Y) - \mathrm{H}(Y|X)$. Then,

- $\mathrm{H}(Y|X=0) = \mathrm{H}(Y|X=1) = \mathsf{H}_b(p) \implies \mathrm{H}(Y|X) = \mathsf{H}_b(p)$.
- $\mathrm{H}(Y) \leq \log 3$, with equality iff $Y$ is uniform.

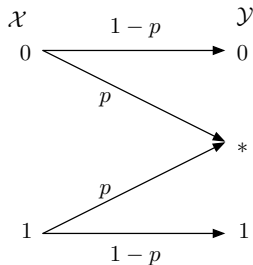**Question**: Is it possible to choose a $\mathsf{P}_X$ such that $Y$ is uniform?

**Ans**: No. So, we cannot say that $\max_{\mathsf{P}_X} \mathrm{H}(Y) = \log 3$.

Instead, we can start with $I(X;Y) = H(X) - H(X|Y)$. Then, we have the *reverse* channel law $P_{X|Y} = \begin{bmatrix} 1 & 0 \\ \alpha & 1-\alpha \\ 0 & 1 \end{bmatrix}$, where $\alpha := \Pr\{X = 0\} \equiv P_X(0)$.

- $H(X|Y=0) = H(X|Y=1) = 0$, $H(X|Y=*) = H_b(\alpha) = H(X)$
  $\implies H(X|Y) = \Pr\{Y = *\} H(X) = p H(X)$.

- $H(X) \leq 1$, with equality iff $X$ is uniform.

Hence, $C_{\text{BEC}} = \max_{P_X} (1-p) H(X) = 1 - p$.

# Erasure channel

We can generalize BEC to the following erasure channel:

- Input $\mathcal{X}$, output $\mathcal{Y} = \mathcal{X} \cup \{*\}$.

- Channel law $P_{Y|X} = \begin{cases} 1-p, & y = x \\ p, & y = * \\ 0, & \text{otherwise} \end{cases}$

A motivation for this model is from networking, where the erasure $*$ models the "packet drop".

---

### Exercise 5

Show that the capacity of the erasure channel is

$$C_{\text{EC}} = (1-p) \log |\mathcal{X}|.$$

---

# Symmetric channel

In computing the capacity of BSC, we observe that

1. $\mathrm{H}(Y|X) = \mathsf{H}_\mathsf{b}(p)$ regardless of $\mathsf{P}_X$. Why? Because all rows of $\mathsf{P}_{Y|X}$ are permutations of a same probability vector $(p, 1-p)$.

2. $\mathrm{H}(Y) = \log|\mathcal{Y}|$ can be attained, that is, $Y$ can be made uniform by choosing $X$ to be uniform. Why? Because all columns of $\mathsf{P}_{Y|X}$, have the same sum $\sum_x \mathsf{P}_{Y|X}(y|x)$.

## Definition 4 (Symmetric Channel)

A symmetric channel is a channel with channel law $\mathsf{P}_{Y|X}$ satisfying (1) all rows of $\mathsf{P}_{Y|X}$ are permutations of a same probability vector $\boldsymbol{p}$, and (2) all columns of $\mathsf{P}_{Y|X}$, have the same sum $\sum_x \mathsf{P}_{Y|X}(y|x)$.

## Exercise 6

Show that the capacity of a symmetric channel is $\log|\mathcal{Y}| - \mathrm{H}(\boldsymbol{p})$.

# Computing DMC capacity by convex optimization

For a DMC $(\mathcal{X}, \mathrm{P}_{Y|X}, \mathcal{Y})$, since $\mathrm{I}(X;Y)$ is a concave function of $\mathrm{P}_X$ for fixed $\mathrm{P}_{Y|X}$, to find the channel capacity of this DMC, we only need to solve a convex program in the capacity formula (1).

There are hence several avenues towards computing the channel capacity.

1. Try to solve it analytically, either by inspection (using structures of the channel law $\mathrm{P}_{Y|X}$), or write down the KKT conditions to find the optimal solution.

2. Call for fast algorithms to solve convex programs over a probability simplex (*conditional gradient descent* for example).

3. Leverage the problem structure and devise specific iterative algorithms.

In the following, let us introduce a widely used algorithm (Blahut-Arimoto algorithm) due to Blahut and Arimoto, which belongs to the third category.

# Blahut[1]-Arimoto[2] algorithm

**Proposition 5 (A variational characterization of mutual information)**

*Let $(X, Y) \sim \mathsf{P}_{X,Y} = \mathsf{P}_X \mathsf{P}_{Y|X} = \mathsf{P}_Y \mathsf{P}_{X|Y}$. Then,*

$$\mathrm{I}(X; Y) = \max_{\mathsf{Q}_{X|Y}} \left\{ \mathsf{E}_{\mathsf{P}_{X,Y}} \left[ \log \frac{\mathsf{Q}_{X|Y}(X|Y)}{\mathsf{P}_X(X)} \right] \right\},$$

*with $\mathsf{P}_{X|Y}$ being the maximizer.*

As a result, for a DMC $(\mathcal{X}, \mathsf{P}_{Y|X}, \mathcal{Y})$, its capacity can be written as a double maximization as follows:

$$\mathrm{C} = \max_{\mathsf{P}} \max_{\mathsf{Q}} \left\{ \sum_x \sum_y \mathsf{P}(x) \mathsf{P}_{Y|X}(y|x) \log \frac{\mathsf{Q}(x|y)}{\mathsf{P}(x)} \right\}. \tag{2}$$

[1] R. E. Blahut, "Computation of channel capacity and rate distortion functions," IEEE Transactions on Information Theory, vol. 18, no. 4, 1972.
[2] S. Arimoto, "An algorithm for calculating the capacity of arbitrary discrete memoryless channels," IEEE Transactions on Information Theory, vol. 18, no. 1, 1972.

**pf**: Goal: show $I(X;Y) \geq E_{P_{X,Y}}\left[\log \frac{Q_{X|Y}(X|Y)}{P_X(X)}\right]$ with equality if $Q_{X|Y} = P_{X|Y}$.

Note that $I(X;Y) = E_{P_{X,Y}}\left[\log \frac{P_{X|Y}(X|Y)}{P_X(X)}\right]$, and hence

$$\text{RHS} - \text{LHS} = E_{P_{X,Y}}\left[\log \frac{Q_{X|Y}(X|Y)}{P_{X|Y}(X|Y)}\right] \leq \log\left(E_{P_{X,Y}}\left[\frac{Q_{X|Y}(X|Y)}{P_{X|Y}(X|Y)}\right]\right),$$

where the last inequality is due to concavity of $\log(\cdot)$ and Jensen's inequality. Hence, it holds with equality if $Q_{X|Y}(x|y) = P_{X|Y}(x|y)$ for all $x, y$.

Finally, note that

$$\begin{aligned}
E_{P_{X,Y}}\left[\frac{Q_{X|Y}(X|Y)}{P_{X|Y}(X|Y)}\right] &= \sum_y P_Y(y) \sum_x P_{X|Y}(x|y)\frac{Q_{X|Y}(x|y)}{P_{X|Y}(x|y)} \\
&= \sum_y P_Y(y) \sum_x Q_{X|Y}(x|y) \leq 1,
\end{aligned}$$

with equality if $\text{supp}_{Q_{X|Y}(\cdot|y)} = \text{supp}_{P_{X|Y}(\cdot|y)}$ for all $y \in \text{supp}_{P_Y}$.

The proof is complete. $\qquad\square$

(2) says that the capacity of a DMC $(\mathcal{X}, \mathsf{P}_{Y|X}, \mathcal{Y})$ can be expressed as a double maximum of a *bounded* concave function $h(\mathsf{P}, \mathsf{Q})$, where

$$h(\mathsf{P}, \mathsf{Q}) := \sum_x \sum_y \mathsf{P}(x) \mathsf{P}_{Y|X}(y|x) \log \frac{\mathsf{Q}(x|y)}{\mathsf{P}(x)},$$

jointly concave in $\mathsf{P}$ and $\mathsf{Q}$. Below is a natural alternating algorithm which gives a strictly increasing sequence of $h$'s[3].

0 Initialize with some $\mathsf{P}^{(0)} \in \mathcal{P}_+(\mathcal{X})$. Set $l = 1$.

1 $\mathsf{Q}^{(l)} \leftarrow \arg\max_{\mathsf{Q}} h(\mathsf{P}^{(l-1)}, \mathsf{Q})$. By Proposition 5,

$$\mathsf{Q}^{(l)}(x|y) = \left( \mathsf{P}^{(l-1)}(x) \mathsf{P}_{Y|X}(y|x) \right) \Big/ \left( \sum_{x \in \mathcal{X}} \mathsf{P}^{(l-1)}(x) \mathsf{P}_{Y|X}(y|x) \right).$$

2 $\mathsf{P}^{(l)} \leftarrow \arg\max_{\mathsf{P} \in \mathcal{P}_+(\mathcal{X})} h(\mathsf{P}, \mathsf{Q}^{(l)})$. By solving the KKT conditions, we get

$$\mathsf{P}^{(l)}(x) = \left( \prod_y \mathsf{Q}^{(l)}(x|y)^{\mathsf{P}_{Y|X}(y|x)} \right) \Big/ \left( \sum_{x \in \mathcal{X}} \prod_y \mathsf{Q}^{(l)}(x|y)^{\mathsf{P}_{Y|X}(y|x)} \right).$$

3 $l \leftarrow l + 1$ and return to Step 1.

---

[3]Check Ch. 9 of R. Yeung, Information Theory and Network Coding, Springer, 2008, for convergence guarantee.

# Fano's inequality

**Lemma 1 (Fano's Inequality)**

$\mathrm{H}(U|V) \leq \mathsf{H_b}(\mathsf{P_e}) + \mathsf{P_e} \log|\mathcal{U}|$, *where* $\mathsf{P_e} := \Pr\{U \neq V\}$.

**pf**: Let $E := \mathbb{1}\{U \neq V\}$, the indicator function of $\{U \neq V\}$. $E \sim \mathrm{Ber}(\mathsf{P_e})$.

Using chain rule and the non-negativity of conditional entropy, we have

$$\mathrm{H}(U|V) \leq \mathrm{H}(U, E|V) = \mathrm{H}(E|V) + \mathrm{H}(U|V, E).$$

Note that $\mathrm{H}(E|V) \leq \mathrm{H}(E) = \mathsf{H_b}(\mathsf{P_e})$, and

$$\mathrm{H}(U|V, E) = \underbrace{\Pr\{E = 1\}}_{=\mathsf{P_e}} \underbrace{\mathrm{H}(U|V, E = 1)}_{\leq \log|\mathcal{U}|} + \Pr\{E = 0\} \underbrace{\mathrm{H}(U|V, E = 0)}_{=0, \because U=V}$$

Hence, $\mathrm{H}(U|V) \leq \mathsf{H_b}(\mathsf{P_e}) + \mathsf{P_e} \log|\mathcal{U}|$. $\qquad\square$

## Corollary 2 (Lower Bound on Error Probability)

$$\mathsf{P_e} \geq \frac{\mathrm{H}(U|V) - 1}{\log|\mathcal{U}|}.$$

**pf**: From Lemma 1 and $\mathsf{H_b}(\mathsf{P_e}) \leq 1$, we have

$$\mathrm{H}(U|V) \leq \mathsf{H_b}(\mathsf{P_e}) + \mathsf{P_e} \log|\mathcal{U}| \leq 1 + \mathsf{P_e} \log|\mathcal{U}|. \qquad \square$$

### Exercise 7

Show that Lemma 1 can be sharpened as follows

$$\mathrm{H}(U|V) \leq \mathsf{H_b}(\mathsf{P_e}) + \mathsf{P_e} \log\left(|\mathcal{U}| - 1\right),$$

if $U, V$ both take values in $\mathcal{U}$.

Fano's inequality is widely used in proving converse parts (impossibility results) of coding theorems in information theory.

# Proof of the (weak) converse

We would like to show that for every sequence of $(n, \lceil nR \rceil)$ codes such that $P_e^{(n)} \to 0$ as $n \to \infty$, the rate $R \leq C^I := \max_{P_X} I(X; Y)$.

**pf**: Let $k = \lceil nR \rceil$. Note that $W \sim \text{Unif}\{1, 2, ..., 2^k\}$ and hence $k = H(W)$.

$$nR \leq \lceil nR \rceil = H(W) = I\left(W; \hat{W}\right) + H\left(W \middle| \hat{W}\right) \quad (3)$$

$$\leq I(W; Y^n) + \left(1 + P_e^{(n)} \log\left(2^k\right)\right) \quad (4)$$

- (3) is due to $\lceil nR \rceil \geq nR$ and the chain rule.
- (4) is due to $W - Y^n - \hat{W}$ and Fano's inequality ($|\mathcal{W}| = 2^k$).

Expanding $I(W; Y^n)$ in (4) by the chain rule and using the fact that $\lceil nR \rceil \leq nR + 1$, we reach the following upper bound on $R$:

$$R \leq \tfrac{1}{n} \sum_{i=1}^{n} I\left(W; Y_i \middle| Y^{i-1}\right) + \tfrac{1}{n} \left(1 + P_e^{(n)} (nR + 1)\right)$$

Set $\varepsilon_n := \frac{1}{n} \left( 1 + \mathsf{P}_{\mathsf{e}}^{(n)} \left( n\mathrm{R} + 1 \right) \right)$. $\varepsilon_n \to 0$ as $n \to \infty$ because $\lim_{n\to\infty} \mathsf{P}_{\mathsf{e}}^{(n)} = 0$.

Next: relate $\sum_{i=1}^n \mathrm{I}\left(W; Y_i \middle| Y^{i-1}\right)$ to $\mathrm{I}(X; Y)$:

$$\mathrm{I}\left(W; Y_i \middle| Y^{i-1}\right) \leq \mathrm{I}\left(W, Y^{i-1}; Y_i\right) \leq \mathrm{I}\left(W, Y^{i-1}, X_i; Y_i\right) \overset{(*)}{=} \mathrm{I}(X_i; Y_i).$$

Here $(*)$ is due to the property of a DMC:

$$\mathsf{P}(y_i | x^i, y^{i-1}, w) = \mathsf{P}(y_i | x^i, y^{i-1}) = \mathsf{P}(y_i | x_i)$$
$$\Longrightarrow Y_i - X_i - \left(W, X^{i-1}, Y^{i-1}\right)$$
$$\Longrightarrow Y_i - X_i - \left(W, Y^{i-1}\right).$$

As a result, $\forall\, i$, $\mathrm{I}\left(W; Y_i \middle| Y^{i-1}\right) \leq \mathrm{I}(X_i; Y_i) \leq \max_{\mathsf{P}_X} \mathrm{I}(X; Y) =: \mathrm{C}^{\mathrm{I}}$.

Hence, we have $\forall\, n \in \mathbb{N}$,

$$\mathrm{R} \leq \frac{1}{n} \sum_{i=1}^n \mathrm{I}\left(W; Y_i \middle| Y^{i-1}\right) + \varepsilon_n \leq \mathrm{C}^{\mathrm{I}} + \varepsilon_n.$$

Taking $n \to \infty$, we conclude that $\mathrm{R} \leq \mathrm{C}^{\mathrm{I}}$ if it is achievable. $\qquad\square$

# Channel coding with feedback: problem setup



1. An $(n, k)$ channel code consists of
   - an encoding function (encoder) $\mathrm{enc}_n : \mathcal{W} \times \mathcal{Y}^{n-1} \to \mathcal{X}^n$ that maps each message $w$ to a length $n$ codeword $x^n$, and the $i$-th symbol $x_i$ is a function of $(w, y^{i-1})$ for all $i \in \{1, ..., n\}$.
   - a decoding function (decoder) $\mathrm{dec}_n : \mathcal{Y}^n \to \mathcal{W}$ that maps a channel output sequence $y^n$ to a reconstructed message $\hat{w}$.

2. The error probability is defined as $\mathsf{P}_e^{(n)} := \Pr\left\{ W \neq \hat{W} \right\}$.

3. A rate $\mathrm{R}$ is said to be *achievable* if there exist a sequence of $(n, \lceil n\mathrm{R} \rceil)$ codes such that $\mathsf{P}_e^{(n)} \to 0$ as $n \to \infty$.

   The Shannon capacity is defined as $\mathrm{C} := \sup\{\mathrm{R} \,|\, \mathrm{R} : \text{achievable}\}$.

# Feedback capacity

## Theorem 5 (Channel Coding Theorem for DMC with Feedback)

*The capacity of the DMC $\left(\mathcal{X}, \mathrm{P}_{Y|X}, \mathcal{Y}\right)$ with feedback is given by* (1)*, the same as that without feedback. In other words, feedback does not increase the channel capacity for DMC.*

Proof of the weaker version (feedback does not improve Shannon capacity) is immediate because in the converse proof of channel coding theorem without feedback, we do not make use of the assumption that there is no feedback. In other words, the proof is identical even with feedback.

Strong converse remains to hold (also due to Wolfowitz), that is, feedback does not improve $\epsilon$-capacity either.

# Overview

To prove the achievability of Theorem 4, need to show the statement below:

$$\forall \mathrm{R} < \mathrm{C}^{\mathrm{I}}, \mathrm{R} \geq 0, \ \exists \text{ a sequence of } (n, \lceil n\mathrm{R} \rceil) \text{codes such that } \lim_{n \to \infty} \mathsf{P}_{\mathsf{e}}^{(n)} = 0.$$

**Reflection**: there are two possible ways to prove the **existence** of objects satisfying some desirable properties:

1. Explicitly construct an object and show that the properties hold.
   (Feinstein's achievability proof)

2. Assume that no objects satisfy the properties, and show contradiction.
   (Shannon's achievability proof)

The achievability proof of the channel coding theorem presented in this lecture (and in most information theory textbooks) is of the second kind, and in fact belongs to the so-called **probabilistic method**.

# Coding over a noisy channel

Let us set up a few notations related to coding over a noisy channel.

1. Codebook $\mathcal{C} = \left\{ x^n(1), x^n(2), ..., x^n(2^k) \right\}$ consists of the $2^k$ *codewords* and is the range of the encoding function $\mathrm{enc}_n$, that is, $\mathcal{C} = \mathrm{enc}_n(\mathcal{W})$.

2. Error probability of message $w$: $\lambda_w := \Pr\{\hat{W} \neq W \mid W = w\}$.

   The average error probability $\mathsf{P}_e^{(n)} = \frac{1}{2^k} \sum_{w=1}^{2^k} \lambda_w$ depends on the decoding algorithm, the channel law, and the codebook $\mathcal{C}$.

3. Maximum likelihood decoder is the optimal decoder minimizing $\mathsf{P}_e^{(n)}$ when the *random* message follows uniform prior:
   $$\hat{w}_{\mathrm{ML}}(y^n) = \arg\max_{w \in \{1, ..., 2^k\}} \mathsf{P}(y^n | x^n(w)).$$

In principle, one can derive the ML decoding rule and compute $\mathsf{P}_{e,\mathrm{ML}}^{(n)}(\mathcal{C})$ for a given codebook $\mathcal{C}$.

However, there are some challenges toward proving the channel coding theorem with this approach.

# Challenges and work-arounds

First, for a given codebook $\mathcal{C}$, the expression of $\mathrm{P}_{\mathrm{e,ML}}^{(n)}(\mathcal{C})$ is intractable in general, and it is hard to obtain insights regarding the asymptotic behaviors (because symbols of $X^n \equiv x^n(W) \sim \mathrm{Unif}(\mathcal{C})$ is correlated over time).

Second, it is unclear how to construct good codebooks with guarantee of vanishing error probability.

In summary, to prove the achievability part of the channel coding theorem, there are two main challenges we shall overcome:

**1 How to show the existence of good codebooks?**

   Circumvent the issue of explicit construction with a random coding argument (a kind of the probabilistic method)

**2 How to analyze the error probability?**

   Circumvent the issue of ML decoding error analysis with a suboptimal decoder and derive upper bounds on its error probability.

# Proof program

**1 Random Codebook Generation**:

Generate an ensemble of codebooks according to certain probability distribution. Hence, the codebook becomes a random object.

Represent the codebook $\mathcal{C} = \left\{x^n(1), x^n(2), ..., x^n(2^k)\right\}$ as a $2^k \times n$ matrix $\mathbf{c}$, where each row corresponds to a codeword:

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

**2 Error Probability Analysis**:

Goal: Show that $\lim_{n\to\infty} \mathsf{E}_{\mathbf{C}}\left[\mathsf{P}_{e,\mathrm{ML}}^{(n)}(\mathbf{C})\right] = 0$ and conclude that there exists a codebook $\mathbf{c}$ such that error probability $\mathsf{P}_{e,\mathrm{ML}}^{(n)}(\mathbf{c}) \to 0$ as $n \to \infty$.

Moreover, to simplify analysis, we shall introduce suboptimal decoders and give a tractable upper bound of error probability using union bound.

# Random codebook generation

A simple way is to i.i.d. generate $2^k$ codewords, and hence the distribution of each codeword $P(x^n) = \prod_{i=1}^{n} P_X(x_i)$. In other words, the elements of the random codebook matrix $\mathbf{C}$ will be i.i.d. distributed according to $P_X$ ($P_X$ is later chosen to be the maximizer in Shannon's capacity formula (1)):

$$\mathbf{C} = \begin{bmatrix} X_1(1) & X_2(1) & \cdots & X_n(1) \\ \vdots & \vdots & \ddots & \vdots \\ X_1(2^k) & X_2(2^k) & \cdots & X_n(2^k) \end{bmatrix},$$

and hence

$$P(\mathbf{c}) \equiv \Pr\{\mathbf{C} = \mathbf{c}\} = \prod_{w=1}^{2^k} \prod_{i=1}^{n} P_X(x_i(w)). \quad \text{(we may write } \mathbf{C} \sim P_X^{\otimes 2^k n} \text{ for short)}$$

It turns out the symmetry in such codebook ensemble distribution helps simplify the analysis greatly, as we will see later.

# Encoding and decoding

For a realization $\mathbf{c}$ of the codebook ensemble $\mathbf{C} \sim \mathsf{P}_X^{\otimes 2^k n}$, we describe the encoding and decoding below:

**Encoding**: for a message $w \in \{1, ..., 2^k\}$, choose the $w$-th row of the codebook matrix $\mathbf{c}$ and send it out.

**Decoding**: ideally one would like to use the following ML decoding rule:

$$\hat{w}_{\mathrm{ML}}(y^n) = \underset{w \in \{1, ..., 2^k\}}{\arg \max} \ \mathsf{P}(y^n | x^n(w))$$

However, the performance of ML decoder is usually not easy to track, as mentioned before. Instead, we introduce a suboptimal decoder based on robust typical sequences as follows:

$$\boxed{\hat{w}_{\mathrm{T}} = \text{a unique } w \text{ such that } (x^n(w), y^n) \in \mathcal{T}_\varepsilon^{(n)}(X, Y).}$$

Here $(X, Y) \sim \mathsf{P}_X \mathsf{P}_{Y|X}$.

# Robust typical sequence

A robust typical sequence is a sequence whose empirical PMF (type) is close to the true per-letter distribution.

For a sequence $x^n$, its type is given by the frequency of occurrence of a symbol in $x^n$: $\hat{\mathsf{P}}_{x^n}(a) := \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\{x_i = a\}, \ \forall a \in \mathcal{X}$.

Due to LLN, $\hat{\mathsf{P}}_{X^n}(a) \xrightarrow{\mathrm{p}} \mathsf{P}_X(a)$ for all $a \in \mathcal{X}$ as $n \to \infty$, if $X_i \overset{\text{i.i.d.}}{\sim} \mathsf{P}_X, \ i = 1, ..., n$.

---

**Definition 5 (Robust Typical Sequence)**

For $\varepsilon \in (0, 1)$, a sequence $x^n$ is called $\varepsilon$-robust typical with respect to random variable $X \sim \mathsf{P}_X$ if $\left| \hat{\mathsf{P}}_{x^n}(a) - \mathsf{P}_X(a) \right| \le \varepsilon \mathsf{P}_X(a), \ \forall a \in \mathcal{X}$.

The $\varepsilon$-typical set $\boxed{\mathcal{T}_\varepsilon^{(n)}(X) := \{x^n \in \mathcal{X}^n \,|\, x^n \text{ is } \varepsilon\text{-typical w.r.t. rv } X\}}$.

---

Equivalent notation: $\mathcal{T}_\varepsilon^{(n)}(\mathsf{P}_X)$.

**Note**: Robust typical sequences also satisfy AEP for memoryless sequences.

**Note**: When the context is clear, we may write "$\mathcal{T}_{\varepsilon}^{(n)}$" instead of "$\mathcal{T}_{\varepsilon}^{(n)}(X)$".

### Example 4

Consider a random bit sequence generated i.i.d. based on $\mathrm{Ber}\left(\frac{1}{2}\right)$. Let us set $\varepsilon = 0.2$ and $n = 10$. What is $\mathcal{T}_{\varepsilon}^{(n)}$? How large is the typical set?

**sol**: Based on the definition, a $n$-sequence $x^n$ is $\varepsilon$-typical iff

$$\hat{\mathsf{P}}_{x^n}(0) \in [0.4, 0.6] \text{ and } \hat{\mathsf{P}}_{x^n}(1) \in [0.4, 0.6].$$

In other words, the # of "0"s in the sequence should be 4, 5, or 6. Hence, $\mathcal{T}_{\varepsilon}^{(n)}$ consists of all length-10 sequences with 4, 5, or 6 "0"s.

The size of $\mathcal{T}_{\varepsilon}^{(n)}$ is $\binom{10}{4} + \binom{10}{5} + \binom{10}{6} = 714$.

# Error probability analysis

Since the ML decoder is optimal, we can analyze the performance of the typicality decoder and use it as an upper bound.

Hence, our goal is turned to showing that if $R < C^I$, $\lim\limits_{n \to \infty} E_C\left[P_{e,T}^{(n)}(C)\right] = 0$.

**1** Leverage symmetry in the codebook ensemble to simplify $E_C\left[P_{e,T}^{(n)}(C)\right]$ and argue that we can focus on the error probability of the first random codeword $X^n(1)$ (its randomness comes from the random $C$):

$$E_C\left[P_{e,T}^{(n)}(C)\right] = E_C\left[2^{-k}\sum_{w=1}^{2^k}\lambda_w(C)\right] = 2^{-k}\sum_w E_C\left[\lambda_w(C)\right]$$

$$= 2^{-k}\sum_w E_C\left[\lambda_1(C)\right] = E_C\left[\lambda_1(C)\right]$$

$$= \Pr\{\hat{w}_T(Y^n) \neq 1 \mid W = 1\}$$

In the last line, the probability is induced by the random $X^n(1)$ and $Y^n$.

**2** For notational simplicity, use $\mathcal{E}$ to denote the error event $\{\hat{w}_{\mathrm{T}}(Y^n) \neq 1\}$. Our next focus is to upper bound $\Pr\{\mathcal{E} \mid W = 1\} := \mathsf{P}_1(\mathcal{E})$.

The trick here is to distinguish $\mathcal{E}$ into two different kinds of errors:

$$\mathcal{E} = \mathcal{E}_{\mathrm{a}} \cup \mathcal{E}_{\mathrm{t}}, \text{ where } \begin{cases} \mathcal{E}_{\mathrm{a}} := \left\{ (X^n(1), Y^n) \notin \mathcal{T}_\varepsilon^{(n)}(\mathsf{P}_{X,Y}) \right\} \\ \mathcal{E}_{\mathrm{t}} := \left\{ (X^n(w), Y^n) \in \mathcal{T}_\varepsilon^{(n)} \text{ for some } w \neq 1 \right\} \end{cases}$$

The core is whether or not the pair of sequences $(X^n(w), Y^n)$ are $\varepsilon$-typical with respect to $\mathsf{P}_{X,Y} = \mathsf{P}_X \mathsf{P}_{Y|X}$.

Let us define $\mathcal{A}_w := \left\{ (X^n(w), Y^n) \in \mathcal{T}_\varepsilon^{(n)} \right\}$.

We can then rewrite $\mathcal{E}_{\mathrm{a}} = \mathcal{A}_1^c$, $\mathcal{E}_{\mathrm{t}} = \cup_{w \neq 1} \mathcal{A}_w$, and hence

$$\mathcal{E} = \mathcal{E}_{\mathrm{a}} \cup \mathcal{E}_{\mathrm{t}} = \mathcal{A}_1^c \cup (\cup_{w \neq 1} \mathcal{A}_w).$$

**3** We are now ready to apply the union bound:

$$\mathsf{P}_1\{\mathcal{E}\} = \mathsf{P}_1\{\mathcal{A}_1^c \cup (\cup_{w \neq 1} \mathcal{A}_w)\} \leq \mathsf{P}_1\{\mathcal{A}_1^c\} + \sum_{w=2}^{2^k} \mathsf{P}_1\{\mathcal{A}_w\}.$$

Next, we shall develop upper bounds on
- the probability that the actually sent codeword $X^n(1)$ and the received $Y^n$ are not typical.
- the probability that some other (random) codeword $X^n (\neq 1)$ and the received $Y^n$ are typical.

---

### Lemma 2 (A Key Lemma)

$\mathsf{P}_1\{\mathcal{A}_1\} \geq 1 - \varepsilon$ for $n$ *large enough, and*

$$\mathsf{P}_1\{\mathcal{A}_w\} \leq 2^{-n(\mathrm{I}(X;Y) - \delta(\varepsilon))}$$

*for all $w \neq 1$, where $\delta(\varepsilon) \to 0$ as $\varepsilon \to 0$.*

4 Finally, let us put all the above together and apply Lemma 2:

$$\mathsf{E_C}\big[\mathsf{P}_{\mathsf{e,T}}^{(n)}(\mathbf{C})\big] = \Pr\{\mathcal{E} \mid W = 1\} := \mathsf{P}_1\{\mathcal{E}\} \leq \mathsf{P}_1\{\mathcal{A}_1^c\} + \sum_{w=2}^{2^k} \mathsf{P}_1\{\mathcal{A}_w\}$$

$$\leq \varepsilon + \sum_{w=2}^{2^k} 2^{-n(\mathrm{I}(X;Y) - \delta(\varepsilon))} \leq \varepsilon + 2^{-n\big(\mathrm{I}(X;Y) - \delta(\varepsilon) - \mathrm{R} - \frac{1}{n}\big)}.$$

As long as $\mathrm{R} < \mathrm{I}(X;Y) - \delta(\varepsilon)$, we are able to make $\Pr\{\mathcal{E}\} \leq 2\varepsilon$ for $n$ large enough, which implies $\lim_{n \to \infty} \mathsf{E_C}\big[\mathsf{P}_{\mathsf{e,T}}^{(n)}(\mathbf{C})\big] = 0$.

We have shown that as long as $\mathrm{R} < \mathrm{I}(X;Y) - \delta(\varepsilon)$,

$$\lim_{n \to \infty} \mathsf{E_C}\big[\mathsf{P}_{\mathsf{e,T}}^{(n)}(\mathbf{C})\big] = 0,$$

and hence there exists a realization of codebook $\mathbf{c}$ such that $\lim_{n \to \infty} \mathsf{P}_{\mathsf{e,T}}^{(n)}(\mathbf{c}) = 0$.

Finally, taking the codebook generating distribution $\mathsf{P}_X = \arg\max_{\mathsf{P}_X} \mathrm{I}(X;Y)$, we conclude that $\forall\, \mathrm{R} < \mathrm{C}^{\mathrm{I}} = \max_{\mathsf{P}_X} \mathrm{I}(X;Y)$, $\mathrm{R}$ is achievable by choosing an appropriate $\varepsilon$. $\qquad\square$

# Proof of Lemma 2 (1): recap of typicality

Recall: by definition, an $\varepsilon$-typical sequence $(x^n, y^n)$ shall satisfy

$$\left|\hat{\mathsf{P}}_{x^n,y^n}(a,b) - \mathsf{P}_{X,Y}(a,b)\right| \le \varepsilon\,\mathsf{P}_{X,Y}(a,b), \ \forall\,(a,b) \in \mathcal{X} \times \mathcal{Y}.$$

(Note: we can think of $(X,Y)$ as a r.v. and apply the same definition of typicality!)

Hence, if $(X^n, Y^n) \sim \mathsf{P}_{X,Y}^{\otimes n}$, then we have

0   $(x^n, y^n) \in \mathcal{T}_\varepsilon^{(n)}(\mathsf{P}_{X,Y}) \implies x^n \in \mathcal{T}_\varepsilon^{(n)}(\mathsf{P}_X), \ y^n \in \mathcal{T}_\varepsilon^{(n)}(\mathsf{P}_Y).$

1   $\forall\,(x^n, y^n) \in \mathcal{T}_\varepsilon^{(n)}, \ \left|\frac{1}{n}\log\mathsf{P}(x^n, y^n) + \mathrm{H}(X,Y)\right| \le \delta(\varepsilon), \ \delta(\varepsilon) = \varepsilon\mathrm{H}(X,Y).$

2   $\mathsf{P}(\mathcal{T}_\varepsilon^{(n)}) \ge 1 - \varepsilon$ for $n$ large enough.

3   $|\mathcal{T}_\varepsilon^{(n)}| \le 2^{n(\mathrm{H}(X,Y)+\delta(\varepsilon))}.$

4   $|\mathcal{T}_\varepsilon^{(n)}| \ge (1 - \varepsilon)2^{n(\mathrm{H}(X,Y)-\delta(\varepsilon))}$ for $n$ large enough.

# Proof of Lemma 2 (2): typical with the true codeword

Let us first consider $P_1 \{\mathcal{A}_1\} = \Pr\left\{(X^n(1), Y^n) \in \mathcal{T}_\varepsilon^{(n)} \,\middle|\, W = 1\right\}$.

We are averaging over a random codebook ensemble $\mathbf{C}$, and the random codebook is generated element-by-element i.i.d. based on $P_X$.

DMC without feedback implies $P(y^n|x^n) = \prod_{i=1}^n P_{Y|X}(y_i|x_i)$.

Hence, given $W = 1$, $(X^n(1), Y^n)$ has the following joint distribution:

$$
\begin{aligned}
P(x^n, y^n) &= P(x^n) \cdot P(y^n|x^n) = \prod_{i=1}^n P_X(x_i) \cdot \prod_{i=1}^n P_{Y|X}(y_i|x_i) \\
&= \prod_{i=1}^n P_{X,Y}(x_i, y_i)
\end{aligned}
$$

By Property 2 (LLN) in the previous slide, we see that for $n$ large enough,

$$
P_1 \{\mathcal{A}_1\} = \Pr\left\{(X^n(1), Y^n) \in \mathcal{T}_\varepsilon^{(n)} \,\middle|\, W = 1\right\} \geq 1 - \varepsilon.
$$

# Proof of Lemma 2 (3): typical with a wrong codeword

Consider $P_1\{\mathcal{A}_w\} = \Pr\left\{(X^n(w), Y^n) \in \mathcal{T}_\varepsilon^{(n)} \,\middle|\, W = 1\right\}$ for $w \neq 1$.

Note that we are averaging over a random codebook ensemble $\mathbf{C}$, and the random codebook is generated element-by-element i.i.d. based on $P_X$.

$X^n(1)$ and $X^n(w)$ have the same distribution $P_X^{\otimes n}$ but they are *independent*.

Due to DMC, $(X^n(1), Y^n) \perp\!\!\!\perp X^n(w)$. Hence, $Y^n \perp\!\!\!\perp X^n(w)$, and

$$P_1\{\mathcal{A}_w\} = \sum_{(x^n, y^n) \in \mathcal{T}_\varepsilon^{(n)}} P(x^n) \cdot P(y^n)$$

$$\leq \underbrace{2^{n(1+\varepsilon)H(X,Y)}}_{\substack{\text{cardinality upper} \\ \text{bound on typical set}}} \cdot \underbrace{2^{-n(1-\varepsilon)H(X)}}_{\substack{\text{upper bound on prob.} \\ \text{of a typical sequence}}} \cdot \underbrace{2^{-n(1-\varepsilon)H(Y)}}_{\substack{\text{upper bound on prob.} \\ \text{of a typical sequence}}}$$

$$= 2^{-n(I(X;Y) - \delta(\varepsilon))},$$

where $\delta(\varepsilon) = \varepsilon(H(X,Y) + H(X) + H(Y)) \to 0$ as $\varepsilon \to 0$. $\qquad\square$

# Some reflections

**Reflection 1**: **Mutual independence of codewords**

In the random coding argument of the proof, $2^k n$ elements of the codebook matrix $\mathbf{C}$ are generated i.i.d., and hence the $2^k$ rows $\left\{ X^n(1), ..., X^n(2^k) \right\}$ are *mutually independent*. In the proof, only pairwise independence is required: $X^n(1) \perp\!\!\!\perp X^n(w), \ \forall w \neq 1$.

**Reflection 2**: **Typicality decoder**

We use typicality decoder other than the optimal ML decoder to find tractable upper bounds on the error probability. There are other suboptimal decoders can be used. For example, the following threshold decoder can also work:

$$\boxed{\hat{w}_{\text{th}} := \text{a unique } w \text{ such that } \iota\left(x^n(w); y^n\right) > \beta}\,,$$

where $\iota\left(x^n; y^n\right) := \log \frac{\mathsf{P}(x^n, y^n)}{\mathsf{P}(x^n)\mathsf{P}(y^n)} = \sum_{k=1}^n \log \frac{\mathsf{P}_{Y|X}(y_k|x_k)}{\mathsf{P}_X(x_k)}$, called the *information density function*, and $\beta := n(\mathrm{I}(X; Y) - \varepsilon)$.

# Joint source-channel coding: problem setup



**Source** model: DMS $P_S$ with entropy $H(P_S)$.

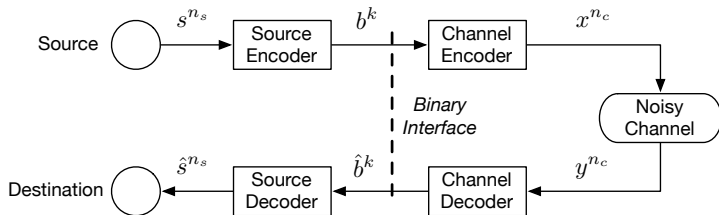**Channel** model: DMC $P_{Y|X}$ with channel capacity $C(P_{Y|X})$.

1. A $(n_s, n_c)$ joint source-channel code consists of
   - $\mathrm{enc}_{n_c} : \mathcal{S}^{n_s} \to \mathcal{X}^{n_c}$ that maps source $s^{n_s}$ to a codeword $x^{n_c}$.
   - $\mathrm{dec}_{n_c} : \mathcal{Y}^{n_c} \to \mathcal{S}^{n_s}$ that maps a channel output $y^{n_c}$ to a reconstructed $\hat{s}^{n_s}$.

2. The error probability is defined as $P_e^{(n_c)} := \Pr\{S^{n_s} \neq \hat{S}^{n_s}\}$.

3. $R$ is achievable if there exist a sequence of $(\lceil n_c R \rceil, n_c)$ codes such that $\lim_{n_c \to \infty} P_e^{(n_c)} = 0$.

# Source-channel separation theorem

## Theorem 6 (Source-Channel Separation)

*If $R < C(P_{Y|X})/H(P_S)$, then $R$ is achievable, i.e., lossless reconstruction of source $\{S_i\}$ is possible.*

*Conversely, if $R > C(P_{Y|X})/H(P_S)$, then $R$ is not achievable, i.e., lossless reconstruction is impossible.*

**pf**: (Achievability Part):

- Choose an $(n_s, \lfloor n_s R_s \rfloor)$ lossless source code with $R_s = H(P_S) + \varepsilon_s$.
- Choose an $(n_c, \lceil n_c R_c \rceil)$ channel code with $R_c = C(P_{Y|X}) - \varepsilon_c$.

By the channel coding theorem, the binary sequence $b^k$ lying at the digital interface between the source and the channel coders can be decoded with vanishing error probability.

By the lossless source coding theorem, the source sequences can be reconstructed with vanishing error probability as long as $b^k$ can be successfully decoded by the channel decoder.

Concatenate the above two codes together, we see that if

$$n_s R_s < n_c R_c \iff \frac{n_s}{n_c} < \frac{R_c}{R_s} = \frac{C - \varepsilon_c}{H(P_S) + \varepsilon_s},$$

the separation scheme is able to reconstruct the source sequence with vanishing error probability. Since $\varepsilon_s, \varepsilon_c$ can be made arbitrarily small, as long as $R < \frac{C}{H(P_S)}$, it is achievable. $\qquad\square$

(Converse Part):

We shall prove that for all achievable $\mathrm{R}$, $\mathrm{R} \leq \frac{\mathrm{C}}{\mathrm{H}(\mathsf{P}_S)}$.

$$
\begin{aligned}
&n_s \mathrm{H}(\mathsf{P}_S) \\
&= \mathrm{H}(S^{n_s}) = \mathrm{I}\left(S^{n_s}; \hat{S}^{n_s}\right) + \mathrm{H}\left(S^{n_s} \middle| \hat{S}^{n_s}\right) \hspace{3.5cm} (5) \\
&\leq \mathrm{I}(S^{n_s}; Y^{n_c}) + \left(1 + \mathsf{P}_{\mathsf{e}}^{(n_c)} n_s \log|\mathcal{S}|\right) \hspace{3.3cm} (6) \\
&\leq \sum_{k=1}^{n_c} \mathrm{I}\left(S^{n_s}; Y_k \middle| Y^{k-1}\right) + \left(1 + \mathsf{P}_{\mathsf{e}}^{(n_c)} n_s \log|\mathcal{S}|\right) \leq n_c\left(\mathrm{C} + \varepsilon_{n_c}\right), \hspace{0.5cm} (7)
\end{aligned}
$$

- (5) is due to the property of DMS and chain rule.
- (6) is due to $S^{n_s} - Y^{n_c} - \hat{S}^{n_s}$ and Fano's inequality.
- (7) is due to similar steps as in the channel coding converse proof.

Since $\varepsilon_{n_c} \to 0$ as $n_c \to \infty$, $\mathrm{R} \leq \frac{n_s}{n_c} \leq \frac{\mathrm{C}}{\mathrm{H}(\mathsf{P}_S)}$ if $\mathrm{R}$ is achievable. $\qquad \square$

# Summary: noisy channel coding and capacity

- For DMC $(\mathcal{X}, \mathsf{P}_{Y|X}, \mathcal{Y})$ with or without feedback: $\mathrm{C} = \max_{\mathsf{P}_X} \mathrm{I}(X;Y)$.

- Symmetric channel capacity $= \log|\mathcal{Y}| - \mathrm{H}(\boldsymbol{p})$, where $\boldsymbol{p}$ permutes all rows of $\mathsf{P}_{Y|X}$. Erasure channel capacity $= (1-p)\log|\mathcal{X}|$.

- Blahut-Arimoto algorithm for computing channel capacity.

- Weak converse: Fano's inequality, data processing, DMC assumption.

- Achievability: random coding argument, typicality decoder.

- Feedback does not increase the capacity of DMC.

- Joint source-channel coding theorem:

$$\mathrm{R} < \frac{\mathrm{C}}{\mathrm{H}(\mathsf{P}_S)} \implies \mathrm{R} \text{ is achievable}$$

$$\mathrm{R} > \frac{\mathrm{C}}{\mathrm{H}(\mathsf{P}_S)} \implies \mathrm{R} \text{ is not achievable}$$

- Source-channel separation is optimal.