# Lec 1: Entropy & Mutual Information

Yanjun Han

<u>Entropy</u>. For a discrete RV $X$ taking value in $\mathcal{X}$ with pmf $p$, its entropy $H(X)$ <span style="color:blue">(or $H(p)$)</span> is defined as

$$H(X) = \sum_{x \in \mathcal{X}} p(x) \log \frac{1}{p(x)} .$$

<u>Remarks</u>:  1. $0 \le H(X) \le \log |\mathcal{X}|$ $\left( H(X) \le \log \sum_{x \in \mathcal{X}} p(x) \cdot \frac{1}{p(x)} = \log |\mathcal{X}| \text{ by Jensen} \right)$

2. $H(X)$ can be <u>finite</u> or <u>infinite</u> when $|\mathcal{X}| = \infty$

3. For continuous (or general) RVs, need to find a measure $\mu$ s.t. $X$ has a density $f$ w.r.t. $\mu$, and define the <u>differential entropy</u>

$$h(X) = \int_{\mathcal{X}} f(x) \log \frac{1}{f(x)} \, d\mu(x)$$

<span style="color:blue">↑ usually lower-case h</span>   <span style="color:blue">↑ value depends on the choice of $\mu$</span>

4. Base of $\log$:  for IT applications (only this lecture), take
$$\log = \log_2 \quad (\text{bits});$$
For other applications (all later lectures), $\log = \log_e$ (nats).

Why entropy? Shannon (1948) shows that entropy characterizes the <u>fundamental limit</u> of <u>source coding</u>.

<u>Source coding problem</u> <span style="color:red">(for the i.i.d. case)</span>

Given:  <span style="color:blue">①</span> an input alphabet $\mathcal{X}$ <span style="color:blue">(e.g. all English letters $\{a, b, \cdots, z\}$)</span>
<span style="color:blue">②</span> a known pmf $p$ on $\mathcal{X}$ <span style="color:blue">(i.e. the source distribution)</span>

Target:  find a map (i.e. code) $f: \mathcal{X} \to \{0,1\}^* := \bigcup_{n=1}^{\infty} \{0,1\}^n$, such that
① it's <u>uniquely decodable</u>, i.e. based on the concatenation $(f(x_1), \cdots, f(x_m))$, one can uniquely decode $m$ and $(x_1, \cdots, x_m) \in \mathcal{X}^m$
② the expected codelength $\mathbb{E}[\ell(f(X))] = \sum_{x \in \mathcal{X}} p(x) \ell(f(x))$ is minimized
<span style="color:blue">$\ell(\cdot)$: length of the codeword (in bits)</span>

<u>Example</u>. If $X = \{a, b, c\}$ and $p = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4})$

(a) The code $a \rightarrow 0$, $b \rightarrow 10$, $c \rightarrow 11$ is uniquely decodable,

  e.g.  $100 10 11$ decodes into $babc$

(b) The code $a \rightarrow 0$, $b \rightarrow 1$, $c \rightarrow 10$ is <u>NOT</u> uniquely decodable.

  e.g.  $10$ decodes into either $c$ or $ba$

(c) The code $a \rightarrow 10$, $b \rightarrow 0$, $c \rightarrow 11$ is uniquely decodable and has a

  smaller expected codelength $2 \cdot \frac{1}{4} + 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{4} = 1.5$ bits $< 1.75$ bits for (a).

Given a length profile $\{l_x\}_{x \in X}$ is there a uniquely decodable code $f$
with $l(f(x)) = l_x$ ?

---

<u>Theorem</u> ( Kraft - McMillan)

  A necessary and sufficient condition is

$$\sum_{x \in X} 2^{-l_x} \leq 1. \quad \text{(Kraft inequality)}$$

---

<u>Pf</u>. (Sufficiency) First note that for a full

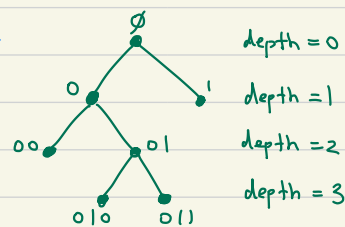  binary tree (i.e. a binary tree where each node

  has 0 or 2 children), then

$$\sum_{\substack{\text{leaf node} \\ v}} 2^{-depth(v)} = 1.$$

Because $\sum_{x \in X} 2^{-l_x} \leq 1$, one can construct a full

binary tree s.t. $X \subseteq \{$all leaf nodes$\}$

  and $depth(x) = l_x$, $\forall x \in X$.

Now use the coding scheme in the example, which results in a <u>prefix code</u>,
i.e. no codeword is a prefix of the other. Easy to show that (Exercise)
prefix codes are uniquely decodable.

depth = 0
depth = 1
depth = 2
depth = 3

A full binary tree, with
codewords $\{1, 00, 010, 011\}$

(Necessity) WLOG assume $|X| < \infty$ and $\ell_{max} := \max_{x \in X} \ell_x < \infty$.

Use a tensor power trick: for uniquely decodable code $f$.

$$\left( \sum_{x \in X} 2^{-\ell(f(x))} \right)^m = \sum_{x_1, \cdots, x_m \in X} 2^{-(\ell(f(x_1)) + \cdots + \ell(f(x_m)))}$$

$$= \sum_{x_1, \cdots, x_m \in X} 2^{-\ell(\overbrace{(f(x_1), \cdots, f(x_m))}^{\text{concatenation}})}$$

$$= \sum_{\ell=1}^{m\ell_{max}} 2^{-\ell} \{\# \text{ of concatenated codewords of total length } \ell\}$$

$$\leq \sum_{\ell=1}^{m\ell_{max}} 2^{-\ell} \cdot 2^{\ell} \quad (\text{by uniquely decodable assumption})$$

$$= m\ell_{max}$$

$$\implies \sum_{x \in X} 2^{-\ell(f(x))} \leq (m\ell_{max})^{1/m} \xrightarrow{m \to \infty} 1. \qquad \boxed{20}$$

Using Kraft inequality, we obtain the following characterization of the smallest expected codelength.

---

Thm ( Source coding theorem for uniquely decodable code )

$$H(X) \leq \min_{\substack{\text{uniquely decodable} \\ f}} \mathbb{E}[\ell(f(X))] < H(X) + 1$$

---

Pf. (Upper bound) $\ell_x = \lceil \log_2 \frac{1}{p(x)} \rceil$ satisfies Kraft inequality, and

$$\sum_{x \in X} p(x) \ell_x < \sum_{x \in X} p(x) \left( \log_2 \frac{1}{p(x)} + 1 \right) = H(x) + 1.$$

(Lower bound) Easy to show via Lagrangian multiplier that

$$\begin{cases} \min_{\ell \in \mathbb{R}_+^{|x|}} \sum_x p(x) \ell_x \\ \text{s.t.} \sum_x 2^{-\ell_x} \leq 1 \end{cases} = \sum_x p(x) \log_2 \frac{1}{p(x)} = H(X) \qquad \boxed{4}$$

**Remark:** 1. The gap between $H(X)$ and $H(X)+1$ could be significant (e.g. when $H(X) = 1.5$ bits). However, in practice, the alphabet $X$ is usually "super-symbols", e.g. $X = \{a, \cdots, z\}^{256}$ instead of $\{a, \cdots, z\}$. In such cases, $H(X) \gg 1$ bit.

2. Information theory is usually good at proving "robust" results even if a small error probability can be tolerated; in contrast, the above combinatorial argument fails to do so. See more details below.

<u>Asymptotic equipartition property (AEP)</u>.

Another way to write the entropy is

$$H(X) = \mathbb{E}_{X \sim p}\left[\log \frac{1}{p(X)}\right]$$

<span style="color:red">(* Warning: Some of you might not be used to seeing the distribution $p$ to appear in BOTH the expectation AND the function )</span>

Therefore, if $X_1, \cdots, X_n \overset{i.i.d.}{\sim} p$, LLN leads to (if $H(X) < \infty$)

$$\frac{1}{n} \log \frac{1}{p(X_1, \cdots, X_n)} = \frac{1}{n} \sum_{i=1}^{n} \log \frac{1}{p(X_i)} \overset{a.s.}{\longrightarrow} \mathbb{E}\left[\log \frac{1}{p(X)}\right] = H(X),$$
$$\text{as } n \to \infty.$$

$$\Rightarrow \forall \varepsilon > 0, \quad \mathbb{P}\left( \underbrace{p(X_1, \cdots, X_n) \in [2^{-n(H(X)+\varepsilon)}, 2^{-n(H(X)-\varepsilon)}]}\right) \longrightarrow 1 \text{ as } n \to \infty.$$

$$\text{call this set } T_n^\varepsilon \text{ (typical set)}$$

$\Rightarrow$ 

> **Thm (AEP).** The typical set $T_n^\varepsilon$ satisfies that
> ① $\mathbb{P}\left( (X_1, \cdots, X_n) \in T_n^\varepsilon \right) \longrightarrow 1$, as $n \to \infty$.
> ② $(1 - o(1)) \, 2^{n(H(X) - \varepsilon)} \leq |T_n^\varepsilon| \leq 2^{n(H(X) + \varepsilon)}$

In other words, AEP states that for $X_1, \dots, X_n \overset{i.i.d}{\sim} P$, the joint distribution of $X_1, \dots, X_n$ is "roughly" a uniform distribution over $\doteq 2^{nH(P)}$ typical sequences.

## Source coding theorem with error probability.

Diagram :  $\boxed{X_1, \dots, X_n \overset{i.i.d}{\sim} P}$ $\xrightarrow{\text{encoder}}$ $\boxed{Y \in \{0,1\}^*}$ $\xrightarrow{\text{decoder}}$ $\boxed{(\hat{X}_1, \dots, \hat{X}_n)}$

with a block error guarantee $\mathbb{P}((X_1, \dots, X_n) \neq (\hat{X}_1, \dots, \hat{X}_n)) \leq \delta$.

> **Thm.** ① Achievability : $\exists$ (encoder, decoder) s.t. $\frac{1}{n}\mathbb{E}[\ell(Y)] \leq H(P) + o(1)$
> and $\delta = o(1)$.
> ② Converse : if $\delta = o(1)$, then ANY (encoder, decoder) pair satisfies
> $\frac{1}{n}\mathbb{E}[\ell(Y)] \geq H(P) - o(1)$.

**Pf.** (Achievability) Consider an encoder-decoder pair that enumerates all typical sequences in $T_n^\epsilon$ and ignores all others. Then by AEP,

$$\text{error prob.} = \mathbb{P}((X_1, \dots, X_n) \notin T_n^\epsilon) \to 0$$
$$\ell(Y) \leq \log_2 |T_n^\epsilon| \leq n(H(P) + \epsilon) \quad \text{deterministically.}$$

Since $\epsilon > 0$ is arbitrary, the achievability follows.

(Converse) Fix any $\epsilon > 0$. Define two sets

$$A = \{(X_1, \dots, X_n) : \ell(Y) > n(H(P) - 2\epsilon)\}$$
$$B = \{(X_1, \dots, X_n) : (X_1, \dots, X_n) = (\hat{X}_1, \dots, \hat{X}_n)\}.$$

then

$$\mathbb{P}(T_n^\epsilon \cap B) \geq 1 - \delta - o(1) \quad \text{by AEP and union bound.}$$

Moreover,

$$|T_n^\varepsilon \cap B \cap A^c| = \left|\{(x_1, \cdots, x_n) \in T_n^\varepsilon \cap B : \ell(Y(x_1, \cdots, x_n)) \le n(H(P) - 2\varepsilon)\}\right|$$

$$\le \left|\{y \in \{0,1\}^* : \ell(y) \le n(H(P) - 2\varepsilon)\}\right|$$

       ↑
by defn. of $B$, if $(x_1, \cdots, x_n), (x_1', \cdots, x_n') \in B$ are different,
one must have $Y(x_1, \cdots \to x_n) \ne Y(x_1', \cdots, x_n')$

$$= \sum_{k=1}^{n(H(P) - 2\varepsilon)} 2^k < 2 \cdot 2^{n(H(P) - 2\varepsilon))}$$

$$\Longrightarrow \quad \mathbb{P}(T_n^\varepsilon \cap B \cap A^c) \le 2^{-n(H(P) - \varepsilon)} \cdot |T_n^\varepsilon \cap B \cap A^c| < 2 \cdot 2^{-n\varepsilon}.$$

           ↑
      by AEP

Therefore, $\quad \mathbb{P}(T_n^\varepsilon \cap A \cap B) \ge 1 - \delta - o(1) - 2 \cdot 2^{-n\varepsilon} = 1 - o(1)$

$$\Longrightarrow \frac{1}{n} \mathbb{E}[\ell(Y)] \ge (H(P) - 2\varepsilon) \cdot \mathbb{P}(A) \ge (1 - o(1)) \cdot (H(P) - 2\varepsilon)$$

by Markov's inequality. Since $\varepsilon > 0$ is arbitrary, the converse follows.     🔲

## Joint entropy and mutual information.

Similar to $\quad H(X) = \mathbb{E}_X\left[\log \frac{1}{P(X)}\right]$, can also define

$$H(X, Y) = \mathbb{E}_{X,Y}\left[\log \frac{1}{P(X,Y)}\right] \quad \text{(joint entropy)}$$

$$H(Y|X) = \mathbb{E}_{X,Y}\left[\log \frac{1}{P(Y|X)}\right]$$

$$\qquad\qquad = H(X, Y) - H(X) \quad \text{(conditional entropy)}$$

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

$$\qquad\qquad = H(Y) - H(Y|X)$$

$$\qquad\qquad = \mathbb{E}_{X,Y}\left[\log \frac{P(X,Y)}{P(X)P(Y)}\right] \quad \text{(mutual information)}$$

**Pf.**  There is a one-line proof using convexity / KL divergence (next lecture),
     but let's present a proof using typicality / AEP that will be useful later.

Define   $T_n^{\varepsilon}(X) = \left\{ (x^n, y^n) : \left| \frac{1}{n} \sum_{i=1}^{n} \log_2 \frac{1}{P_X(X_i)} - H(X) \right| \leq \varepsilon \right\}$

     $T_n^{\varepsilon}(Y) = \left\{ (x^n, y^n) : \left| \frac{1}{n} \sum_{i=1}^{n} \log_2 \frac{1}{P_Y(Y_i)} - H(Y) \right| \leq \varepsilon \right\}$

     $T_n^{\varepsilon}(X,Y) = \left\{ (x^n, y^n) : \left| \frac{1}{n} \sum_{i=1}^{n} \log_2 \frac{1}{P_{XY}(X,Y)} - H(X,Y) \right| \leq \varepsilon \right\}$.

and   joint typical set   $T_n^{\varepsilon} = T_n^{\varepsilon}(X) \cap T_n^{\varepsilon}(Y) \cap T_n^{\varepsilon}(X,Y)$.

For  $(X_1, Y_1), \cdots, (X_n, Y_n) \overset{i.i.d.}{\sim} P_{XY}$, LLN + union bound  yields

          $\mathbb{P}((X^n, Y^n) \in T_n^{\varepsilon}) \xrightarrow{n \to \infty} 1$,

from which  one  deduces  that    $|T_n^{\varepsilon}| \geq (1-o(1)) \, 2^{n(H(X,Y)-\varepsilon)}$.

Next  draw  $(\tilde{X}_1, \tilde{Y}_1), \cdots, (\tilde{X}_n, \tilde{Y}_n) \overset{i.i.d}{\sim} P_X P_Y$.  Then

     $1 \geq \mathbb{P}((\tilde{X}^n, \tilde{Y}^n) \in T_n^{\varepsilon})$

        $= \sum_{(x^n, y^n) \in T_n^{\varepsilon}} \mathbb{P}(\tilde{X}^n = x^n, \tilde{Y}^n = y^n)$

        $\geq (1-o(1)) \, 2^{n(H(X,Y)-\varepsilon)} \cdot 2^{-n(H(X)+\varepsilon)} \cdot 2^{-n(H(Y)+\varepsilon)}$

        $= (1-o(1)) \, 2^{-n(I(X;Y)+3\varepsilon)}$

$\Rightarrow$  $I(X;Y) + 3\varepsilon \geq 0$,  and  $I(X;Y) \geq 0$ by taking $\varepsilon \to 0^+$.

This is a fundamental inequality to prove other inequalities. e.g.

① $H(X_1, \cdots, X_n) = \sum_{k=1}^{n} H(X_k | X^{k-1}) \leq \sum_{k=1}^{n} H(X_k)$

② If $P_{Y^n | X^n} = \prod_i P_{Y_i | X_i}$, then

$$I(X^n ; Y^n) = H(Y^n) - H(Y^n | X^n)$$

$$= H(Y^n) - \sum_{i=1}^{n} H(Y_i | X_i) \quad \left( H(Y^n | X^n) = \mathbb{E}\left[ \log \frac{1}{P(Y^n | X^n)} \right] \right.$$

$$\left. = \mathbb{E}\left[ \sum_i \log \frac{1}{P(Y_i | X_i)} \right] \right)$$

$$\leq \sum_{i=1}^{n} H(Y_i) - \sum_{i=1}^{n} H(Y_i | X_i)$$

$$= \sum_{i=1}^{n} I(X_i ; Y_i)$$

③ If $P_{X^n} = \prod_i P_{X_i}$, then

$$I(X^n ; Y^n) = H(X^n) - H(X^n | Y^n)$$

$$= \sum_i H(X_i) - H(X^n | Y^n)$$

$$\geq \sum_i H(X_i) - \sum_i H(X_i | Y^n) \quad (\text{by } ①)$$

$$\geq \sum_i H(X_i) - \sum_i H(X_i | Y_i) \quad (\text{conditioning reduces}$$
$$\text{entropy})$$

$$= \sum_i I(X_i ; Y_i)$$

Remark: All inequalities that can be shown via
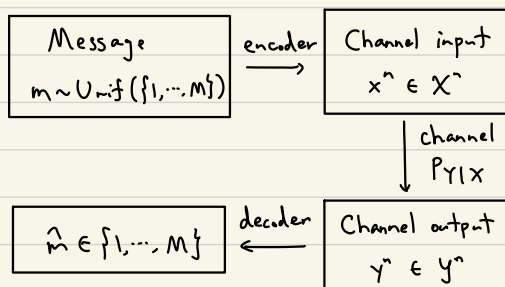    ① monotonicity :    $H(X) \leq H(X, Y)$
    ② submodularity :    $H(X_A) + H(X_B) \geq H(X_{A \cup B}) + H(X_{A \cap B})$
are called   <u>Shannon-type inequalities</u>.

Why mutual information? Shannon (1948) shows that it characterizes
the fundamental limit of <u>communications/channel coding</u> and <u>lossy compression</u>
                      (later this lecture)          (skipped; related to "mutual info method"
                                                  for statistics later)

# Channel coding problem.

Diagram :

| Message $m \sim \text{Unif}(\{1, \cdots, M\})$ | $\xrightarrow{\text{encoder}}$ | Channel input $x^n \in X^n$ |
|---|---|---|

$\downarrow$ channel $P_{Y|X}$ → known and given by nature
→ n channel uses are independent
i.e. $P_{Y^n|X^n} = \prod_i P_{Y_i|X_i}$

| $\hat{m} \in \{1, \cdots, M\}$ | $\xleftarrow{\text{decoder}}$ | Channel output $Y^n \in Y^n$ |
|---|---|---|

Goal :   Given a (block) error probability guarantee $\mathbb{P}(m \neq \hat{m}) \leq \delta$,
aim to send as many messages as possible, or equivalently,
maximize the <u>rate</u> of communication

$$R_n = \frac{\log M}{n} \quad \text{(bits per channel use)}$$

---

<u>Defn</u> (channel capacity).

$$C = C(P_{Y|X}) = \max_{P_X} I(X; Y), \quad \text{with } P_{XY} = P_X P_{Y|X}$$
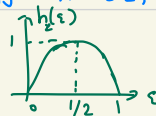
---

(In other words, given the transition probability $P_{Y|X}$ from $X$ to $Y$,
design an input distribution $P_X$ s.t. $I(X; Y)$ is maximized )

<u>Examples</u>.  ① Binary symmetric channel (BSC):

$$P_{Y|X}: \quad X \quad \begin{array}{c} 0 \\ 1 \end{array} \begin{array}{cc} \overset{0}{\phantom{x}} & \overset{1}{\phantom{x}} \\ \begin{bmatrix} 1-\varepsilon & \varepsilon \\ \varepsilon & 1-\varepsilon \end{bmatrix} \end{array}$$

$$I(X; Y) = H(Y) - H(Y|X) \leq 1 - h_2(\varepsilon), \quad \text{with equality iff } P_X = [\tfrac{1}{2}, \tfrac{1}{2}].$$

↑ binary entropy function
$h_2(\varepsilon) = \varepsilon \log_2 \frac{1}{\varepsilon} + (1-\varepsilon) \log_2 \frac{1}{1-\varepsilon}$.

② Binary erasure channel (BEC):

$$P_{Y|X}: \quad \begin{array}{c} \\ X \end{array} \begin{array}{c} \\ 0 \\ 1 \end{array} \begin{array}{c} \overset{Y}{\phantom{x}} \\ \begin{array}{ccc} 0 & 1 & \perp \end{array} \\ \begin{bmatrix} 1-\varepsilon & 0 & \varepsilon \\ 0 & 1-\varepsilon & \varepsilon \end{bmatrix} \end{array}$$

$$\begin{aligned} I(X;Y) &= H(X) - H(X|Y) \\ &= H(X) - \underbrace{P(Y \neq \perp) \, H(X|Y \neq \perp)}_{=0} - P(Y=\perp) \underbrace{H(X|Y=\perp)}_{= H(X)} \\ &= (1-\varepsilon) H(X) \leq 1-\varepsilon, \quad \text{with equality iff } P_X = [\tfrac{1}{2}, \tfrac{1}{2}]. \end{aligned}$$

---

**Thm (Shannon's channel coding theorem)** Fix any $\varepsilon > 0$.

① Achievability: if $R_n < C - \varepsilon$, then $\exists$ (encoder, decoder) s.t.
$$P(m \neq \hat{m}) \longrightarrow 0 \quad \text{as} \quad n \to \infty.$$
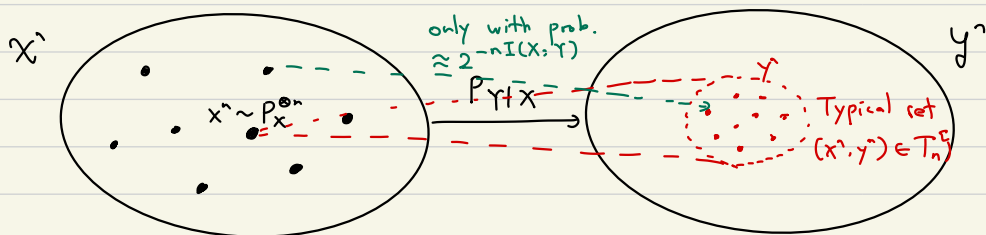
② (Weak) converse: if $R_n > C + \varepsilon$, then $\forall$ (encoder, decoder),
$$\liminf_{n \to \infty} P(m \neq \hat{m}) > 0.$$

(Strong converse: $\liminf\limits_{n \to \infty} P(m \neq \hat{m}) = 1$ ; see Lec 4)

---

In other words, the maximum rate of communication is (asymptotically) the channel capacity!

Achievability: random coding & typicality.

<u>Random codebook</u>: generate $X_{(1)}^n, \cdots, X_{(m)}^n \overset{i.i.d.}{\sim} P_X^{\otimes n}$

<u>Encoder</u>: for message $m \in [M]$, send $X_{(m)}^n$

<u>Decoder</u>: find the unique $\hat{m} \in [M]$ s.t. $(X_{(\hat{m})}^n, Y^n)$ is joint typical) (see defn. on page 8); if none or not unique, report failure.

<u>Analysis</u>: WLOG assume that the true message is $m = 1$.

Then $\hat{m} = m$ if:

① $(X_{(1)}^n, Y^n)$ is joint typical;

② <u>none</u> of $(X_{(2)}^n, Y^n), \cdots, (X_{(m)}^n, Y^n)$ is joint typical.

By LLN, $\mathbb{P}(①) = 1 - o(1)$.

Reversing the analysis on Page 8, since $(X_{(2)}^n, Y^n) \sim P_X^{\otimes n} \otimes P_Y^{\otimes n}$ (independent !!),

$$\mathbb{P}\left( (X_{(2)}^n, Y^n) \text{ joint typical} \right) \leq 2^{-n\left( I(X;Y) - 3\varepsilon \right)},$$

so union bound gives $\mathbb{P}(②) \geq 1 - M \cdot 2^{-n\left( I(X;Y) - 3\varepsilon \right)}$. If $\log_2 M < n(I(X;Y) - 4\varepsilon)$ then

$$\mathbb{P}(②) \geq 1 - e^{-n\varepsilon} = 1 - o(1).$$

Therefore, $\mathbb{P}(\hat{m} = 1) \geq \mathbb{P}(① \text{ and } ②) = 1 - o(1)$. ▣

<span style="color:red">Remark: 1. Random coding was a remarkable idea at the time, when algebraic codes were more popular. This also motivated the entire field of probabilistic methods.

2. This coding scheme is computationally expensive. First efficient codes which attains the Shannon limit were found in 2000's, including the spatially coupled LDPC code and polar code.</span>

## Weak converse : Fano's inequality.

(How IT-based ideas are robust to errors)

**Lemma.** (Data processing inequality for MI)

If $X - Y - Z$ forms a Markov chain (i.e. $P_{XYZ} = P_X P_{Y|X} P_{Z|Y}$). then

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

**Pf.** Shannon-type inequalities:

$$I(X; Y) - I(X; Z) = H(X|Z) - H(X|Y)$$
$$= H(X|Z) - H(X|Y, Z) \quad (\text{By Markov})$$
$$= I(X; Y|Z) \geq 0. \qquad \boxed{a}$$

**Thm** (Fano's inequality)   If $X \sim \text{Unif}([M])$.

Then
$$P(X \neq Y) \geq 1 - \frac{I(X; Y) + \log 2}{\log M}.$$

**Pf.** Let $E = \mathbb{1}(X \neq Y)$. Then

$$H(X|Y) = H(X|Y, E) + \underbrace{I(X; E|Y)}_{\leq H(E) \leq \log 2}$$

$$\leq P(E=1) \underbrace{H(X|Y, E=1)}_{\leq H(X) = \log M} + P(E=0) \underbrace{H(X|Y, E=0)}_{= 0} + \log 2$$

$$\leq P(X \neq Y) \cdot \log M + \log 2.$$

On the other hand.

$$H(X|Y) = H(X) - I(X; Y)$$
$$= \log M - I(X; Y), \quad (X \sim \text{Unif}[M])$$

rearranging yields the claim.                     $\boxed{a}$

(In Lec 2, we'll see more "principled" proofs of Fano's inequality)

To apply Fano's inequality, if $R_n > C + \varepsilon$,

$$\mathbb{P}(m \neq \hat{m}) \geq 1 - \frac{I(m; \hat{m}) + \log 2}{\log M}$$

$$\geq 1 - \frac{I(X^n; Y^n) + \log 2}{\log M} \quad \text{(Markov structure } m - X^n - Y^n - \hat{m})$$

$$\geq 1 - \frac{\sum_{i=1}^{n} I(X_i; Y_i) + \log 2}{\log M} \quad \text{(inequality ② on Page 9 )}$$

$$\geq 1 - \frac{nC + \log 2}{\log M} \quad \text{(defn. of } C)$$

$$\geq 1 - \frac{nC + \log 2}{n(C + \varepsilon)} \xrightarrow{n \to \infty} \frac{\varepsilon}{C + \varepsilon} > 0,$$

establishing the weak converse.

[q]