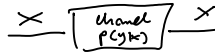


3.1 Mutual information

X, Y two random variables



how much do we learn about X once we know Y
intuitively this is how much information channel has transmitted

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

$$I(X:Y) \geq 0 \quad (=0 \text{ for independent variables})$$

Def: Conditional mutual information

$$I(X:Y|Z) = H(X|Z) - H(X|YZ)$$

Fact

(Chain rule for I :)

$$I(X_1, X_2 : Y) = I(X_1 : Y) + I(X_2 : Y | X_1)$$

Proof:

$$L = H(X_1, X_2) + H(Y) - H(X_1, X_2, Y)$$

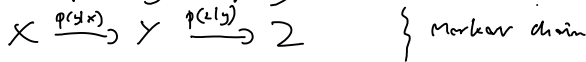
$$R = H(X_1) + H(Y) - H(X_1, Y) + H(X_2 | X_1) + H(Y | X_1) - H(X_2, Y | X_1)$$

$$\{ H(X_2, Y | X_1) = H(X_2 | Y X_1) + H(Y | X_1)$$

$$- H(X_1, X_2, Y) = -H(X_1, Y) - H(X_2 | Y X_1)$$

$$0 = 0$$

3.2 Data processing inequality



$$p(x, y, z) = p(z|y) p(y|x) p(x)$$

$$\text{(i) } I(X:Y) \geq I(X:Z) \quad \left\{ \begin{array}{l} \text{alternatively} \\ H(X|Y) \leq H(X|Z) \end{array} \right.$$

$$\text{(ii) } I(Y:Z) \geq I(X:Z) \quad \left\{ H(Z|Y) \leq H(Z|X) \right.$$

Proof:

$$\begin{aligned} \text{(i) } I(X:Y) &= I(Y:X) + I(Z:X|Y) \\ &= I(Z:X) + I(Y:X|Z) \end{aligned}$$

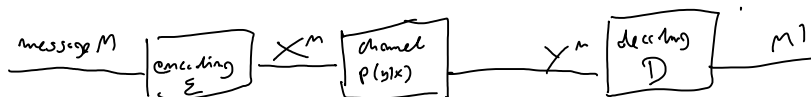
$$I(Z:X|Y) = H(Z|X) - H(Z|XY) = 0$$

↑
not relevant

$$I(Y:X) = I(Z:X) + I(Y:X|Z) \Rightarrow I(Y:X) \geq I(Z:X)$$

(ii) analogously

3.3 Shannon channel theorem (the greatest achievement of information theory)



$$\text{code} = (M, n)$$

$$\begin{aligned} \text{Max error probability } \lambda_n &= \max_{m \in M} \Pr \left\{ D(Y^n) \neq m \mid X^n = E(m) \right\} = \\ &= \sum_{y^n} p(y^n | x^n(m)) \cdot \mathbb{1} \left[D(y^n) \neq m, 1, 0 \right] \end{aligned}$$

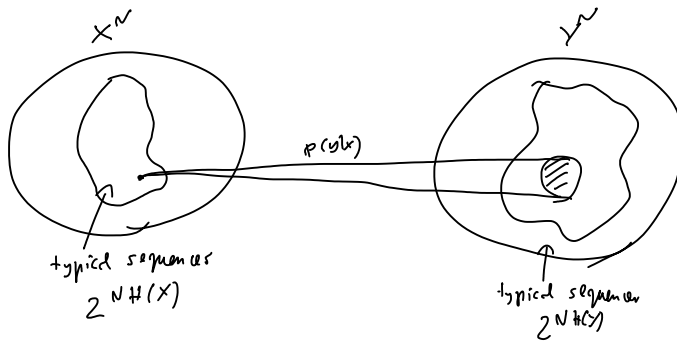
Definition Rate of the code $R = \frac{\log |M|}{n}$

Definition A rate R is achievable if there exist a sequence of $(2^{nR}, n)$ codes where maximal probability $\lambda_n \rightarrow 0$ as $n \rightarrow \infty$

of error $\lim_{n \rightarrow \infty} \rightarrow 0$.

Theorem: R is achievable iff $R \leq C = \max_{p_X} I(X; Y)$

Intuition:

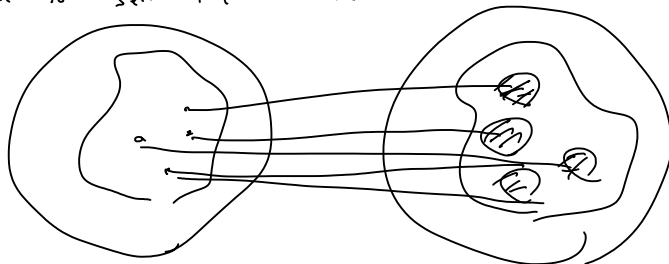


one typical sequence in X^N will be transformed to $2^{NH(Y|X)}$

$N p(x)$ way to prepare x random probability measure
 $2^{N p(x) H(Y|X)}$ unique strings in $N p(x)$

to some other knowledge x . Many wires:
 $2^{\sum p(x) H(Y|X)} = 2^{N H(Y|X)}$ typ sequences

how many different input sequences can be used to send information:



$$\frac{2^{NH(X)}}{2^{NH(Y|X)}} = 2^{NI(X;Y)}$$

Mutual information tells us about channel capacity

Strictly:

For every rate $R < \max_{p_X} I(X; Y)$ there exists a

(1) coding $(2^{nR}, n)$ with probability of error arbitrarily small for $n \rightarrow \infty$.
 number of messages \hookrightarrow number of symbols transmitted

(2) Conversely if $(2^{nR}, n)$ exists $\Rightarrow R \leq C$

Definition: Jointly typical sequences T_ϵ^n

(x^n, y^n) are jointly ϵ -typical with respect to $p(x, y)$ iff

- $|\frac{1}{n} \log p(x^n) - H(X)| \leq \epsilon$
- $|\frac{1}{n} \log p(y^n) - H(Y)| \leq \epsilon$
- $|\frac{1}{n} \log p(x^n, y^n) - H(X, Y)| \leq \epsilon$ $p(x^n, y^n) = \prod_{i=1}^n p(x_i, y_i)$

Properties \rightarrow asymptotically for large n

\rightarrow a.s. $(\text{number of } T_\epsilon^n) \approx 2^{nI(X;Y)}$

(c) $\Pr[(x^m, y^m) \in T_\epsilon^m] \geq 1 - \delta$

(c') $(1 - \delta) 2^{m(H(x, y) - \epsilon)} \leq |T_\epsilon^m| \leq 2^{m(H(x, y) + \epsilon)}$

(c'') If \tilde{x}^m, \tilde{y}^m are independent variables with the same marginals as x^m, y^m : $p(\tilde{x}^m, \tilde{y}^m) = p(x^m)p(y^m)$ then

$$(1 - \delta) 2^{-m(I(x, y) + 3\epsilon)} \leq \Pr[(\tilde{x}^m, \tilde{y}^m) \in T_\epsilon^m] \leq 2^{-m(I(x, y) - 3\epsilon)}$$

Proofs

(i) By law of large numbers for sufficiently large m we have

$$\Pr[|-\frac{1}{m} \log p(x^m) - H(x)| > \epsilon] \leq \frac{\delta}{2}$$

$$\Pr[|-\frac{1}{m} \log p(y^m) - H(y)| > \epsilon] \leq \frac{\delta}{2} \Rightarrow \Pr[(x^m, y^m) \notin T_\epsilon^m] \leq \delta$$

$$\Pr[|-\frac{1}{m} \log p(x^m y^m) - H(x, y)| > \epsilon] \leq \frac{\delta}{2}$$

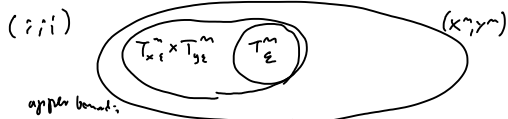
(ii) $\Pr(x^m, y^m) \geq 2^{m(H(x, y) + \epsilon)} \quad \sum_{x^m, y^m \in T_\epsilon^m} p(x^m, y^m) \leq 1$

$$|T_\epsilon^m| \cdot 2^{-m(H(x, y) + \epsilon)} \leq 1 \quad |T_\epsilon^m| \leq 2^{m(H(x, y) + \epsilon)}$$

$$\Pr(x^m, y^m) \leq 2^{-m(H(x, y) - \epsilon)}$$

$$1 = \sum_{x^m, y^m} p(x^m, y^m) = \sum_{(x^m, y^m) \in T_\epsilon^m} p(x^m, y^m) + \sum_{(x^m, y^m) \notin T_\epsilon^m} p(x^m, y^m) \leq |T_\epsilon^m| 2^{-m(H(x, y) - \epsilon)} + \delta$$

$$|T_\epsilon^m| \geq (1 - \delta) 2^{m(H(x, y) - \epsilon)}$$



upper bound:

$$\Pr((\tilde{x}^m, \tilde{y}^m) \in T_\epsilon^m) = \sum_{(x^m, y^m) \in T_\epsilon^m} p(x^m) p(y^m)$$

$$\leq \underbrace{2^{-m(H(x) - \epsilon)} 2^{-m(H(y) - \epsilon)}}_{\text{upper bound on probability of } p(x^m), p(y^m)} \cdot \underbrace{2^{m(H(x, y) + \epsilon)}}_{\substack{\text{upper bound} \\ \text{on } |T_\epsilon^m|}} = 2^{-m(I(x, y) - 3\epsilon)}$$

lower bound:

$$\Pr((\tilde{x}^m, \tilde{y}^m) \in T_\epsilon^m) = \sum_{(x^m, y^m) \in T_\epsilon^m} p(x^m) p(y^m) \geq (1 - \delta) 2^{m(H(x, y) - \epsilon)} \cdot 2^{-m(H(x) + \epsilon)} 2^{-m(H(y) + \epsilon)}$$

$$= (1 - \delta) 2^{-m(I(x, y) + 3\epsilon)}$$

Proof of channel coding theorem

(a) Let us fix $p(x)$. We generate 2^{mR} random sequences $x^m \rightarrow$ we obtain a certain $\{2^{mR}, m\}$ code \mathcal{C} . Probability to generate a particular code

$$P(\mathcal{C}) = \prod_{i=1}^{mR} \prod_{j=1}^m p(x_{ij}(m))$$

We choose a message $m \rightarrow x^m(m) = (x_1(m), \dots, x_m(m))$

receives receives y^m with probability $p(y^m | x^m(m)) = \prod_{i=1}^m p(y_i | x_i(m))$

He decides a message m' provided

such that $(x^m(m'), y^m)$ is a jointly typical sequence and there is no other w for which $(x^m(w), y^m)$ is typical, otherwise an error is declared.

What is the probability of error?

probability of error averaged over all choice of codes for a single message m : (by symmetry nothing depends on m)

in the

Th: $H(P_e) + P_e \log(|X|-1) \geq H(X|Y)$

Proof

Let E be random variable $= \begin{cases} 0 & X = \tilde{X} \\ 1 & X \neq \tilde{X} \end{cases}$

by chain rule:

$$H(E|X|Y) = H(E|Y) + H(X|E|Y)$$

$$= H(X|Y) + H(E|X|Y)$$

since $g(y)$ is deterministic

$$H(E|Y) \leq H(E) = H(P_e)$$

$$H(X|E|Y) = \sum_{e=0,y} p(e,y) H(X|e=0,y) + \sum_{e=1,y} p(e,y) H(X|e=1,y) \leq P_e \log(|X|-1)$$

$$H(P_e) + P_e \log(|X|-1) \geq H(E|Y) + H(X|E|Y) = H(X|Y)$$

\Downarrow

$$1 + P_e \log |X| \geq H(X|Y)$$

$$P_e \geq \frac{H(X|Y) - 1}{\log |X|}$$

Using Fano inequality:

$$\epsilon \geq \frac{H(M|Y^m) - 1}{mR}$$

$$\epsilon mR + 1 \geq H(M|Y^m)$$

By data processing inequality we have

$$\epsilon mR + 1 \geq H(X^m|Y^m)$$

- M be a uniform distribution

$$mR = H(M) = I(M; Y^m) + H(M|Y^m) \leq I(X^m(M); Y^m) + \epsilon mR + 1$$

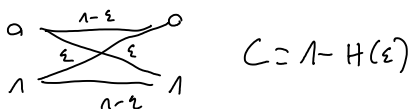
$$< mC + \epsilon mR + 1$$

$$\leq C + \epsilon R + \frac{1}{m} \xrightarrow{m \rightarrow \infty} R \leq C$$

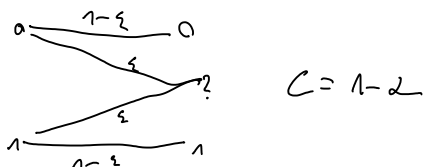
Shannon coding not really practical

3.4 Examples \uparrow

3.4.1 Binary symmetric channel ($\epsilon = 10\%$)



3.4.2 Binary erasure channel



3.5 Hamming codes (practical but suboptimal code)

