

Class no 7:

Joint entropy in 2 RVs

$$H(X, Y) = \sum_{(x, y) \in \text{supp}(P_{X, Y})} p(x, y) \log_2 \frac{1}{p(x, y)}$$

$X \times Y \triangleq \{(x, y) : x \in X, y \in Y\}$   
 "x" indicates Cartesian product

Note:

Suppose:  $X = Y$ . Then  $H(X, Y) = H(X) = H(Y)$

If  $X, Y$  independent  $H(X, Y) = H(X) + H(Y)$

Joint entropy in RVs  $X_1, \dots, X_n$

$$H(X_1, \dots, X_n)$$

$$\triangleq \sum_{(x_1, \dots, x_n) \in \text{supp}(P_{X_1, \dots, X_n})} p(x_1, \dots, x_n) \log_2 \frac{1}{p(x_1, \dots, x_n)}$$

$$P_{X, Y} : X \times Y \rightarrow [0, 1]$$

Joint distribution of  $X$  &  $Y$

$$\sum_{(x, y) \in X \times Y} p_{X, Y}(x, y) = 1$$

Chain Rule for Joint Entropy:

Lemma:

$$H(X_1, \dots, X_n) = H(X_1) + H(X_2 | X_1) + H(X_3 | X_1, X_2) + \dots + H(X_n | X_1, \dots, X_{n-1})$$

entropy of  $X_3$  given  $(X_1, X_2)$

where

$$H(X, Y | U, V) \text{ Conditional Joint entropy of } X \& Y \text{ on } U \& V$$

→ (see next page) ..

Recall:

$$H(X | Y) = \sum_{y \in \text{supp}(P_Y)} P_Y(y) \underbrace{H(X | Y = y)}_{\text{entropy of } X \text{ w/ distribution } P_{X|Y}(x|y)}$$

$$H(X | Y = y) = \sum_{x \in \text{supp}(P_{X|Y=y})} p(x|y) \log_2 \frac{1}{p(x|y)}$$

$$H(\underline{X}, \underline{Y} | \underline{U}, \underline{V}) \triangleq \sum_{\substack{(u,v) \in \text{supp}(P_{U,V}) \\ \downarrow \\ \text{small}(u,v)}} P_{U,V}(u,v) H(X, Y | U=u, V=v)$$

where

$$H(X, Y | U=u, V=v) \triangleq \sum_{(x,y) \in \text{supp}(P_{X,Y|U=u,V=v})} p(x,y|u,v) \log_2 \frac{1}{p(x,y|u,v)}$$

Note that  $P(x,y|u,v)$  is a <sup>valid</sup> joint distribution on RVs  $X$  &  $Y$ .

This can be extended to any number of variables before & after the conditioning.

Note:

How to understand  $p(x_1, x_2, x_3 | y_1, y_2)$ ?

Recall:

$$p(x_1, x_2) = \underset{x_1, x_2}{p(x_1)} \underset{x_2, x_1}{p(x_2|x_1)}$$

$$\left[ \begin{array}{l} p_{Y|X}(y|x) \\ \triangleq \frac{p(x,y)}{p(x)} \end{array} \right]$$

$$\rightarrow p(x_1, x_2, x_3 | y_1, y_2) = \frac{p(x_1, x_2 | y_1, y_2)}{p(x_3 | x_1, x_2, y_1, y_2)}$$

$$= \frac{p(x_3, y_2 | y_1) p(x_1, x_2 | x_3, y_1, y_2)}{p(y_2 | y_1)}$$

$$= \frac{p(\underbrace{x_3, y_2, x_1, x_2}_{\text{joint}} | \underline{y_1})}{p(\underline{y_2} | \underline{y_1})}$$

$$= P(x_1, x_2, x_3 | y_2, y_1)$$

Now we prove the chain rule:

Proof: Note the chain rule with the joint distribution  $P(x_1, \dots, x_n)$

$$= P(x_1) P(x_2, \dots, x_n | \underline{x_1})$$

$$= P(x_1) P(x_2 | x_1) P(x_3, \dots, x_n | x_1, x_2)$$

$$= P(x_1) P(x_2 | x_1) P(x_3 | x_1, x_2) P(x_4, \dots, x_n | x_1, x_2, x_3)$$

$$= \dots$$

$$= P(x_1) P(x_2 | x_1) P(x_3 | x_1, x_2) \dots P(x_n | x_1, \dots, x_{n-1})$$

Use defn of joint entropy  $H(X_1, \dots, X_n)$  to complete the proof [Exercise]

Class no 8:

$X \rightarrow$  entropy of  $H(X) \Rightarrow$  avg uncertainty abt  $X$

Observer ( $R_X$ ) sees ' $X$ '

$\rightarrow$  How much uncertainty does  $\text{Obs}(R_X)$  have abt  $X$ ?

Uncertainty abt  $X$  after  $\text{Obs } X = H(X|X)$

$= 0$

$\rightarrow$  What is the original unest in  $X$  (before observing  $X$ )  $= H(X)$   
 $\hookrightarrow$  (assume that  $P_X$  is known)

$\Rightarrow$  Reduction in <sup>of X</sup> avg uncertainty achieved by observing X

$$= H(X) - H(X|X) = H(X)$$

$\uparrow$  original                       $\uparrow$  uncer. after obs

Suppose  $X, Y$  are related to each other  
 (captured by a given joint prob. distribution

$$P_{X,Y}(x,y) = P(X=x, Y=y), \forall x,y$$

$\downarrow$  (obs know this)

Reduction in uncer. of  $X$  after observing  $Y$

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

$\uparrow$  avg uncertainty in  $X$  (left in  $X$ )  
 after obs.  $Y$ .

"Information gained abt  $X$  after observing  $Y$ "

$\downarrow$   
 "MUTUAL INFORMATION BETWEEN  $X$  &  $Y$ ",  $I(X;Y)$

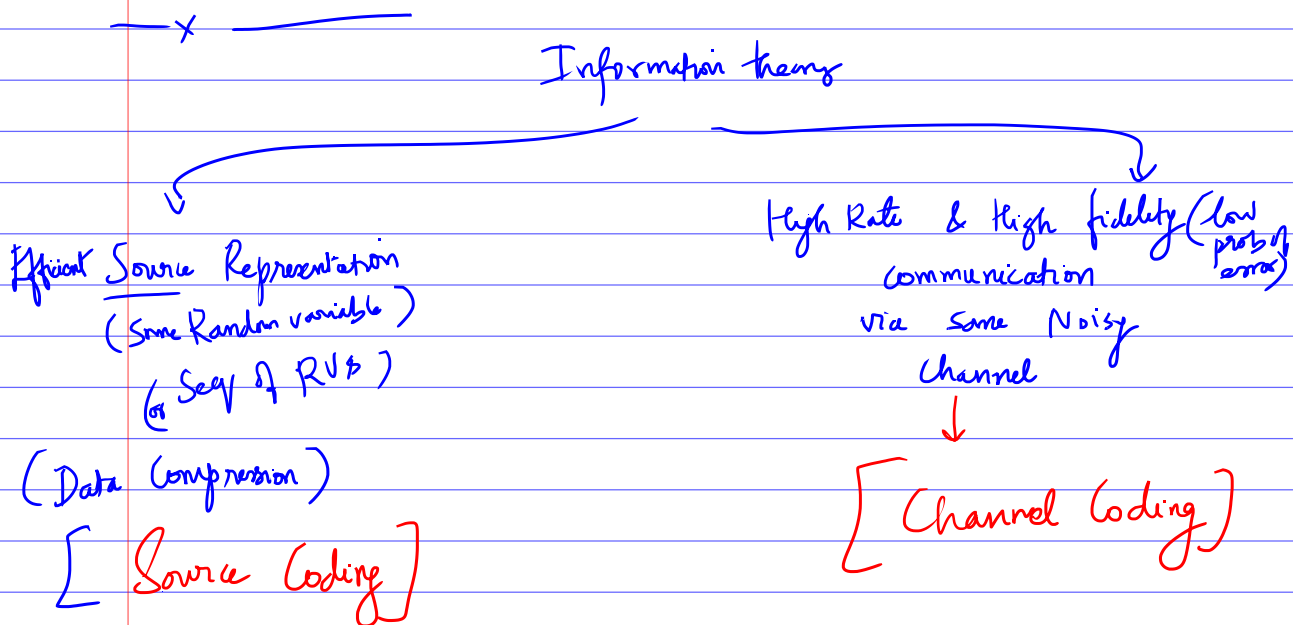
$$I(X;Y) \triangleq H(X) - H(X|Y).$$

Easy to show that  $H(X) - H(X|Y) = H(Y) - H(Y|X)$

$\rightarrow$  also equivalent to Info gained abt  $Y$  after observing  $X$ .

(1)  $I(X; Y) \leq \min(H(X), H(Y))$   
 ↑  
 is important because it tells that the information gained depends on 'grosser' 'quantized' observation.

(2)  $I(X; Y) \geq 0$  (Exercise: Represent  $I(X; Y)$  as a relative entropy between two prob distributions)



Overview of Source Coding:

Suppose we have  $X \in \{a, b\}$  a binary source with prob distribution  $P_X$

Suppose obs observes one instance of  $X$ , then wants to store (or communicate) it through some noise-free medium, which can carry/store only  $\{0, 1\}$  (bits)

In general we need 1 bit.  $a \rightarrow 0$   
 $b \rightarrow 1$   
(example)

But if  $R_x$  knows  $P_X$ , & it happens that

$$P_X(b) = 1, \quad P_X(a) = 0.$$

$\downarrow$   $\downarrow$   
 $(P(X=b))$   $P(X=a)$

then  $R_x$  need not even 'Read' or 'Receive' the encoded value to 'know'  $X$ . It can simply declare the value of  $X$  to be  $b$  & this will be correct / error-free with probability  $= 1$ .

$\Rightarrow$  Our code is not Required at all  
 $\Rightarrow$  we need a 0-length code.

Now, Suppose we are allowing for a small probability of error  $\Rightarrow P(\text{error}) \leq \epsilon$  for small  $\epsilon \in [0, 1)$

Now that if length of 'code' (stored symbol)  $= 1$ , then  $P(\text{error}) = 0$ .

Can we have length  $= 0$ , & some prob distribution for  $X$

(Ans: Yes! for which  $P(\text{error}) \leq \epsilon$ ?)

$\rightarrow$  Suppose  $P(X=b) = 1-\epsilon, P(X=a) = \epsilon$ .

For this suppose we use 0-length code,

what is  $P(\text{error}) = ?$

$\rightarrow$  "Declare  $X=b$  always" at decoder,  
this gives us  $P(\text{error}) \leq \epsilon$ .