

C0781 / QIC 890 Lec 08, Oct 4, 2016

Topic 3: classical communication via classical channels

Let  $X, Y$  be two rv's with joint distribution  $p(x, y)$   
and sample space  $\Omega_X \times \Omega_Y$

①  $H(X, Y) = - \sum_{x, y} p(x, y) \log p(x, y)$  (as defined earlier)

②  $\forall y \in \Omega_Y$ , let  $p(y) = \sum_x p(x, y)$  the marginal dist<sup>n</sup> on  $Y$ .

The conditional distribution of  $X$  given  $Y=y$ :

$$q_y(x) = \frac{p(x, y)}{p(y)} \quad \text{if } p(y) \neq 0$$

Def [Conditional entropy]

Using the above notations, the entropy of  $X$  conditioned on  $Y$  is:

$$H(X|Y) := \underbrace{\sum_y p(y)}_{\text{Averaged over } Y} \underbrace{H(q_y)}_{\text{Entropy of } X \text{ given } Y=y}$$

Thm (Chain rule)

$$H(X, Y) = H(Y) + H(X|Y).$$

Pf:  $H(X, Y) = - \sum_x \sum_y p(x, y) \log p(x, y)$   $\swarrow$   $p(y) q_y(x)$

$$= - \sum_y \underbrace{\sum_x p(x, y)}_{p(y)} \log p(y) - \underbrace{\sum_y \sum_x p(y) q_y(x) \log q_y(x)}_{H(q_y)}$$
$$= H(Y) + H(X|Y)$$

(or:  $H(X, Y) = H(Y, X)$ ,  $H(X, Y) = H(X) + H(Y|X)$ )

Def: [Relative entropy] [Kullback Leibler distance]

Let  $p(x), q(x)$  be two distributions on  $\Omega$ .

$$D(p||q) := \sum_{x \in \Omega} p(x) \log \frac{p(x)}{q(x)}$$

NB if  $\text{supp}(p) \not\subseteq \text{supp}(q)$  then  $D(p||q) = \infty$ ,  $0 \log \frac{0}{q} = 0$ ,  $p \log \frac{p}{q} = \infty$ .

NB:  $D$  non negative,  $D=0 \Leftrightarrow p=q$ , but not sym, no  $\Delta$  ineq (not metric)

\* Not used in our course:

- $D(p||q)$  measures the inefficiency of assuming that the dist<sup>n</sup> is  $q$  when it's  $p$ .
- In hypothesis testing, dist<sup>n</sup> is either  $p$  or  $q$  and  $n$  iid samples are given.

Any algorithm to discriminate  $p, q$ , WLOG is deterministic, and

$$\text{output} \begin{cases} p & \text{if } x^n \in A \\ q & \text{if } x^n \in A^c \end{cases} \text{ for some } A \subseteq \Omega^n.$$

Subject to the constraint  $\text{prob}(\text{output } p | p) \geq 1 - \epsilon$

$$\text{let } \beta(n, \epsilon) := \min_{\substack{A \subseteq \Omega^n \\ P(A) \geq 1 - \epsilon \text{ under } p^n}} q^n(A) \sim 2^{-n D(p||q)}.$$

Pf similar to AEP.

Def [Mutual Info]

$$\begin{aligned} I(X; Y) &:= D(p(x, y) || p(x) \cdot p(y)) \\ &= \sum_{x, y} p(x, y) \log \frac{p(x, y)}{p(x) \cdot p(y)} \end{aligned}$$

## Properties of H, D, & I: (Cover & Thomas (p 2))

(H1) Range:  $0 \leq H(X) \leq \log |X|$

"=" iff  $\exists$  a s.t.  $p(x) = 0 \forall x \neq a$

"=" iff  $p(x) = \frac{1}{|X|} \forall x$

(H2)  $H(X|Y) \geq 0$  (from def & (H1), conv combination of entropies)

(H3)  $H(XY) \geq H(X)$  (from (H2) & chain rule). *The more the merrier / messier*

(P4)  $D(p||q) \geq 0$ .

Pf: We need Jensen's ineq, that  $E f(x) \geq f(E(x))$  for any convex function  $f$  & rv  $X$ , with "=" iff  $X$  const.

$$\begin{aligned} \text{Then: } -D(p||q) &= -\sum_{x \in \text{supp}(p)} p(x) \log \frac{p(x)}{q(x)} \\ &= \sum_{x \in \text{supp}(p)} p(x) \log \frac{q(x)}{p(x)} \\ &\stackrel{\text{JI}}{\leq} \log \sum_{x \in \text{supp}(p)} p(x) \cdot \frac{q(x)}{p(x)} \\ &= 0 \end{aligned}$$

Since  $\log(\cdot)$  is strictly concave,  $D(p||q) = 0$  iff  $\frac{p}{q} = \text{const}$

iff  $p(x) = q(x) \forall x$

(I5)  $I(X;Y) \geq 0$ , "=" iff  $X, Y$  independent.

$$\textcircled{I6} \quad I(X:Y) = H(X) - H(X|Y) \\ = H(Y) - H(Y|X)$$

$$\text{Chain rule} \rightarrow = H(X) + H(Y) - H(XY)$$

Interpretation:  $I(X:Y)$  = decrease in ignorance of  $X$  given  $Y$   
= amount of info of  $X$  carried by  $Y$ .

Note:  $\textcircled{I6}$  follows from the defs of  $H$  &  $I$   
but it is NOT a definition.

$$\begin{aligned} \text{Pf: } I(X:Y) &= \sum_{xy} p(xy) \log \frac{p(xy)}{p(x)p(y)} \\ &= \sum_{xy} p(xy) \log \frac{f_y(x)}{p(x)} \\ &= \sum_x \underbrace{\left( \sum_y p(xy) \right)}_{p(x)} \cdot \log \frac{1}{p(x)} + \underbrace{\sum_{xy} p(y) \cdot f_y(x) \log f_y(x)}_{-H(f_y(x))} \\ &= H(X) - H(X|Y) \end{aligned}$$

$\textcircled{H7}$  Subadditivity (SA)

$$H(XY) \leq H(X) + H(Y), \quad "=" \text{ iff } X, Y \text{ indep.}$$

Pf: follows from  $\textcircled{I5}$  &  $\textcircled{I6}$

$\textcircled{H8}$  Conditioning reduces entropy

$$H(X|Y) \leq H(X), \quad "=" \text{ iff } X, Y \text{ indep.}$$

Pf: from  $\textcircled{I5}$  &  $\textcircled{I6}$

(H9) For 3 rv's  $X, Y, Z$ ,  $H(X|Z) \geq H(X|YZ)$ .

Pf: Let  $p(xyz)$  be the joint dist<sup>n</sup> for  $XYZ$ .

For each  $z$ , let  $q_z(xy) = \frac{p(xyz)}{p(z)}$

where  $p(z) = \sum_{xy} p(xyz)$ .

Let  $\tilde{X}\tilde{Y}$  has distribution  $q_z(xy)$ .

We have  $H(\tilde{X}) \geq H(\tilde{X}|\tilde{Y})$  from previous discussion (\*)

•  $H(\tilde{X}) = \text{Entropy on } \eta_z(x) = \sum_y q_z(xy)$

But  $\eta_z(x) = \sum_y \frac{p(xyz)}{p(z)} = \frac{p(xz)}{p(z)} = p_z(x)$  from  $\sum_y p(xyz)$

$\therefore H(\tilde{X}) = H(X|Z=z)$  evaluated on  $\sum_y p(xyz)$ .

•  $H(\tilde{X}|\tilde{Y}) = \sum_y \left( \sum_x q_z(xy) \right) \cdot H(\tilde{X}|\tilde{Y}=y)$

$H(\tilde{X}|\tilde{Y}=y) = \text{Entropy on } \frac{q_z(xy)}{\sum_x q_z(xy)}$

$$= \frac{\frac{p(xyz)}{p(z)}}{\sum_x \frac{p(xyz)}{p(z)}} = \frac{p(xyz)}{\sum_x p(xyz)} = \text{dist on } X \text{ given } Y=y, Z=z \text{ evaluated on } p(xyz)$$

$$\sum_x q_z(xy) = \sum_x \frac{p(xyz)}{p(z)} = \frac{p(yz)}{p(z)} = \text{prob of } y \text{ given } z=z.$$

$\therefore H(\tilde{X}|\tilde{Y}) = H(X|YZ=z)$

- Now take convex combination of the inequality over  $z$

$$\sum_z p(z) H(\hat{X}) \geq \sum_z p(z) H(\hat{X}|Y)$$

$$\parallel \qquad \parallel$$

$$H(X|Z) \qquad H(X|YZ)$$

(I10) Def: The conditional mutual information

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

$$= H(Y|Z) - H(Y|XZ)$$

By (H9),  $I(X:Y|Z) \geq 0$

(H10) Concavity: mixing increases entropy

Let  $\alpha_k > 0$ ,  $\sum_{k=1}^t \alpha_k = 1$ ,  $p_1, \dots, p_t$  distributions on  $\mathcal{X}$

$$\text{Then } \sum_{k=1}^t \alpha_k H(p_k) \leq H\left(\sum_{k=1}^t \alpha_k p_k\right)$$

Pf: let  $K_X$  be rv's on  $\{1, 2, \dots, t\} \times \mathcal{X}$

$$\text{prob}(k, x) = \alpha_k p_k(x)$$

$$\text{Then } \sum_{k=1}^t \alpha_k H(p_k) = H(X|K)$$

$$H\left(\sum_{k=1}^t \alpha_k p_k\right) = H(X)$$

$\therefore$  (H9) follows from (H8)

(H11) Strong subadditivity SSA or data processing inequality DPI

• Def: We say  $X \rightarrow Y \rightarrow Z$  (a Markov chain)

$$\text{if } p(x, y, z) = p(x) p(y|x) p(z|y)$$

$$\text{(ie if } p(z|y) = p(z|xy))$$

• DPI: If  $X \rightarrow Y \rightarrow Z$ , then  $I(X; Y) \geq I(X; Z)$ .

$$\text{Pf: } I(X; YZ)$$

$$= \underbrace{H(X) + H(Z) - H(XZ)}_{I(X; Z)} + H(XZ) - H(Z) - (H(XYZ) - H(YZ))$$

$$= I(X; Z) + H(X|Z) - H(X|YZ) \stackrel{\text{(H9)}}{\geq} I(X; Z)$$

$$I(X; YZ)$$

$$= \underbrace{H(X) + H(Y) - H(XY)}_{I(X; Y)} + \underbrace{H(XY) - H(Y)}_{H(X|Y)} - \underbrace{[H(XYZ) - H(YZ)]}_{H(Z|XY)}$$

$$= I(X; Y) + H(Z|Y) - H(Z|XY) = I(X; Y)$$

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

$\uparrow$   
Since  $p(z|y) = p(z|xy)$

- In particular, if we process  $Y$  to obtain  $Z$  for any  $X$  that may be correlated with  $Y$ , we have  $p(Z|Y) = p(Z|XY)$ .

So such processing can increase mutual info (from  $I(X:Y)$  to  $I(X:Z)$ ).



### Def [Jointly typical sequence]

Given a distribution  $p(x,y)$ , drawn iid  $n$  times

$x^n y^n$  is  $\delta$ -jointly typical if

$$(a) \left| -\frac{1}{n} \log p(x^n) - H(X) \right| \leq \delta \quad (x^n \text{ typical})$$

$$(b) \left| -\frac{1}{n} \log p(y^n) - H(Y) \right| \leq \delta \quad (y^n \text{ typical})$$

$$(c) \left| -\frac{1}{n} \log p(x^n y^n) - H(X,Y) \right| \leq \delta \quad ((x,y)^n \text{ typical})$$

### Def [Jointly typical set]

$$A_{n,\delta} = \{ x^n y^n \in \Omega_x^n \times \Omega_y^n : x^n y^n \text{ jointly typical} \}$$

Obs: if  $x^n y^n \in A_{n,\delta}$ ,

$$\text{then } p(x^n | y^n) = \frac{p(x^n y^n)}{p(y^n)} \leq \frac{2^{-n(H(X,Y) - \delta)}}{2^{-n(H(Y) + \delta)}} \leq 2^{-n(H(X|Y) - 2\delta)}$$

$$p(x^n | y^n) = \frac{p(x^n y^n)}{p(y^n)} \geq \frac{2^{-n(H(X,Y) + \delta)}}{2^{-n(H(Y) - \delta)}} \geq 2^{-n(H(X|Y) + 2\delta)}$$

NB. For this  $y^n$ , there are  $\approx 2^{nH(X|Y)}$   $\tilde{x}^n$ 's s.t.  $\tilde{x}^n y^n \in A_{n,\delta}$

Pf = similar to that in AEP.

## Thm [Joint AEP]

Using above defs,  $\forall \epsilon > 0 \quad \forall \delta > 0 \quad \exists n_0$  s.t.  $\forall n > n_0$

①  $\Pr(A_{n,\delta}) \geq 1 - \epsilon$

②  $(1 - \epsilon) 2^{n(H(X,Y) - \delta)} \leq |A_{n,\delta}| \leq 2^{n(H(X,Y) + \delta)}$

③ Suppose  $x^n y^n$  is drawn according to the following distribution

$$f(x^n y^n) = p(x^n) \cdot p(y^n)$$

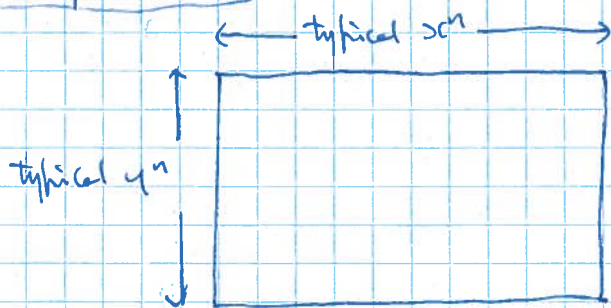
then  $2^{-n(I(X;Y) + 3\delta)} \leq \Pr_f(x^n y^n \in A_{n,\delta}) \leq 2^{-n(I(X;Y) - 3\delta)}$

Pf (in Cover & Thomas (p 8))

Ingredients:

The AEP itself, the union bound, etc.

Keep on board



$$\text{Total} \approx 2^{nH(X)} \times 2^{nH(Y)} \text{ entries}$$

The entry associated with  $x^n y^n$  is

$$\begin{cases} 1 & \text{if } x^n y^n \in A_{n,\delta} \\ 0 & \text{otherwise} \end{cases}$$

① says that tolerating a prob of failure  $\epsilon$ , we can focus on this table

② says there are  $\approx 2^{nH(X,Y)}$  "1"s

Obs says each column has  $\approx 2^{nH(Y|X)}$  "1"s  
 each row  $\approx 2^{nH(X|Y)}$  "1"s

③ says a random entry in the table has prob =  $2^{-nI(X;Y)}$  to be "1".

Example:  $\mathcal{X} = \mathcal{Y} = \{0, 1\}$

$$p(00) = \frac{1}{2}(1-e)$$

$$p(01) = \frac{1}{2}e$$

$$p(10) = \frac{1}{2}e$$

$$p(11) = \frac{1}{2}(1-e)$$

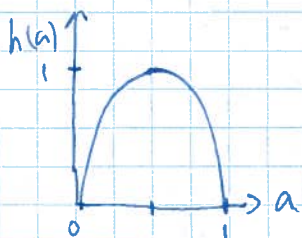
for  $e \in [0, 1]$ , say  $e = 0.1$ .

$$H(X) = H(Y) = 1$$

both marginals uniform

$$\begin{aligned} H(XY) &= 2 \left(-\frac{1}{2}\right) (1-e) \log\left(\frac{1}{2}(1-e)\right) + 2 \left(-\frac{e}{2}\right) \log\left(\frac{e}{2}\right) \\ &= 1 + h(e) = 1.469 \end{aligned}$$

where  $h(a) = -a \log a - (1-a) \log(1-a)$  is the binary entropy function



$$H(X|Y) = H(XY) - H(Y) = 0.469.$$

$$I(X:Y) = H(X) + H(Y) - H(XY) = 0.531$$

$\therefore$  There are  $\approx 2^{n \cdot 1.469}$  jointly typical  $x^n y^n$ 's.

For each typical  $y^n$ , there are  $\approx 2^{n H(X|Y)} = 2^{n \cdot 0.469}$   $\tilde{x}^n$ 's such that  $\tilde{x}^n y^n$  is jointly typical.

For an  $\tilde{x}^n$  chosen randomly for the typical set for  $X_1, \dots, X_n$ ,

$$\text{Prob}(\tilde{x}^n y^n \notin A_{n,\delta}) \approx \frac{2^{n H(X|Y)}}{2^{n H(X)}} = 2^{-n I(X:Y)} \approx 2^{-0.531n}.$$

# App 1 = Distributed source coding

Goal: Sample  $XY$   $n$  times iid, get  $x^n y^n$ .

Give  $x^n$  to Alice,  $y^n$  to Bob.

How many cbits from Alice to Bob is needed for Bob to learn  $x^n$ ?

Ans:  $\approx n(H(X|Y) + \epsilon)$  bits (and let's see what  $\epsilon$  should be)

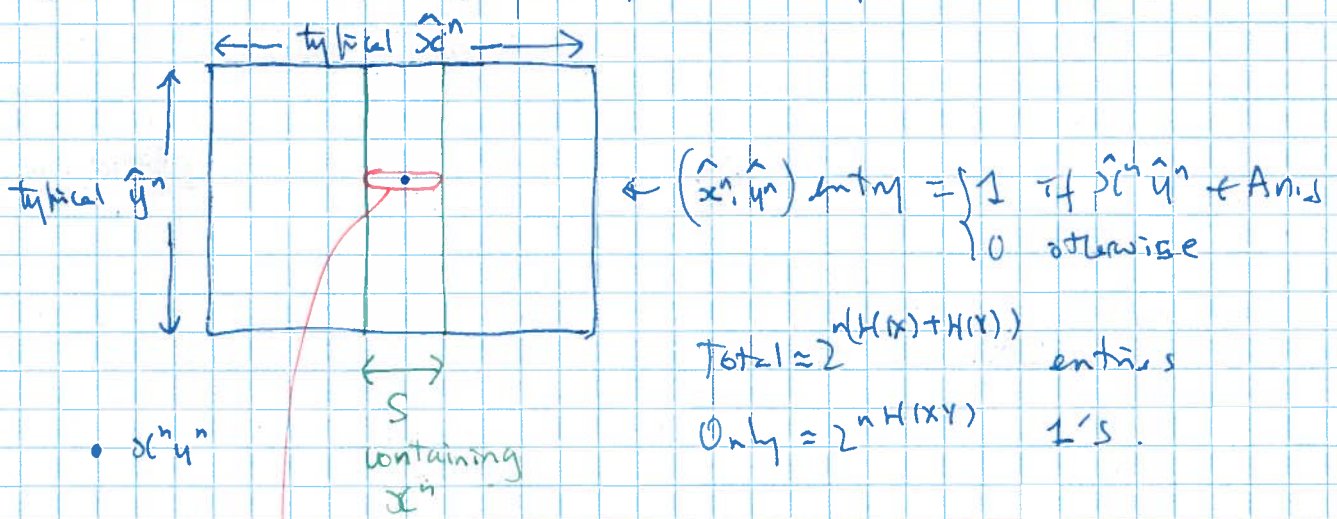
Method: Slepian-Wolf coding. Let  $T = T_{\epsilon}$  for  $X_1, X_2, \dots, X_n$ .

① Alice & Bob pre agree on a partition of  $T$  into  $\approx 2^{n(H(X|Y) + \epsilon)}$  sets each with at most  $\frac{2^{n(H(X) + \delta)}}{2^{n(H(X|Y) + \epsilon)}} \approx 2^{n(I(X;Y) + \delta - \epsilon)}$  elements by AEP.

② When Alice receives  $x^n$ , she tells Bob which of these  $2^{n(H(X|Y) + \epsilon)}$  sets contain  $x^n$ . Call the set  $S$ .

③ When Bob receives the label for  $S$ , he finds all  $\hat{x}^n \in S$  that are jointly typical with  $y^n$ . Call this set  $J$ . Output "ERR" if  $|J| = 0$  or  $|J| > 1$ .

Why this works? For large enough  $n$ , with prob  $\geq 1 - \epsilon$ ,  $x^n y^n \in A_{n,\epsilon}$ .



Bob checks all  $\hat{x}^n y^n$  that are jointly typical.

The next step relies on how  $T$  is partitioned into the sets.

Idea: from all possible partitions, pick one at random (say  $P$ ) and analyse the prob  $P$  works.

In this case,  $S$  contains  $x^n$ , but all other elements are chosen at random. Each such  $\hat{x}^n$  has a prob

$$\approx \frac{2^{nH(X|Y)}}{2^{nH(X)}} \approx 2^{-nI(X:Y)} \text{ to be jointly typical with } y^n.$$

# 1's in the row labeled by  $y^n$  # entries in the row.

More precisely,  $\hat{x}^n$  &  $y^n$  are independent

$$\therefore \text{Prob}_P(\hat{x}^n y^n \in A_{n,\delta}) \leq 2^{-n(I(X:Y) - 3\delta)} \text{ from the JSEP.}$$

$$\therefore \text{Prob}_P(\exists \hat{x}^n \in S, \hat{x}^n \neq x^n, \hat{x}^n y^n \in A_{n,\delta})$$

$$\leq |S| \cdot 2^{-n(I(X:Y) - 3\delta)} \leq 2^{-n(4\delta - \alpha)}$$

union bound over each elt in  $S$

$$\therefore \text{Choosing } \alpha = 5\delta, \text{ prob}(P \text{ fails}) \leq \epsilon + 2^{-n\delta} \text{ (indep of } y^n).$$

$\swarrow$   $\searrow$   
 $x^n y^n \notin A_{n,\delta}$      $S \text{ bad}$

$\therefore \exists$  partition that works, requiring  $n(H(X|Y) + 5\delta)$  bits.