# Introduction to Classical Information Theory

James Schneeloch

Lecture for physics 407

Thursday 12-5-2013

Cover & Thomas Ch. 2

Sections 1-6

- It is a mathematical theory of Communication

  Shannon, 1948 "A mathematical theory of Communication"

  $\longrightarrow$ over 65,000 citations

  ( EPR paradox paper )
  ( has ~11,000 citations )
  ( Bell inequality paper )
  ( has ~7900 citations )

  Uses:

  $\longrightarrow$ Understanding limits of communicating data

  $\longrightarrow$ data compression

  $\longrightarrow$ statistical inference

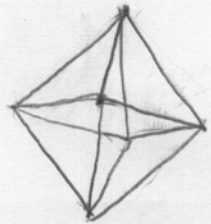  $\longrightarrow$ cryptography

## How to quantify information?

- One symbol can have different meanings
- We quantify information in the context of <u>communication</u> (in bits)
  * We can't say how many bits a symbol has; we can say how many bits it takes to convey that symbol to others

## What is a bit?

- The information conveyed in the answer of a yes/no question.

  ( 20 questions $\longrightarrow$ 20 bits )

- Bits are written as --

  (on or off) (0 or 1) (bright or dark)
  (red or green) (horizontal or vertical), etc...

ex/

Let's say you have an 8-sided _loaded_ die.

Random Variable $\underline{X} = \{X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8\}$

: $X_i$ = "lands on side (i)" | $P(X_i)$ = probability that it lands on side (i)

$$P(\underline{X}) = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}, \frac{1}{64}, \frac{1}{128}, \frac{1}{128} \right\}$$

- Someone rolls the die and hides the outcome.

  - You want to figure out the outcome
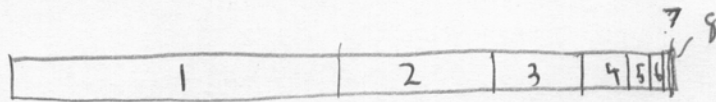
  - You may ask only yes/no questions.

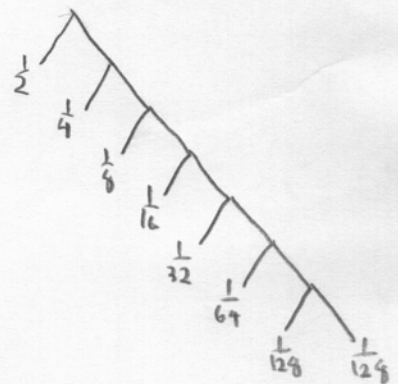① How many questions (bits) does it take __to be sure__ of the outcome of this roll?

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |

$3 = \log_2(8)$

3 bits are needed

② How many questions (bits) does it take __on average__ per roll over many rolls?

| 1 | 2 | 3 | 4 | 5 | 7 8 |

$$\langle \text{\# of questions} \rangle = \frac{1}{2}(1) + \frac{1}{4}(2) + \frac{1}{8}(3)$$
$$+ \frac{1}{16}(4) + \frac{1}{32}(5) + \frac{1}{64}(6)$$
$$+ \frac{2}{128}(7) = \frac{127}{64} \sim 1.98 \text{ bits}$$
$$(\text{less than } 3)$$

* Tangent:

- How many bits does it take to be sure of the thermodynamic microstate of a liter of water at room temperature?

$$S = k_B \ln(\Omega)$$

$$\# \text{ of bits} = \log_2(\Omega) = \frac{\ln(\Omega)}{\ln(2)} = \frac{S}{k_B \ln(2)}$$

$$S_{H_2O} \sim 6 \, \text{J/mol·K} \quad \text{at } 20-25°C$$

there are about 55 mol of $H_2O$ in 1 liter of it.

so $\qquad \# \text{ of bits} \sim 3.5 \times 10^{25} \text{ bits}$

- The highest information transfer rate over an optical fiber is currently about $10^{15}$ bits/s $\left( \substack{\text{with a} \\ \text{12-core fiber}} \right)$

  - At this rate, how long would it take to transfer the information of the microstate of that liter of water?

$$\text{Total time} = \frac{3.5 \times 10^{25} \text{ bits}}{10^{15} \text{ bits/s}} \sim 3.5 \times 10^{10} \text{ seconds}$$

$$\text{or} \sim \underline{1100 \text{ years}}$$

(1 billion seconds is $\sim$ 32 years)

# Entropy

| The Shannon entropy (in bits) | $H_2(\underline{X}) \equiv -\sum_{x_i \in \underline{X}} P(x_i) \log_2(P(x_i))$ |
|---|---|

convention:

$0 \log_b 0 \equiv 0$

since

$\lim_{z \to 0} z \log_b z = 0$

$H_2(\underline{X}) \rightarrow$ The minimum average number of bits needed to communicate the outcome of $\underline{X}$.

$\quad\quad\quad \rightarrow$ A measure of the inherent uncertainty in the outcome of $\underline{X}$.

Entropy can be measured in different bases

$b = 2$ "bits"
$b = 3$ "trits"
$b = e$ "nats"

$$H_b(\underline{X}) = -\sum_{x_i \in \underline{X}} P(x_i) \log_b(P(x_i))$$

Note: A "trit" is the amount of information conveyed by answering a 3-answer question
(more, less, same)
(here, there, neither)

$\left[ \text{Let's just use bits} \rightarrow H(\underline{X}) \Rightarrow H_2(\underline{X}) \text{ (by convention)} \right]$

$\begin{pmatrix} \text{Useful} \\ \text{form} \end{pmatrix} \quad H(\underline{X}) = \left\langle \log\left(\frac{1}{P(\underline{X})}\right) \right\rangle_{P(\underline{X})}$

## Elementary properties

① $H(\underline{X}) \geq 0$.

$\quad P(x_i) \in [0,1] \Rightarrow \log_b\left(\frac{1}{P(x_i)}\right) \geq 0$, so $H(\underline{X}) \geq 0$

② $H_b(\underline{X}) = (\log_b a) H_a(\underline{X})$

$\left( \begin{array}{l} \text{change of base formula} \\ \text{for logarithms} \end{array} \right)$

$\boxed{\log_b P = \log_b a \, \log_a P}$

ex/ Biased coin-toss

$$\underline{X} = \{X_1, X_2\} \quad X_1 = \text{"heads"}$$
$$X_2 = \text{"tails"}$$

$$P(\underline{X}) = \{P, 1-P\}$$

what is $H(\underline{X})$?

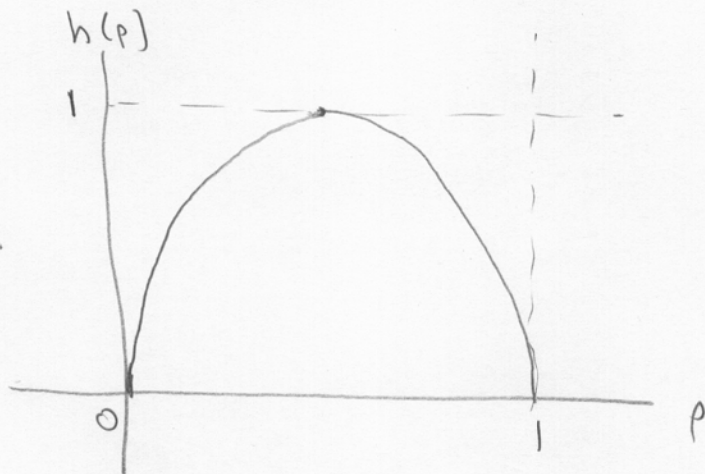$$H(\underline{X}) = -\sum_{x_i \in \underline{X}} P(X_i) \log(P(X_i))$$

the binary entropy function

$$\underline{H(\underline{X}) = -p\log_2(p) - (1-p)\log_2(1-p) \equiv \boxed{h(p)}}$$

As $p \to 0$ or $p \to 1$

$h(p) \to 0$

"A certain outcome requires no bits to communicate."



If $p = \frac{1}{2}$, we need on average $1$ bit per coin toss.

---

ex/. 8-sided loaded die

$$P(\underline{X}) = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}, \frac{1}{64}, \frac{1}{128}, \frac{1}{128} \right\}$$

what is $H(\underline{X})$?

$$H(\underline{X}) = -\frac{1}{2}\log\left(\frac{1}{2}\right) - \frac{1}{4}\log\left(\frac{1}{4}\right) - \sim\!\!\sim\!\!\sim$$

$$= \frac{1}{2}\log(2) + \frac{1}{4}\log(4) + \sim\!\!\sim\!\!\sim$$

$$= \frac{1}{2}(1) + \frac{1}{4}(2) + \frac{1}{8}(3) + \frac{1}{16}(4) + \sim$$

$$= \frac{127}{64} \text{ bits} \quad (\text{just like before})$$

# Joint and Conditional Entropy

$\left(\begin{array}{c}\text{Marginal} \\ \text{Entropy}\end{array}\right)$ $\quad H(\underline{X}) = \left\langle \log\left(\frac{1}{P(\underline{X})}\right)\right\rangle_{P(\underline{X})}$ $\quad$ "Average # of bits you need to communicate the outcome of $\underline{X}$"

## Joint Entropy

$$H(\underline{X},\underline{Y}) \equiv -\sum_{X_i, Y_j \in (\underline{X}\underline{Y})} P(X_i, Y_j) \log\left(P(X_i, Y_j)\right) = \left\langle \log\left(\frac{1}{P(\underline{X},\underline{Y})}\right)\right\rangle_{P(\underline{X},\underline{Y})}$$

$H(\underline{X},\underline{Y}) \rightarrow$ "Average # of bits you need to communicate the outcomes of both $\underline{X}$ and $\underline{Y}$"

Note: If $\underline{X}$ and $\underline{Y}$ are independent, then $H(\underline{X},\underline{Y}) = H(\underline{X}) + H(\underline{Y})$

## Conditional Entropy

$$H(\underline{Y}|\underline{X}) \equiv -\sum_{X_i, Y_j \in (\underline{X},\underline{Y})} P(X_i, Y_j) \log\left(P(Y_j|X_i)\right) = \left\langle \log\left(\frac{1}{P(\underline{Y}|\underline{X})}\right)\right\rangle_{P(\underline{X},\underline{Y})}$$

Note:

$$H(\underline{Y}|\underline{X}) = \sum_{X_i \in \underline{X}} P(X_i)\, H(\underline{Y}|\underline{X} = X_i)$$

where $H(\underline{Y}|\underline{X} = X_i) = -\sum_{Y_j \in \underline{Y}} P(Y_j|X_i) \log\left(P(Y_j|X_i)\right)$

Note: $\quad H(\underline{X}|\underline{Y}) \neq H(\underline{Y}|\underline{X})$

Note: From Bayes' Rule

$$\boxed{H(\underline{X},\underline{Y}) = H(\underline{X}) + H(\underline{Y}|\underline{X})}$$

Chain rule of the joint entropy

because $\log\left(\frac{1}{P(A,B)}\right) = \log\left(\frac{1}{P(A)P(B|A)}\right) = \log\left(\frac{1}{P(A)}\right) + \log\left(\frac{1}{P(B|A)}\right)$

ex) Let $\underline{X},\underline{Y}$ have the following joint distribution

| $\underline{X}$ \ $\underline{Y}$ | a | b | c | d |
|---|---|---|---|---|
| a | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{1}{32}$ | $\frac{1}{32}$ |
| b | $\frac{1}{16}$ | $\frac{1}{8}$ | $\frac{1}{32}$ | $\frac{1}{32}$ |
| c | $\frac{1}{16}$ | $\frac{1}{16}$ | $\frac{1}{16}$ | $\frac{1}{16}$ |
| d | $\frac{1}{4}$ | 0 | 0 | 0 |

marginals

| | a | b | c | d |
|---|---|---|---|---|
| $\underline{X}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ |
| $\underline{Y}$ | $\frac{1}{2}$ | $\frac{1}{4}$ | $\frac{1}{8}$ | $\frac{1}{8}$ |

What are $H(\underline{X}), H(\underline{Y}), H(\underline{X},\underline{Y}), H(\underline{X}|\underline{Y})$, and $H(\underline{Y}|\underline{X})$?

(skip in lecture)

$$H(\underline{X}) = -\sum_{x_i} P(x_i) \log(P(x_i)) = -\frac{1}{4}\log\left(\frac{1}{4}\right) - \frac{1}{4}\log\left(\frac{1}{4}\right) - \frac{1}{4}\log\left(\frac{1}{4}\right) - \frac{1}{4}\log\left(\frac{1}{4}\right)$$

$$H(\underline{x}) = 2 \text{ bits}$$
$$H(\underline{Y}) = \frac{7}{4} \text{ bits}$$

note:
$$H(\underline{X}|\underline{Y}) \neq H(\underline{Y}|\underline{X})$$

$$H(\underline{X},\underline{Y}) = \frac{27}{8} \text{ bits}$$
$$H(\underline{X}|\underline{Y}) = H(\underline{X},\underline{Y}) - H(\underline{X}) = \frac{11}{8} \text{ bits}$$
$$H(\underline{Y}|\underline{X}) = \frac{13}{8} \text{ bits}$$

## Mutual Information

(The Shannon mutual information)

$$H(\underline{X}:\underline{Y}) \equiv \sum_{x_i,Y_j \in \underline{X}\,\underline{Y}} P(X_i,Y_j) \log\left(\frac{P(X_i,Y_j)}{P(X_i)P(Y_j)}\right)$$

$$H(\underline{X}:\underline{Y}) = \left\langle \log\left(\frac{P(\underline{X},\underline{Y})}{P(\underline{X})P(\underline{Y})}\right)\right\rangle_{P(\underline{X},\underline{Y})}$$

$$H(\underline{X}:\underline{Y}) = \left\langle \log\left(\frac{1}{P(\underline{X})}\right) - \log\left(\frac{1}{P(\underline{X}|\underline{Y})}\right)\right\rangle_{P(\underline{X},\underline{Y})}$$

$$\boxed{H(\underline{X}:\underline{Y}) = H(\underline{X}) - H(\underline{X}|\underline{Y})}$$

$$= H(\underline{Y}) - H(\underline{Y}|\underline{X})$$

→ "average number of bits communicated about the outcome of $\underline{Y}$ by communicating the outcome of $\underline{X}$.

NOTE: From $H(\underline{X}|\underline{Y}) = H(\underline{X},\underline{Y}) - H(\underline{Y})$

$$\boxed{H(\underline{X}:\underline{Y}) = H(\underline{X}) + H(\underline{Y}) - H(\underline{X},\underline{Y})}$$

Elementary properties

① $H(\underline{X}:\underline{Y}) = H(\underline{Y}:\underline{X})$

② $H(\underline{X}:\underline{X}) = H(\underline{X})$, because $H(\underline{X},\underline{X}) = H(\underline{X})$
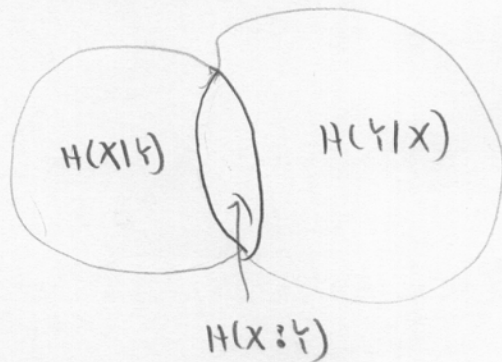
Conditional Mutual Information

$$H(\underline{X}:\underline{Y}|\underline{Z}) \equiv \sum_{x_i, y_j, z_k \in \underline{XYZ}} P(X_i, Y_j, Z_k) \log\left(\frac{P(X_i, Y_j | Z_k)}{P(X_i|Z_k) P(Y_j|Z_k)}\right)$$

$$H(\underline{X}:\underline{Y}|\underline{Z}) = \left\langle \log\left(\frac{P(X,Y|Z)}{P(X|Z) P(Y|Z)}\right)\right\rangle_{P(\underline{X}\underline{Y}\underline{Z})}$$
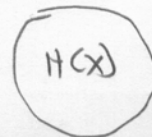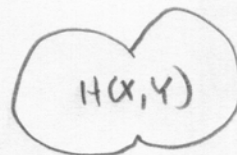
$$\boxed{H(\underline{X}:\underline{Y}|\underline{Z}) = H(\underline{X}|\underline{Z}) + H(\underline{Y}|\underline{Z}) - H(\underline{X},\underline{Y}|\underline{Z})}$$

$\hookrightarrow$ " Average # of bits communicated about the outcome of $\underline{Y}$ by communicating the outcome of $\underline{X}$, given that the outcome of $\underline{Z}$ is known "

(Entropy Venn Diagram) ?



$H(X|Y)$     $H(Y|X)$

$H(X:Y)$

" amount of uncertainty )

$H(X,Y)$

$H(X)$

$H(Y)$

## Relative Entropy  (for any 2 distributions $P(\underline{X}), Q(\underline{X})$)

$$D(P(\underline{X}) \| Q(\underline{X})) \equiv \sum_{X_i \in \underline{X}} P(X_i) \log \left( \frac{P(X_i)}{Q(X_i)} \right)$$

(a.k.a. Kullback-Leibler divergence)

- It is a measure of divergence between two probability distributions

- It is a measure of inefficiency of coding the outcomes accordig to $Q(\underline{X})$, when true distribution is $P(\underline{X})$

Convention?
$$0 \log \left( \frac{0}{0} \right) = 0$$
$$0 \log \left( \frac{0}{q} \right) = 0$$
$$p \log \left( \frac{p}{0} \right) \to \infty$$

Vocab:

code: The system of assignment of code "words" to each outcome of $\underline{X}$.

code word: Sequence of (binary) digits assigned to particular outcome of $\underline{X}$.

- - - - - - - - - - - - - - - - - - - - - -

$H(\underline{X})$ = Minimum average length of (binary) codeword to describe outcome of $\underline{X}$ with distribution $P(\underline{X})$ ← true

★ If we know the true distribution $P(\underline{X})$, we can ideally construct a code of average length $H(\underline{X})$

- If we constructed a code for $Q(\underline{X})$, when the distribution really was $P(\underline{X})$, our average codeword length would be longer by $D(P(\underline{X}) \| Q(\underline{X}))$

- - - - - - - - - - - - - - - - - - - - - -

$D(P(\underline{X}) \| Q(\underline{X}))$ is **not** a distance measa

$$D(P(\underline{X}) \| Q(\underline{X})) \neq D(Q(\underline{X}) \| P(\underline{X}))$$

Look up triangle inequality for distance metrics

$$D(P(\underline{X}) \| Q(\underline{X})) + D(Q(\underline{X}) \| R(\underline{X})) \ngeq D(P(\underline{X}) \| R(\underline{X}))$$

$$D(P(\underline{X}) \| Q(\underline{X})) = \left\langle \log\left(\frac{P(\underline{X})}{Q(\underline{X})}\right) \right\rangle_{P(\underline{X})}$$

Note:

$$\left[\begin{array}{l} \text{The mutual information} \\ \text{is also a relative} \\ \text{entropy} \end{array}\right] \rightarrow H(\underline{X}; \underline{Y}) = D(P(\underline{X}, \underline{Y}) \| P(\underline{X}) P(\underline{Y}))$$

## Conditional Relative Entropy

$$D(P(\underline{Y}|\underline{X}) \| Q(\underline{Y}|\underline{X})) \equiv \sum_{X_i, Y_j \in (\underline{X}\,\underline{Y})} P(X_i; Y_j) \log\left(\frac{P(Y_j|X_i)}{Q(Y_j|X_i)}\right)$$

$$D(P(\underline{Y}|\underline{X}) \| Q(\underline{Y}|\underline{X})) = \left\langle \log\left(\frac{P(\underline{Y}|\underline{X})}{Q(\underline{Y}|\underline{X})}\right) \right\rangle_{P(\underline{X},\underline{Y})} \quad \text{just the}$$

Note: With Bayes' Rule

$$\left\langle \log\left(\frac{P(\underline{X},\underline{Y})}{Q(\underline{X},\underline{Y})}\right) \right\rangle_{P(\underline{X},\underline{Y})} = \left\langle \log\left(\frac{P(\underline{X})P(\underline{Y}|\underline{X})}{Q(\underline{X})Q(\underline{Y}|\underline{X})}\right) \right\rangle_{P(\underline{X},\underline{Y})}$$

$$= \left\langle \log\left(\frac{P(\underline{X})}{Q(\underline{X})}\right) \right\rangle_{P(\underline{X},\underline{Y})} + \left\langle \log\left(\frac{P(\underline{Y}|\underline{X})}{Q(\underline{Y}|\underline{X})}\right) \right\rangle_{P(\underline{X},\underline{Y})}$$

so that

$$D(P(\underline{X},\underline{Y}) \| Q(\underline{X},\underline{Y})) = D(P(\underline{X}) \| Q(\underline{X})) + D(P(\underline{Y}|\underline{X}) \| Q(\underline{Y}|\underline{X}))$$

chain Rule for Relative entropy

# Jensen's Inequality and its Consequences
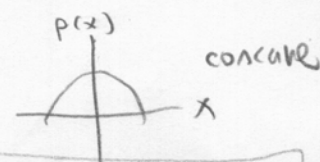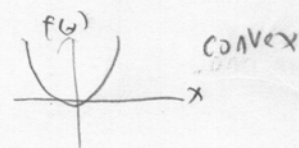
"the inequality about convex functions"

Definition:

- A function $f(x)$ is <u>convex</u> in the interval $x \in [a,b]$
  if for any $(x_1, x_2) \in [a,b]$ and for any $\lambda \in [0,1]$

$$f(\lambda x_1 + (1-\lambda) x_2) \leq \lambda f(x_1) + (1-\lambda) f(x_2)$$

$\rightarrow$ $f(x)$ is <u>strictly convex</u> if equality here
implies $\lambda = 0$ or $1$.

Note:   If $\dfrac{\partial^2 f}{\partial x^2} \geq 0$ over $x \in [a,b]$

then $f(x)$ is convex over $x \in [a,b]$



convex

concave

---

Jensen's Inequality:

If $f(\underline{X})$ is a convex function of random variable $\underline{X}$

then     $\langle f(\underline{X}) \rangle_{p(x)} \geq f(\langle \underline{X} \rangle_{p(x)})$

---

ex/     $\langle x^2 \rangle \geq \langle x \rangle^2$   because $f(x) = x^2$ is a convex function

$\langle -\log(x) \rangle \leq -\log(\langle x \rangle)$     because $f(x) = -\log(x)$ is
a concave function

Consequence:

for any two distributions $P(\underline{X})$ and $Q(\underline{X})$

$$\boxed{D(P(\underline{X}) \| Q(\underline{X})) \geq 0}$$

Information inequality

Proof:

$$-D(P(\underline{X}) \| Q(\underline{X})) = -\sum_{x_i \in \underline{X}} P(x_i) \log\left(\frac{P(x_i)}{Q(x_i)}\right) \qquad \text{concave}$$

$$= \sum_{x_i \in \underline{X}} P(x_i) \log\left(\frac{Q(x_i)}{P(x_i)}\right)$$

with Jensen's inequality

$$\leq \log\left(\sum_{x_i \in \underline{X}} P(x_i) \frac{Q(x_i)}{P(x_i)}\right) = \log\left(\sum_{x_i \in \underline{X}} Q(x_i)\right)$$

$$= \log(1) = 0$$

$$-D(P(\underline{X}) \| Q(\underline{X})) \leq 0$$

so $\underline{D(P(\underline{X}) \| Q(\underline{X})) \geq 0}$

Similarly...

$$\boxed{D(P(\underline{Y}|\underline{X}) \| Q(\underline{Y}|\underline{X})) \geq 0}$$

# Consequences of Information Inequality

$$D(P(\underline{X}) \| Q(\underline{X})) \geq 0 \quad \rightarrow \quad H(\underline{X}) \leq \log(N)$$

since

$$\log(N) + H(\underline{X}) = D(P(\underline{X}) \| u(\underline{X}))$$

$$: u(\underline{X}) = \text{uniform distribution}$$

$$\downarrow$$

$$H(\underline{X} : \underline{Y}) \geq 0$$

$$\rightarrow H(\underline{X} | \underline{Y}) \leq H(\underline{X})$$

$$\rightarrow H(\underline{X}, \underline{Y}) \leq H(\underline{X}) + H(\underline{Y})$$

$$D(P(\underline{Y} | \underline{X}) \| Q(\underline{Y} | \underline{X})) \geq 0 \quad \rightarrow H(\underline{Y} | \underline{X}) \leq \log(N)$$

$$\downarrow$$

$$H(\underline{X} : \underline{Y} | \underline{Z}) \geq 0$$

$$\rightarrow H(\underline{X} | \underline{Y} \underline{Z}) \leq H(\underline{X}, \underline{Y})$$

$$\rightarrow H(\underline{X}, \underline{Y} | \underline{Z}) \leq H(\underline{X} | \underline{Z}) + H(\underline{Y} | \underline{Z})$$

Other useful

## Chain Rules

$$H(\underline{W}\,\underline{X}\,\underline{Y}\,\underline{Z}) = H(\underline{W}) + H(\underline{X} | \underline{W}) + H(\underline{Y} | \underline{W}\,\underline{X}) + H(\underline{Z} | \underline{W}\,\underline{X}\,\underline{Y})$$

$$H(\underline{W} : \underline{X}\,\underline{Y}\,\underline{Z}) = H(\underline{W} : \underline{X}) + H(\underline{W} : \underline{Y} | \underline{X}) + H(\underline{W} : \underline{Z} | \underline{X}\,\underline{Y})$$