

Lecture 9

A Brief Intro to Information Theory

- Motivation



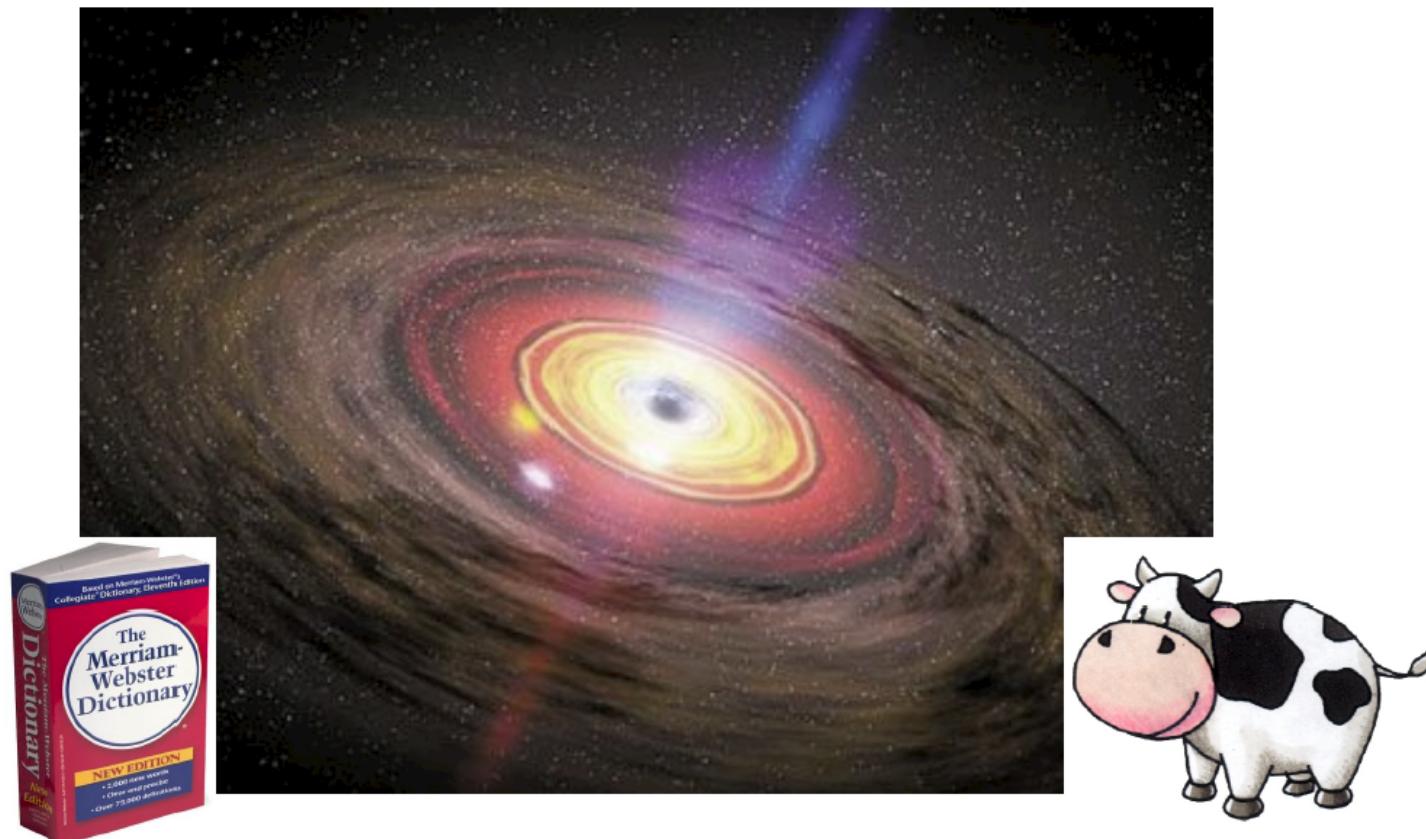
Contents

- Motivation
- Entropy and Mutual Information

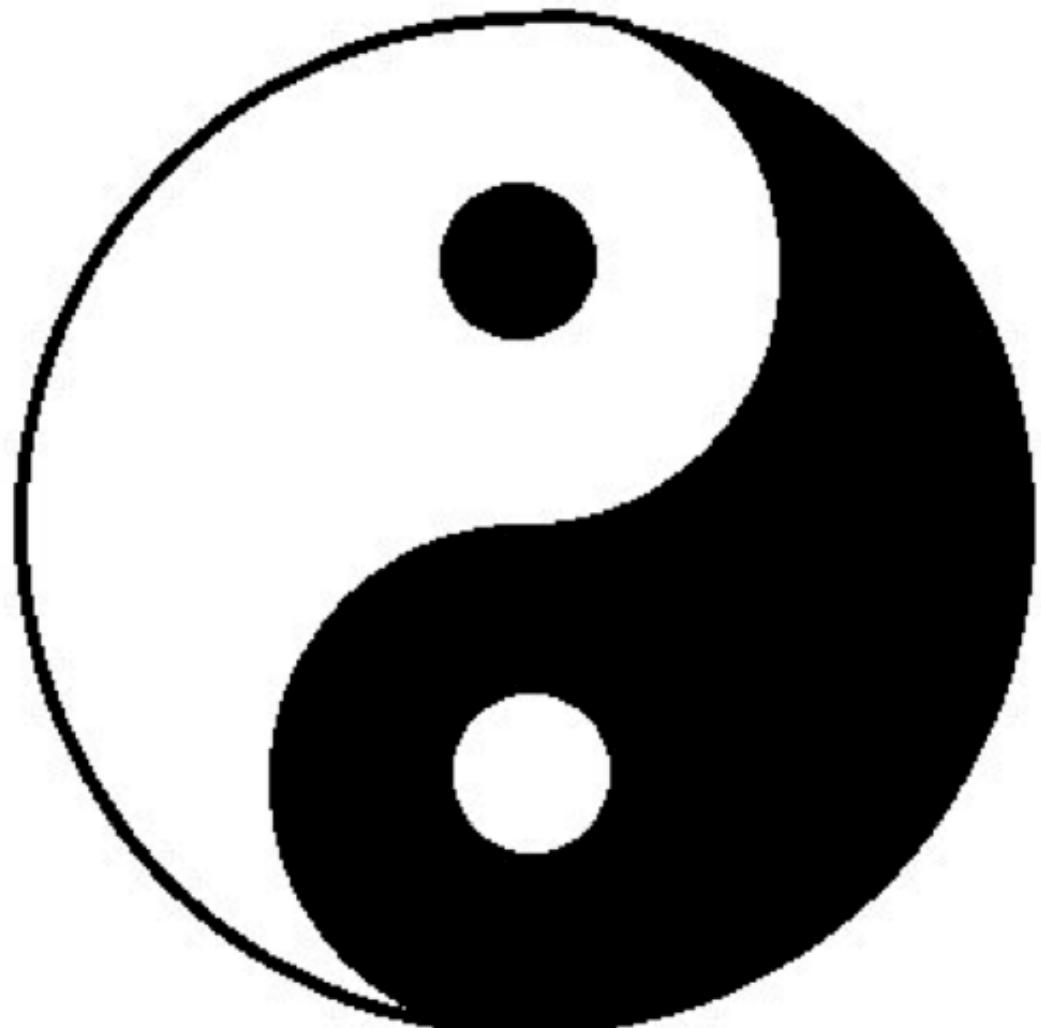
A Thought Experiment

Throw a cow or a dictionary into a black hole,
which has higher information loss?

- Tom Cover



How to quantify information?



Small information content

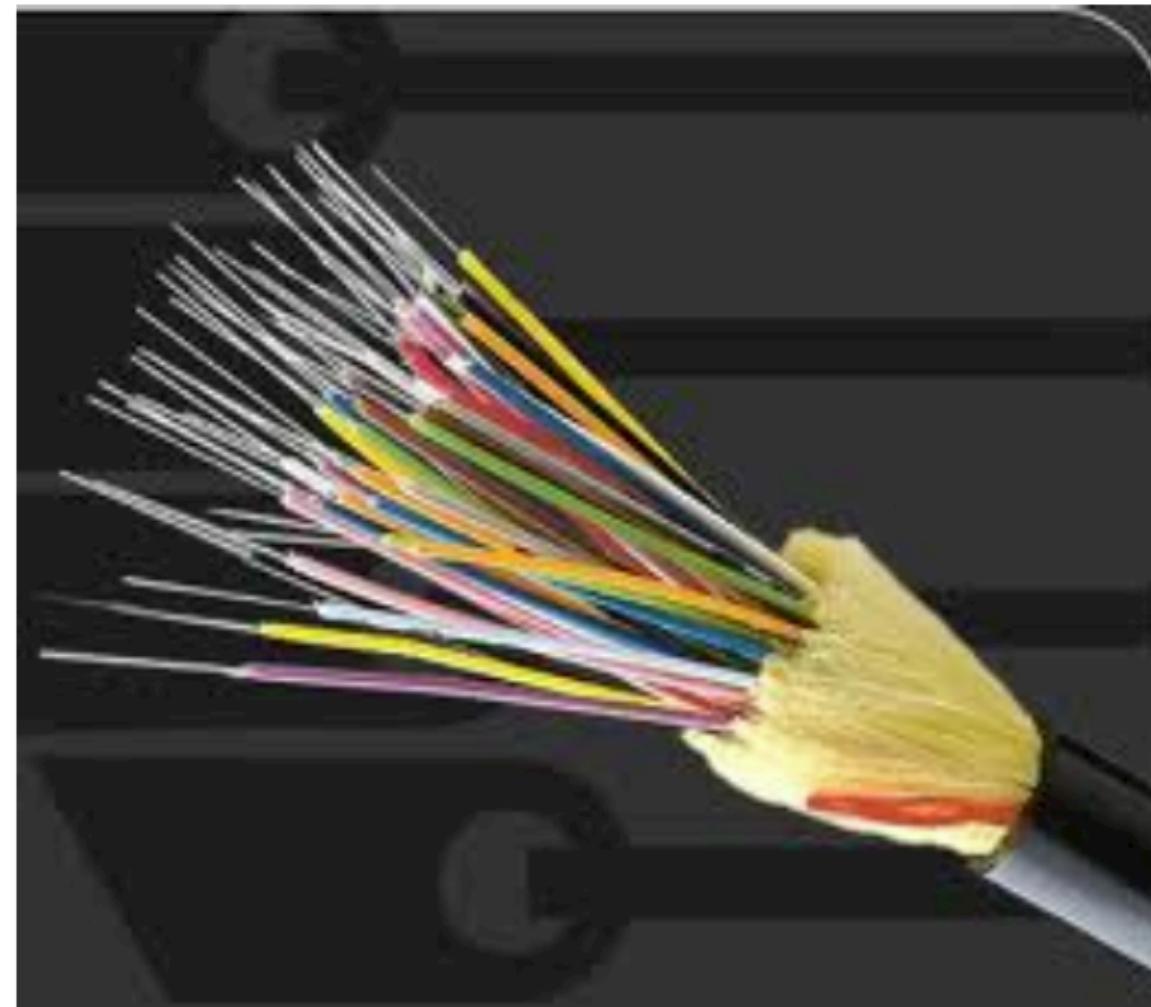


Large information content

What is the fundamental limit of data transfer rate?



WiFi: data rate \sim Mbit/s

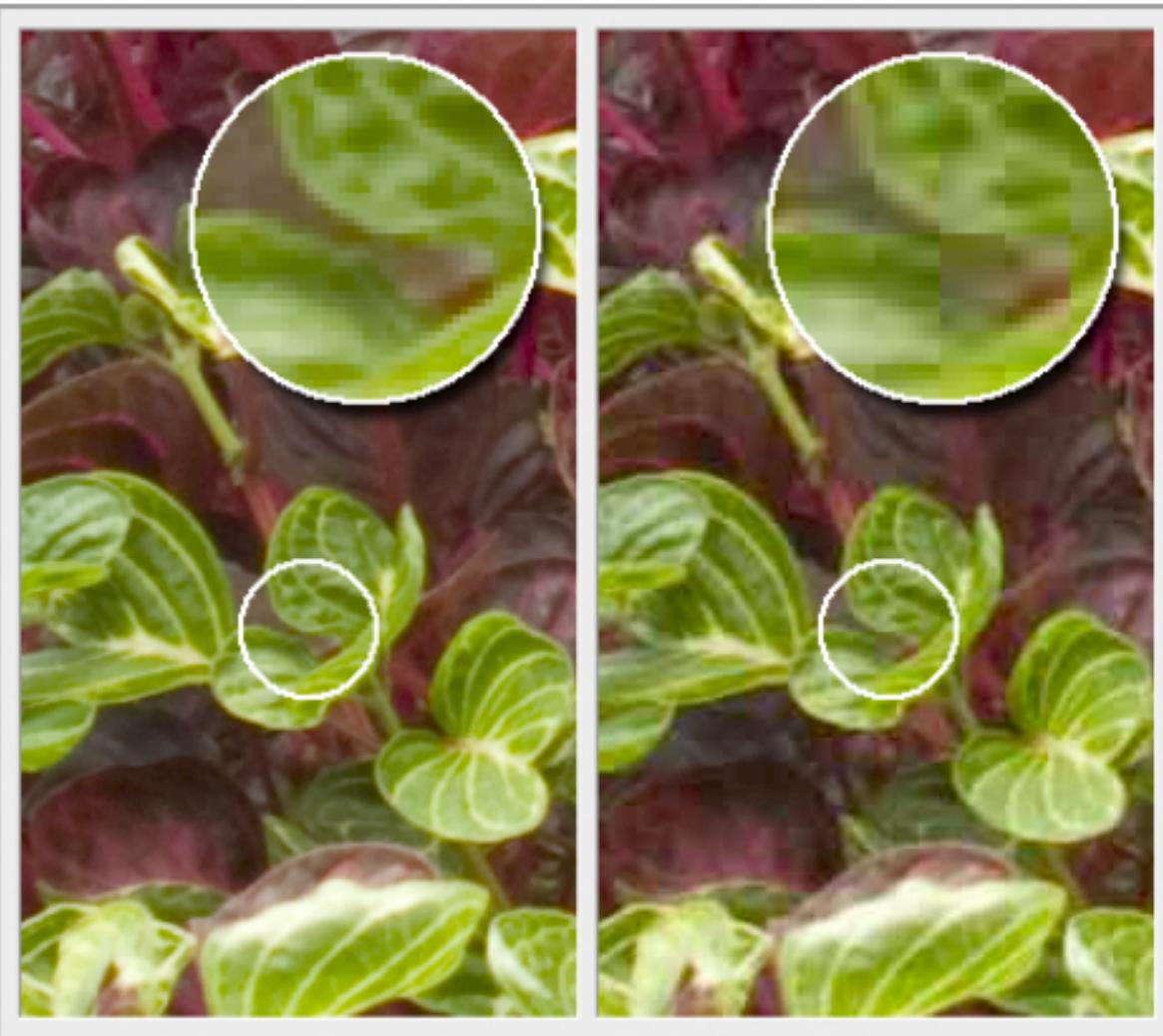


Fiber Optics: data rate \sim Tbit/s

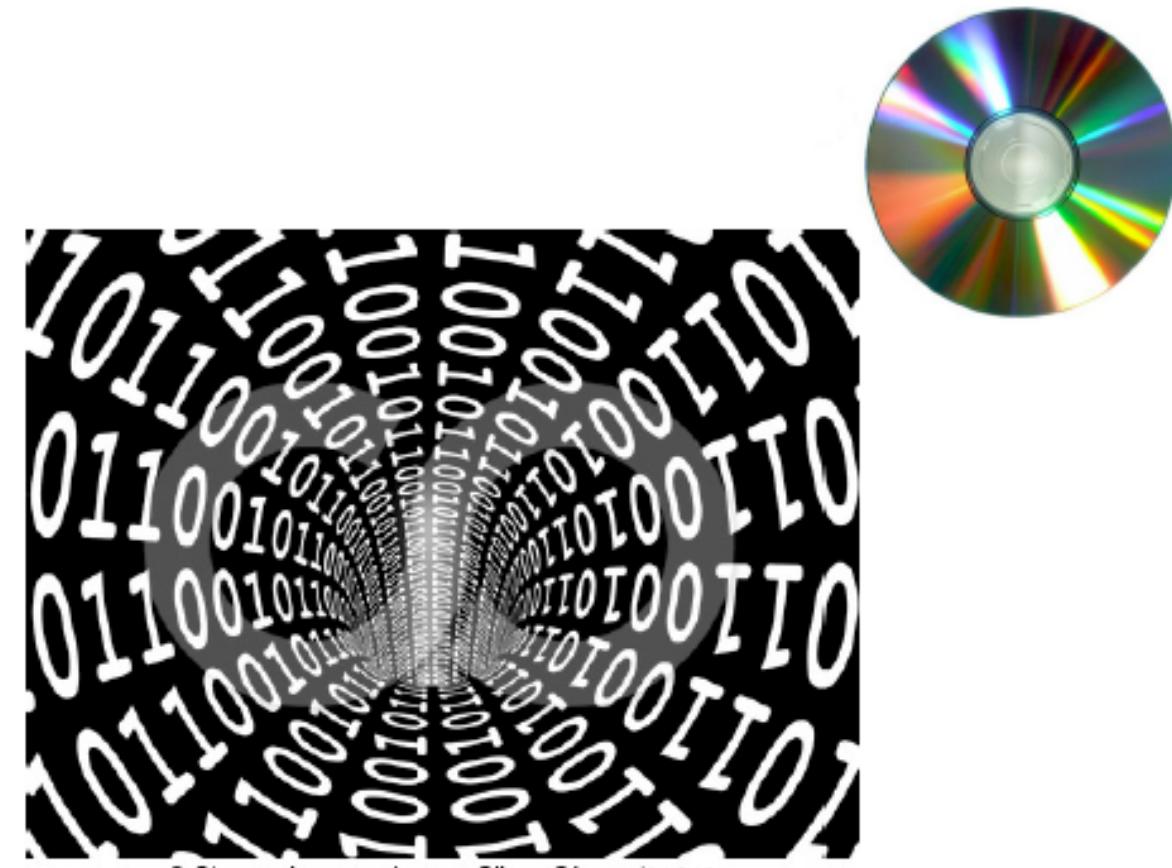
Some people think information theory (IT) is about...



But IT is also about these...

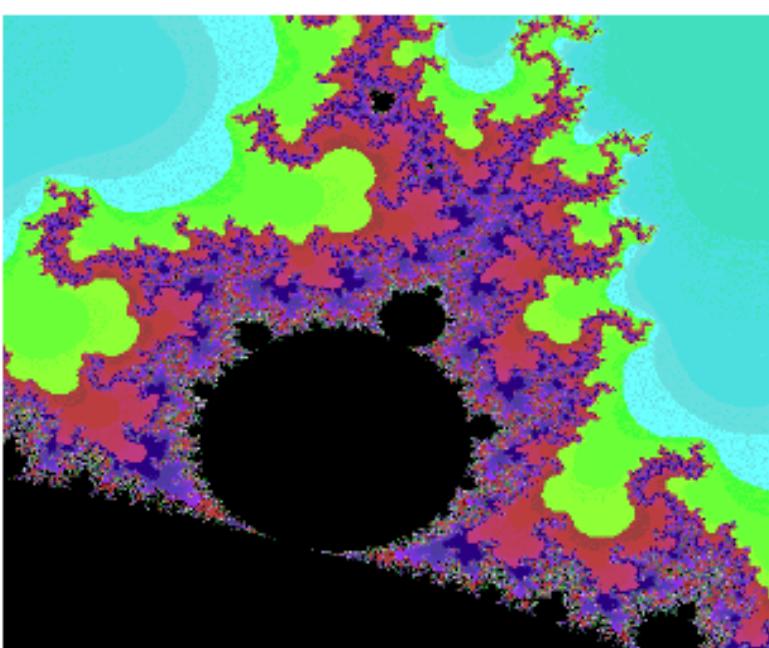


Data Compression

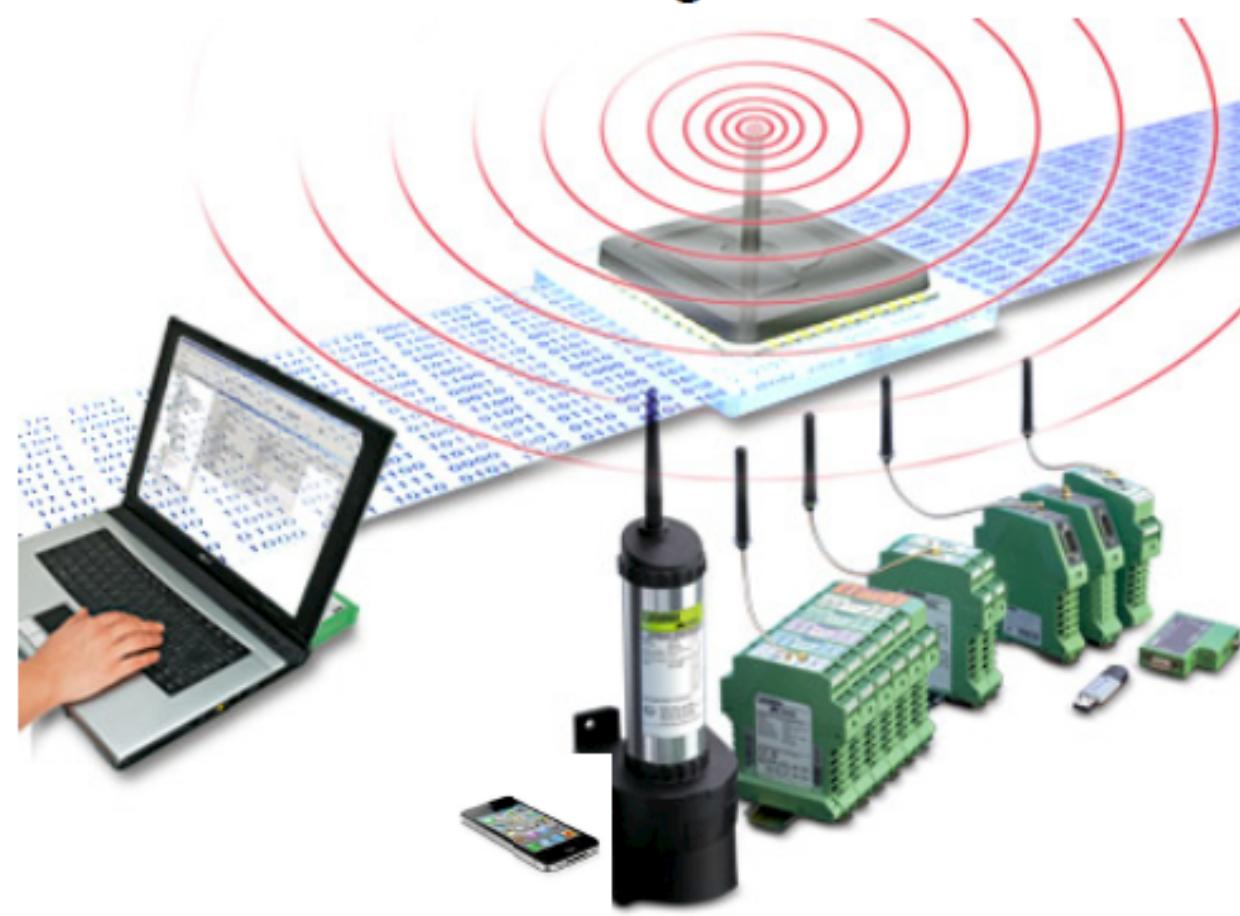


© ShazamImages * www.ClipartOf.com/61816

Coding

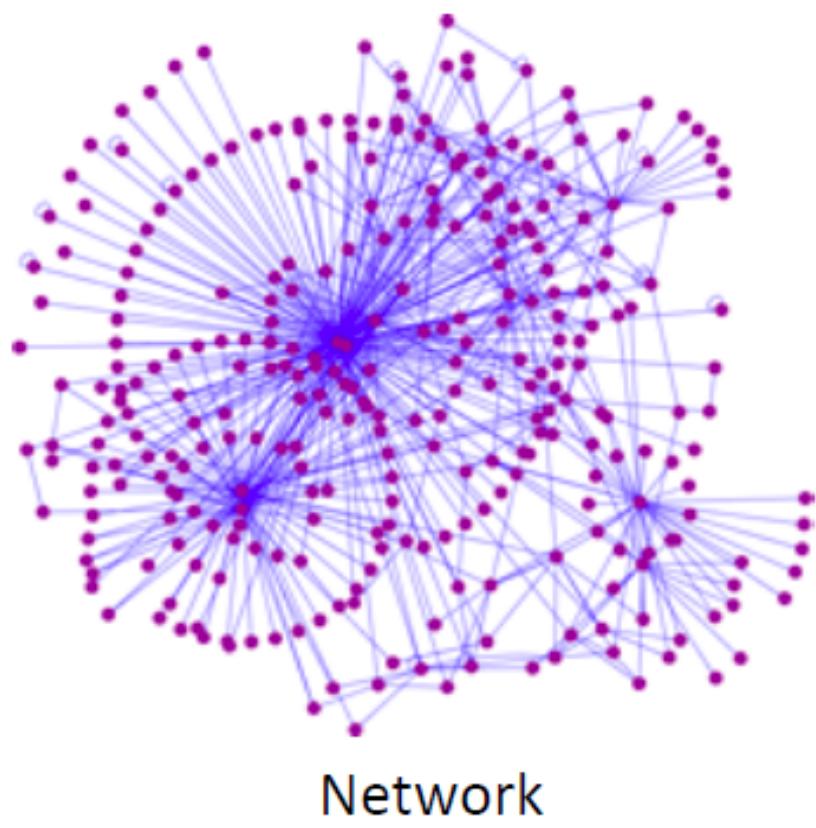


Computation: Kolmogorov Complexity



Data Communication

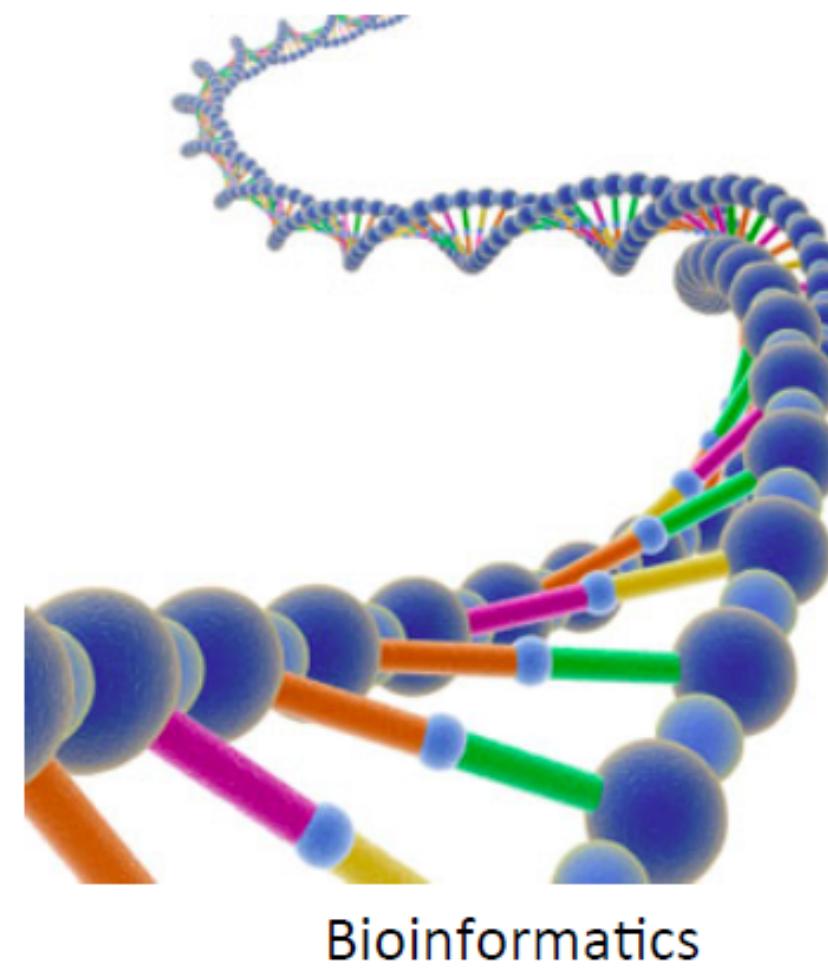
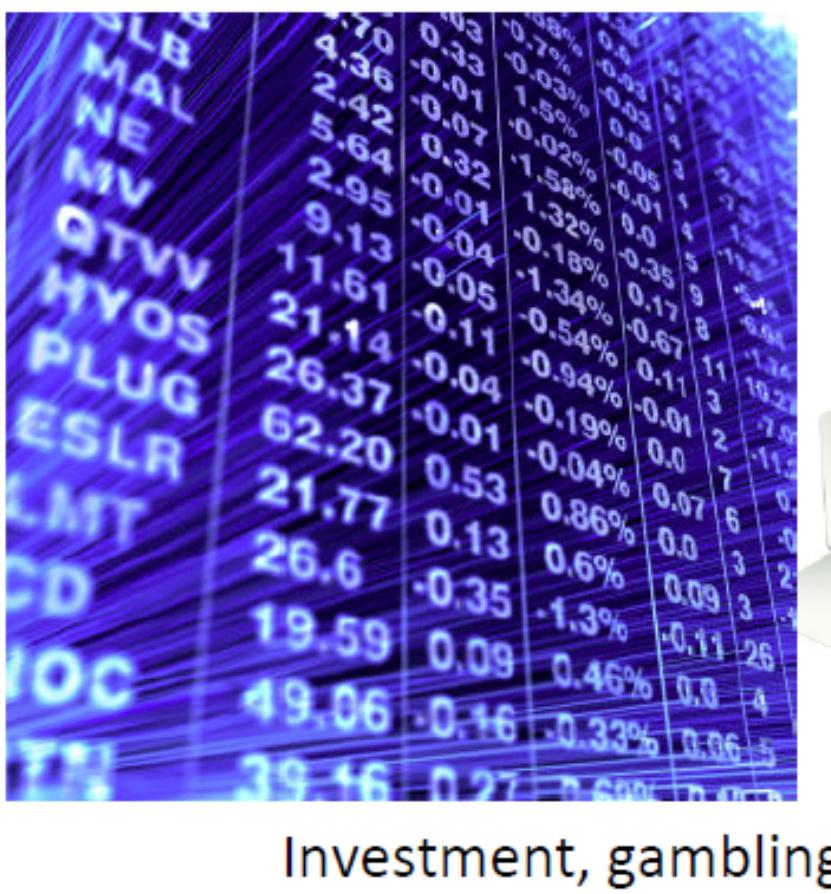
And even these...



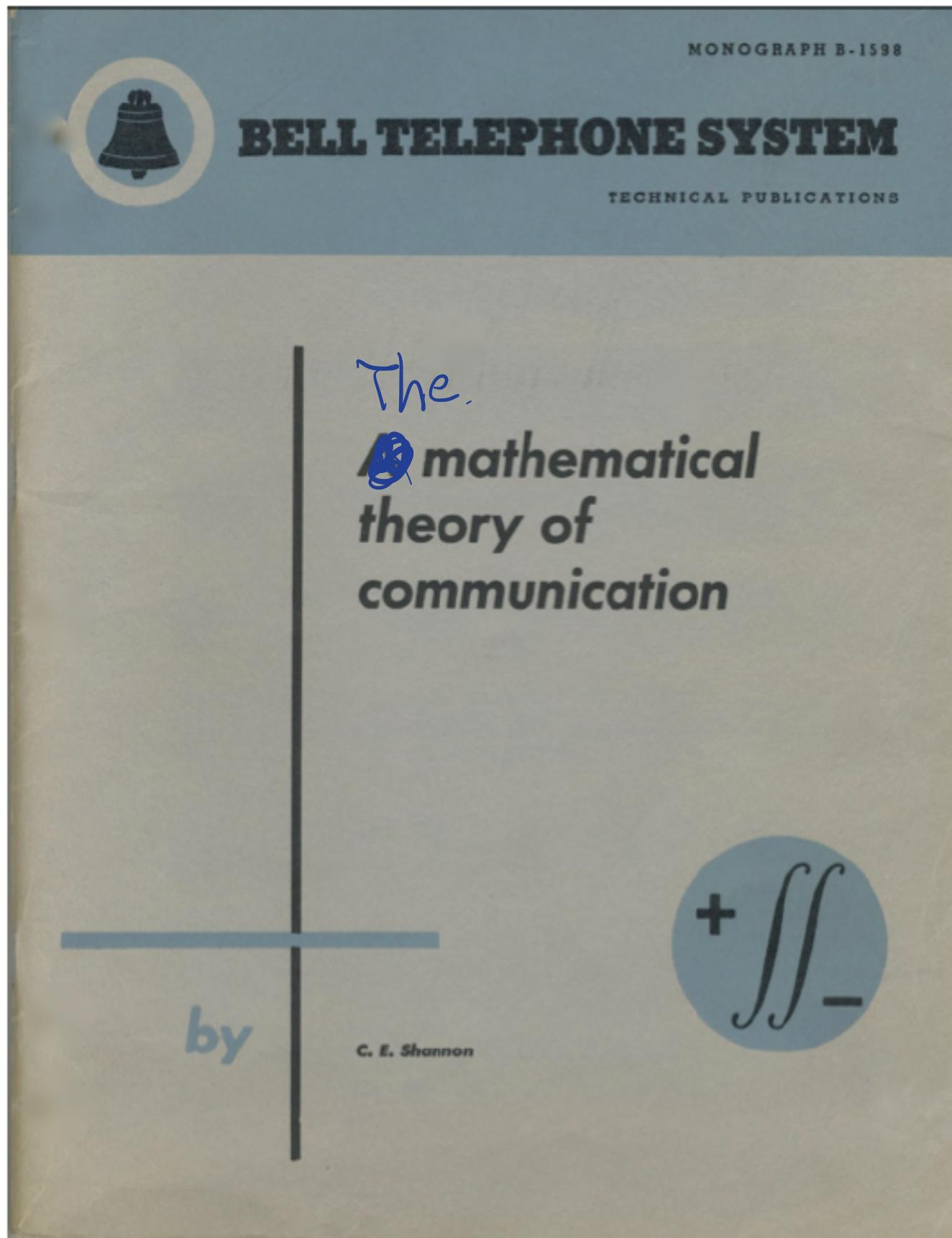
$$y = \Phi x$$

$M \times 1 \quad M \times N \ (M < N) \quad N \times 1$

(Compressed) Sensing



Where IT all begins...



1948, Bell Sys. Tech. Journal



Shannon, 1916 - 2001

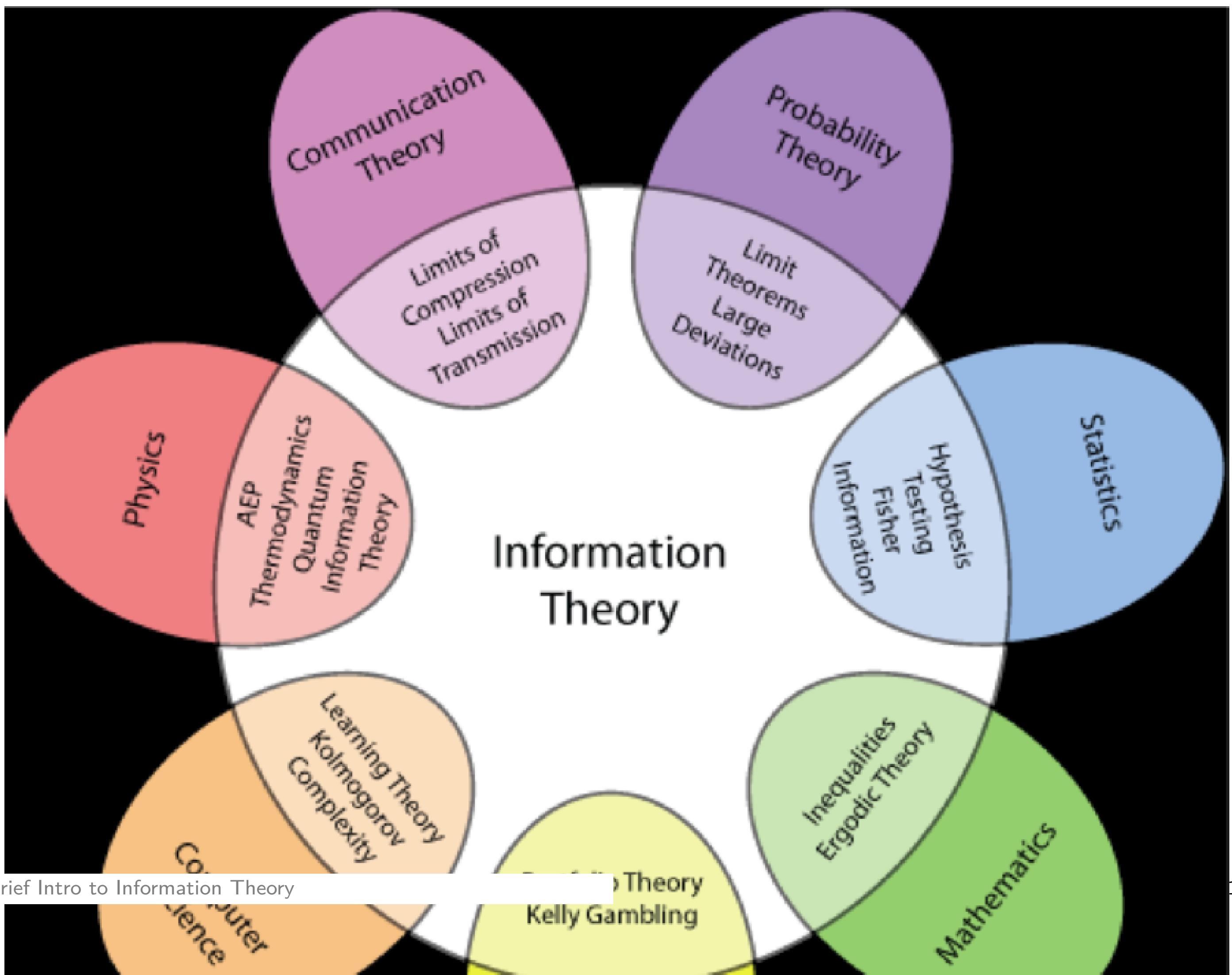
Information Theory

- Shannon's information theory deals with limits on data compression (source coding) and reliable data transmission (channel coding)
 - How much can data be compressed? Entropy
 - How fast can data be reliably transmitted over a noisy channel?
- Two basic "point-to-point" communication theorems (Shannon 1948)
 - **Source coding theorem:** the minimum rate at which data can be *compressed losslessly* is the *entropy rate* of the source
 - **Channel coding theorem:** The maximum rate at which data can be *reliably transmitted* is the *channel capacity* of the channel

Extensions and Applications

- Since Shannon's 1948 paper, many extensions
 - Rate distortion theory
 - Source coding and channel capacity for more complex sources
 - Capacity for more complex channels (multiuser networks)
- Information theory was considered (by most) an esoteric theory with no apparent relation to the "real world"
- Recently, advances in technology (algorithms, hardware, software)
today there are practical schemes for
 - data compression
 - transmission and modulation
 - error correcting coding
 - compressed sensing techniques
 - information security ...

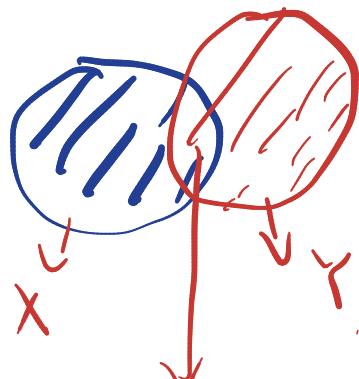
IT encompasses many fields



In this class we will cover the basics

• Nuts and Bolts

- Entropy: uncertainty of a single random variable



$$H(X) = - \sum_x p(x) \log_2 p(x) \text{ (bits)}$$

$X: \text{Support } \{1\}$

$$\Pr(X=1) = 1$$

$X: \text{Support } \{0, 1\}$

$$\Pr(X=0) = \Pr(X=1) = \frac{1}{2}$$

- Conditional Entropy: $H(X|Y)$

Mutual information: reduction in uncertainty due to another random variable

$$H(X) = - \Pr(X=0) \log_2 \Pr(X=0)$$

$$- \Pr(X=1) \log_2 \Pr(X=1)$$

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

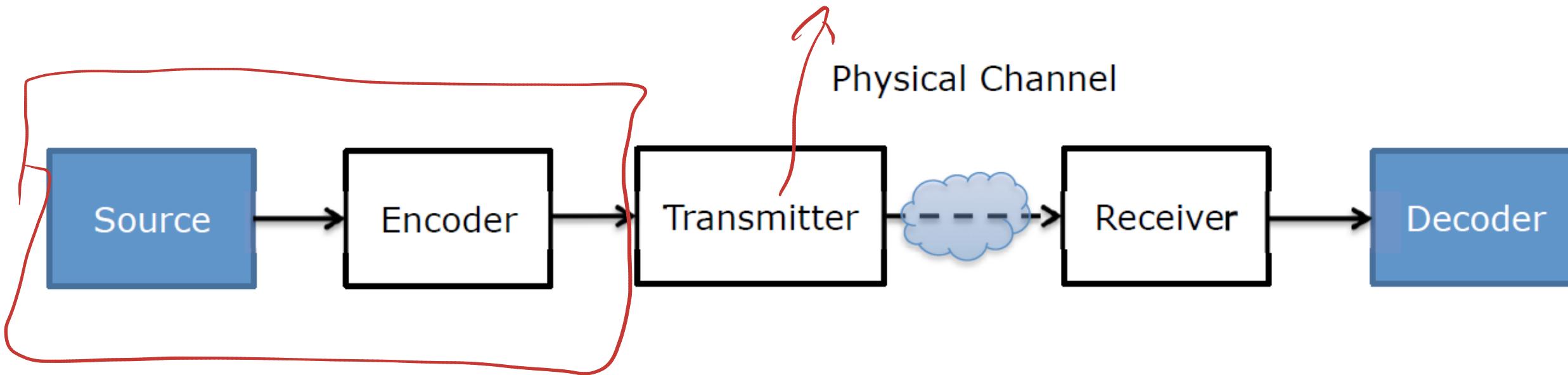
$$= - \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

- Channel capacity $C = \max_{p(x)} I(X; Y)$

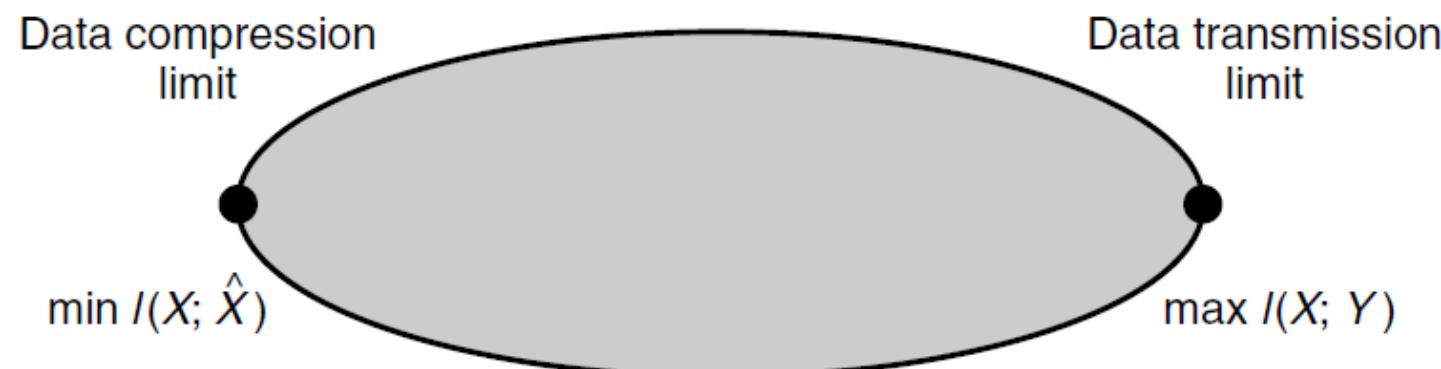
$$= - \log_2 \frac{1}{2} = 1$$

- Relative entropy: $D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$

Fundamental Limits



- Data compression limit (lossless source coding)
- Data transmission limit (channel capacity)
- Tradeoff between rate and distortion (lossy compression)



Important Functionals

- Upper case X, Y, \dots refer to random variables
- \mathcal{X}, \mathcal{Y} alphabet of random variables $\backslash \text{mathcal}\{\mathcal{X}\}$
 $\backslash \text{mathcal}\{\mathcal{Y}\}$
- $p(x) = P(X = x)$
- $p(x, y) = P(X = x, Y = y)$
- Probability density function $f(x)$

Expectation and Variance

- **Expectation:** $\mu = \mathbb{E}\{X\} = \sum xp(x)$
- Why is this of particular interest? It appears in Law of Large Number (LLN): If x_n independent and identically distributed,

$$\frac{1}{N} \sum_{n=1}^N x_n \rightarrow \mathbb{E}\{X\}, \text{ w.p.1}$$

- **Variance:** $\sigma^2 = \mathbb{E}\{(X - \mu)^2\} = \mathbb{E}\{X^2\} - \mu^2$
- Why is this of particular interest? It appears in Central Limit Theorem (CLT):

$$\frac{1}{\sqrt{N\sigma^2}} \sum_{n=1}^N (x_n - \mu) \rightarrow \mathcal{N}(0, 1)$$

Information theory: is it all about theory?

Yes and No.

Yes, it's theory

- Yes, it's theory. We will see many proofs. But it's also in preparation for other subjects
 - Coding theory (Prof. R. Calderbank)
 - Wireless communications
 - Compressed sensing
 - Stochastic network
 - Many proof ideas come in handy in other areas of research

No, it's practical too

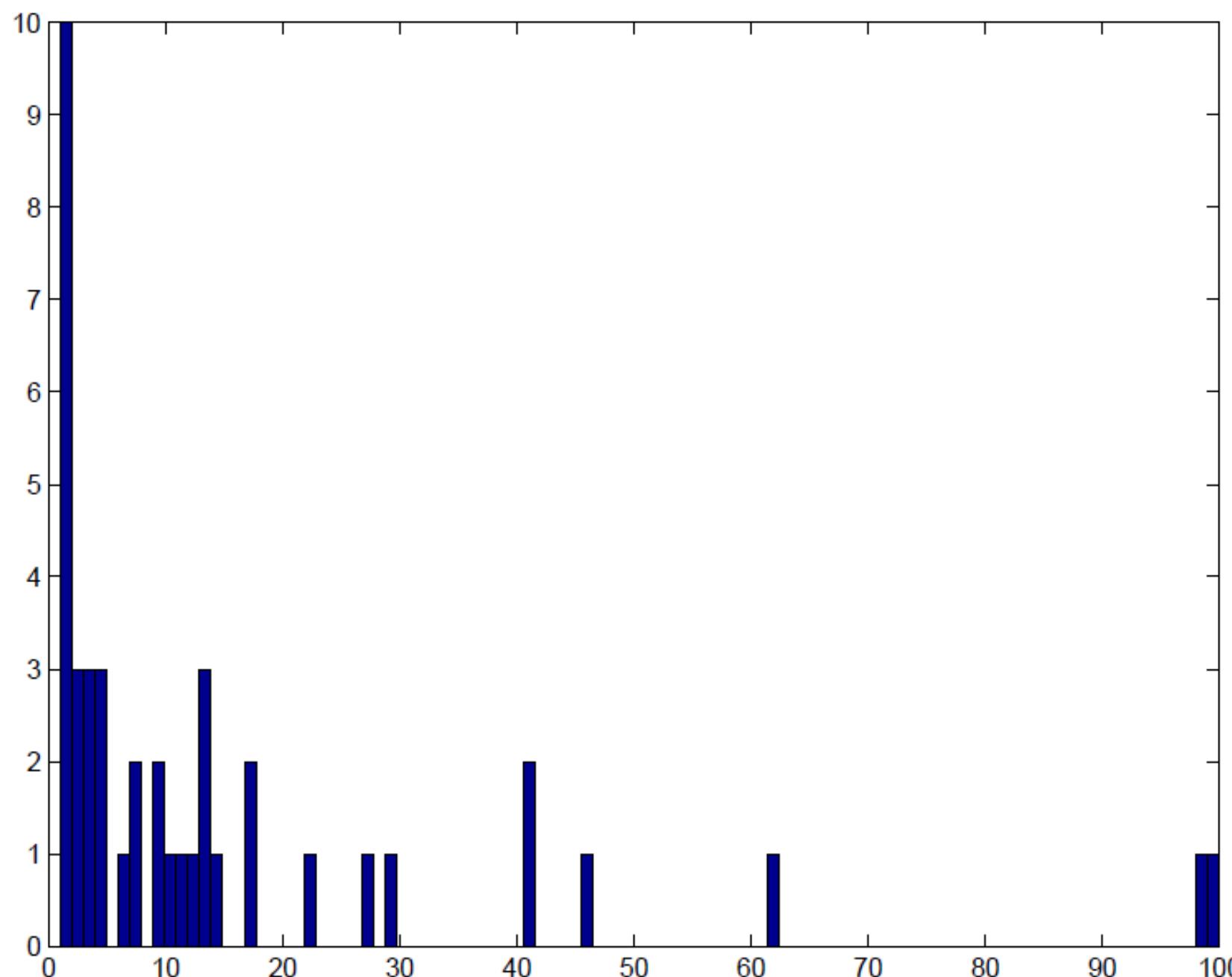
- No. Hopefully you will walk out of this classroom understanding
 - Basic concepts people talk on the streets: entropy, mutual information ...
 - Channel capacity - all wireless guys should know
 - Huffman code (the optimal lossless code)
 - Hamming code (commonly used single error correction code)
 - "Water-filling" - power allocation in all communication systems
 - Rate-distortion function - if you want to talk with data compression guy

Contents

- Motivation
- Entropy and Mutual Information

The winner is:

Eunsu Ryu, with number 6



$$X = \{1, 2, \dots, 8\}$$
$$\frac{1}{8} \dots \frac{1}{8}$$

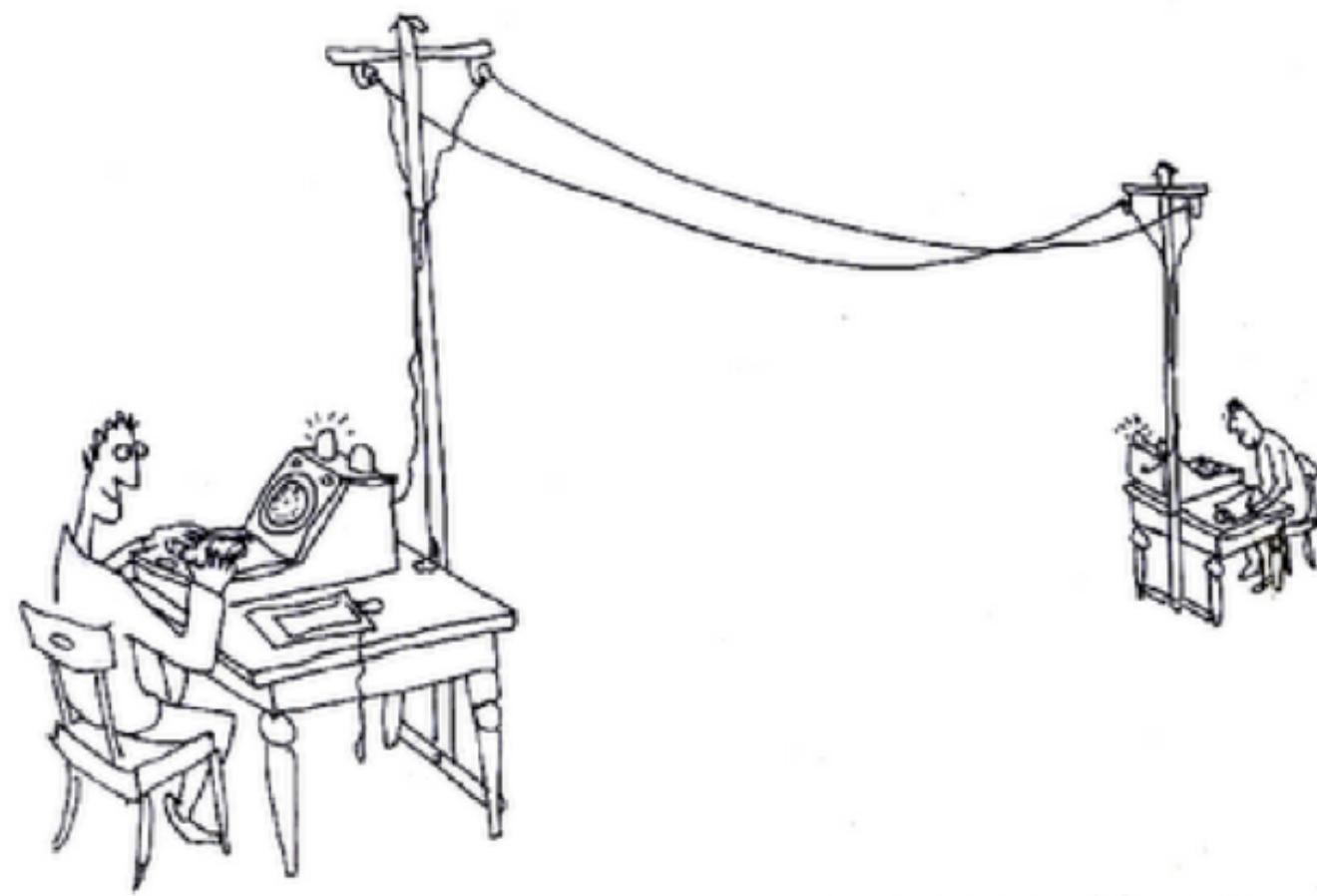
$$E[X] = - \sum_x p(x) \log_2 \frac{1}{8}$$

$$= - \sum_{x=1}^8 \frac{1}{8} \log_2 \frac{1}{8}$$

$$= - \log_2 \frac{1}{8} = 3.$$

A strategy to win the game?

The winner is:



ABC TELEGRAPH EXHIBIT FOR THE NATIONAL ARCHIVES
INVENTIONS EXHIBITION TIM HUNKIN 7/10/05

Which horse won?

Uncertainty measure

- Let X be a random variable taking on a finite number M of different values x_1, \dots, x_M
- What is X : English letter in a file, last digit of Dow-Jones index, result of coin tossing, password
- With probability p_1, \dots, p_M , $p_i > 0$, $\sum_{i=1}^M p_i = 1$
- Question: what is the uncertainty associated with X ?
- Intuitively: a few properties that an uncertainty measure should satisfy
- It should not depend on the way we choose to label the alphabet

Desired properties

- It is a function of p_1, \dots, p_M
- Let this uncertainty measure be

$$H(p_1, \dots, p_M)$$

- **Monotonicity.** Let $f(M) = H(1/M, \dots, 1/M)$. If $M < M'$, then

$$f(M) < f(M')$$

- Picking one person randomly from the classroom should result less possibility than picking a person randomly from the US.

Desired properties (continued)

- **Additivity.** Two independent RV X and Y , each uniformly distributed, alphabet size M and L . The uncertainty for the pair (X, Y) , is ML . However, due to independence, when X is revealed, the uncertainty in Y should not be affected. This means

$$f(ML) - f(M) = f(L)$$

- **Grouping rule** (Problem 2.27 in Text). Dividing the outcomes into two, randomly choose one group, and then randomly pick an element from one group, does not change the number of possible outcomes.

Entropy

- The only function that satisfies the requirements is the entropy function

$$H(p_1, \dots, p_M) = - \sum_{i=1}^M p_i \log_2 p_i$$

- General definition of entropy

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x) \text{ bits}$$

- $0 \log 0 = 0$

Understanding Entropy

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

- Uncertainty in a single random variable

- Can also be written as:

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

- Intuition: $H = \log(\# \text{of outcomes/states})$

- Entropy is a functional of $p(x)$

- Entropy is a lower bound on the number of bits need to represent a RV. E.g.: a RV that has uniform distribution over 32 outcomes

$$= \mathbb{E}_X [-\log_2 p(x)]$$

$$= \mathbb{E}_X [\log_2 \frac{1}{p(x)}]$$

$$\begin{aligned} X & \{1, 2, 3, \dots, 8\} \\ 1 & \rightarrow 001 \quad 4 \rightarrow 100 \\ 2 & \rightarrow 010 \quad 5 \rightarrow 101 \\ 3 & \rightarrow 011 \quad 6 \rightarrow 110 \\ 7 & \rightarrow 111 \\ 8 & \rightarrow 000 \end{aligned}$$

Properties of entropy

- $H(X) \geq 0$

$$H(x) = E_x \left[\log \frac{1}{P(x)} \right] \quad \log \frac{1}{P(x)} \geq 0$$

- Definition, for Bernoulli random variable, $X = 1$ w.p. p ,

$$X = 0 \text{ w.p. } 1 - p$$

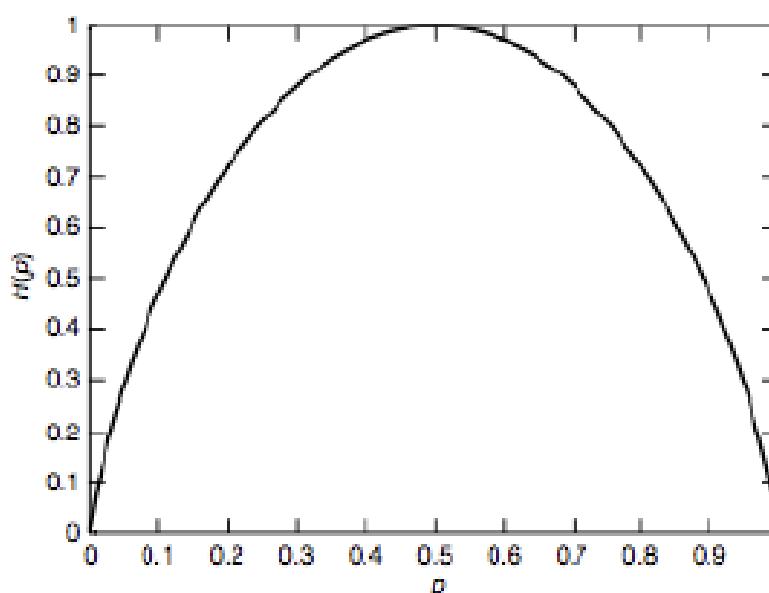
$$H'(p) = -(1-p) \log(1-p) - \left[(1-p) \cdot \frac{1}{1-p} \cdot \gamma \right]$$

$$H(p) = -p \log p - (1-p) \log(1-p)$$

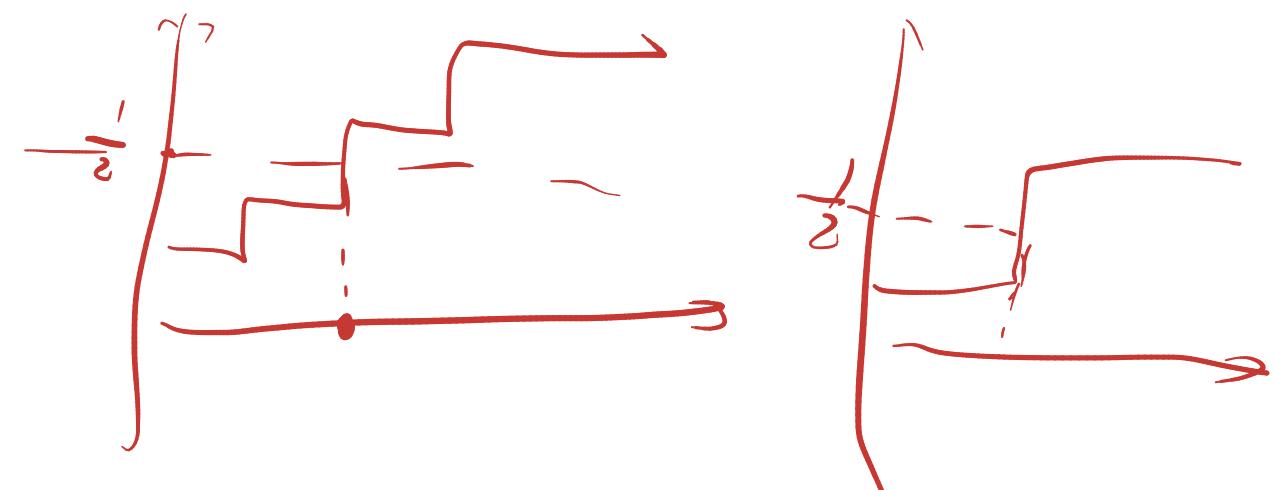
$$- \log(1-p)$$

$$= -1 - \log p + 1 + \log(1-p)$$

$$H''(p) = -\frac{1}{p} - \frac{1}{1-p} \leq 0$$



- **Concave**
- Maximizes at $p = 1/2$
- Example: how to ask questions?



Joint entropy

$$H(X) = -E \log p(x)$$

$$H(Y) = -E \log p(y)$$

$$H(X, Y) = -E \log p(x, y)$$

- Extend the notion to a pair of discrete RVs $(X, Y) \Rightarrow H(X, Y) = -E \log p(x, y)$
- Nothing new: can be considered as a single vector-valued RV
- Useful to measure dependence of two random variables

✓ $H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \underbrace{p(x, y)}_{\text{underbrace}} \log p(x, y)$

✓ $H(X, Y) = -E \log p(X, Y)$

Conditional Entropy

$$(X, Y) \sim p(x, y)$$

$$Y|X \sim p(Y|X)$$

$$Y|X=x \sim p(Y|X=x) = \frac{p(x, y)}{p(x)}, \forall y \in \text{support}(Y)$$

- Conditional entropy: entropy of a RV given another RV. If

$$(X, Y) \sim \underbrace{p(x, y)}$$

$$\boxed{H(Y|X)} = \sum_{x \in \mathcal{X}} p(x) H(Y|X=x)$$

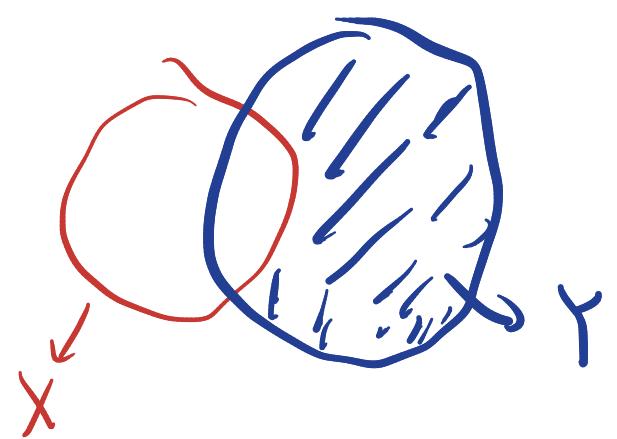
- Various ways of writing this

$$H(Y|X=x) = - \sum_{y \in Y} p(y|x) \log p(y|x)$$

$$H(Y|X) = - \sum_{x \in X} \sum_{y \in Y} p(x) p(y|x) \log p(y|x)$$

$$= - \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x, y) \log p(y|x)$$

$$= - \mathbb{E} \log p(Y|X)$$



Chain rule for entropy

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

- Entropy of a pair of RVs = entropy of one + conditional entropy of the other:

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

- Proof:

$$H(X,Y) = - \sum_{(x,y)} p(x,y) \underbrace{\log p(x,y)}$$

- $H(Y|X) \neq H(X|Y)$

$$= - \sum_{(x,y)} p(x,y) \log p(y|x)$$

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

$$\begin{aligned} &= - \sum_{(x,y)} p(x,y) \log p(x) + \sum_{(x,y)} p(x,y) \log p(y|x) \\ &\quad \nearrow H(Y|X) \\ &- \sum_x \left(\sum_y p(x,y) \right) \log p(x) = - \sum_x p(x) \log p(x) = H(X) \end{aligned}$$

1. If X and Y are independent, then

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



$$= - \sum_x \sum_y p(x,y) \log p(y|x)$$

$$P(y)$$

$$= - \sum_x \sum_y p(x,y) \log P(y)$$

$$= - \sum_y \left(\sum_x p(x,y) \right) \log P(y)$$

$$= - \sum_y p(y) \log P(y)$$

$$= H(Y)$$

2. If Y is deterministic function of X

e.g., $Y = g(X)$

$$\Rightarrow \underline{H(Y|X)} = 0$$

$$- \sum_x \sum_y p(x,y) \log p(y|x)$$

$$p(y|x) = \begin{cases} 0, & \text{if } y \neq g(x) \\ 1, & \text{if } y = g(x) \end{cases}$$

$$- \sum_x \sum_{y:y=g(x)} p(x,y) (\log 1)$$

$$= 0$$

Relative entropy

$$D(P||Q) = \int_{x \in \mathcal{X}} f_P(x) \log \frac{f_P(x)}{f_Q(x)} dx$$

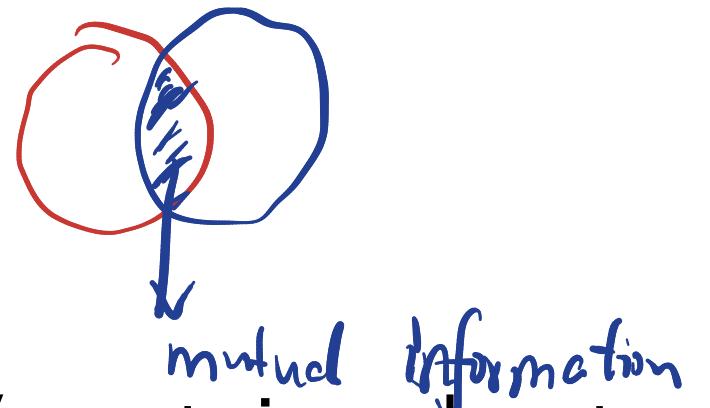
- Measure of distance between two distributions

$$D(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} = \underset{\sim P}{\mathbb{E}_x} \left[\log \frac{p(x)}{q(x)} \right]$$

- Also known as Kullback-Leibler distance in statistics: expected log-likelihood ratio
- A measure of inefficiency of assuming that distribution is q when the true distribution is p
- If we use distribution is q to construct code, we need $H(p) + D(p||q)$ bits on average to describe the RV

Mutual information

(互信息)



- Measure of the amount of information that one RV contains about another RV

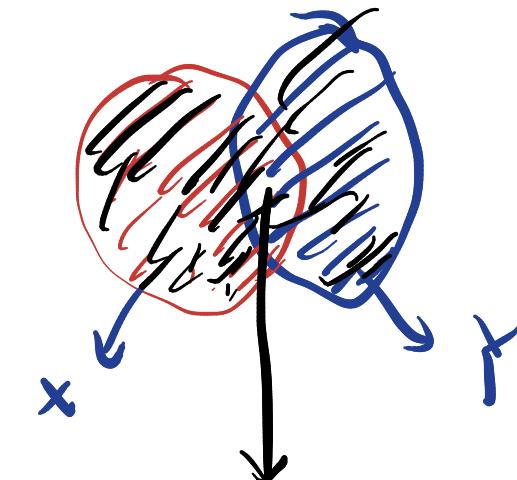
$$I(X;Y) = I(Y;X)$$

$$I(X;Y) = \sum_{x \in X} \sum_{y \in Y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} = D(p(x,y)||p(x)p(y))$$

- Reduction in the uncertainty of one random variable due to the knowledge of the other
- Relationship between entropy and mutual information

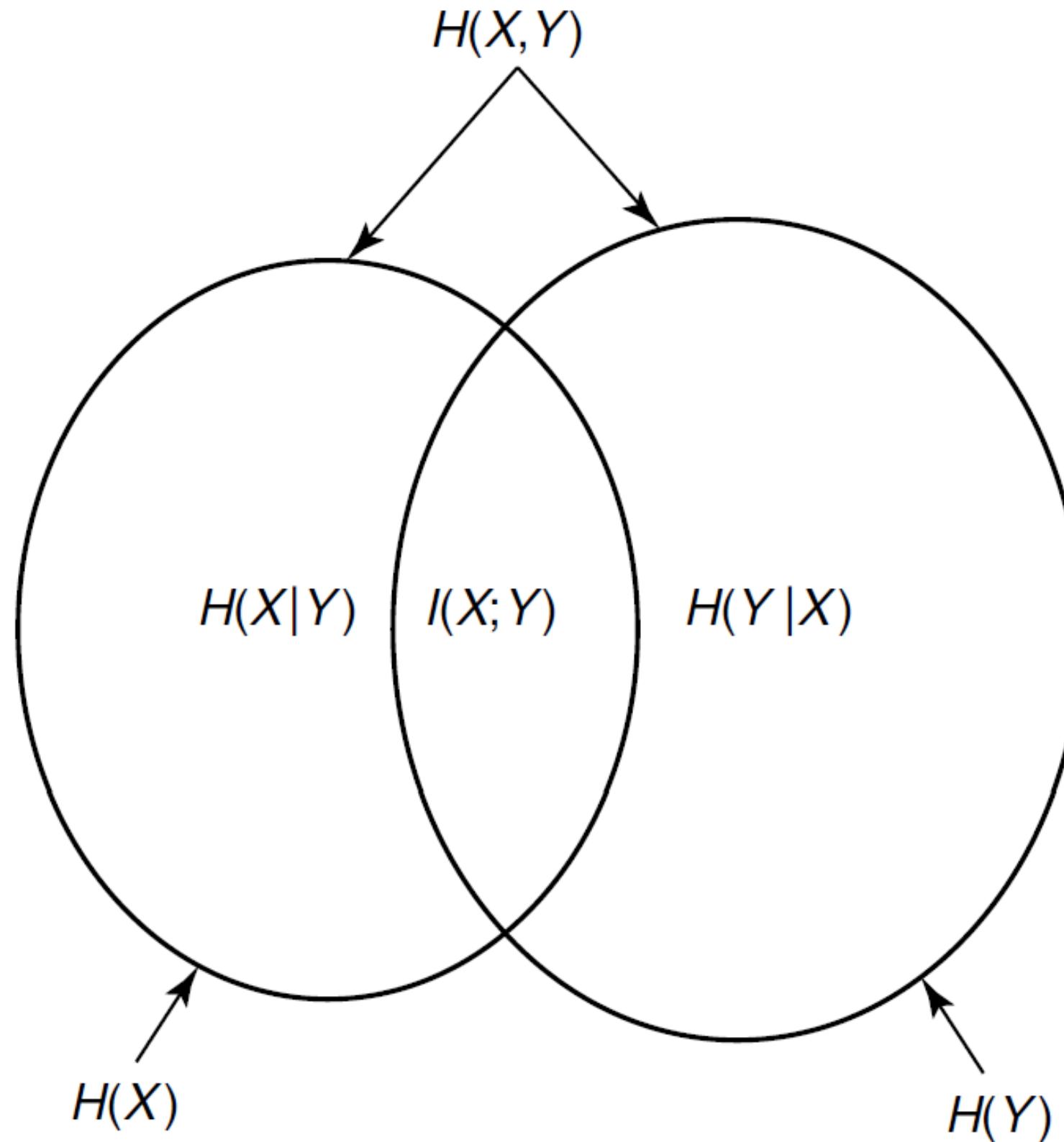
$$\begin{aligned} I(X;Y) &= H(Y) - H(Y|X) \\ &= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} \\ &= \sum_{(x,y)} p(x,y) \log \frac{p(y|x)}{p(y)} = - \sum_{(x,y)} p(x,y) p(y) + \sum_{(x,y)} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} \end{aligned}$$

Mutual information properties



- $I(X; Y) = H(Y) - H(Y|X)$ $I(X; Y)$
- $H(X, Y) = H(X) + H(Y|X) \rightarrow I(X; Y) = \underbrace{H(X)}_{\text{red}} + \underbrace{H(Y)}_{\text{blue}} - \underbrace{H(X, Y)}_{\text{black}}$
- $I(X; X) = H(X) - H(X|X) = H(X)$ Entropy is "self-information"
- Example: calculating mutual information

Venn diagram



$I(X; Y)$ is the intersection of information in X with information in Y

Example: Blood type and skin cancer risk

$$H(X|Y)$$

$$= \sum_y p(y) H(X|Y=y)$$

$$= \sum_{y=1}^4 \frac{1}{4} \cdot H(X|Y=y)$$

Y: chance for
skin cancer

X: blood type

	A	B	AB	O
1 Very Low	1/8	1/16	1/32	1/32
2 Low	1/16	1/8	1/32	1/32
3 Medium	1/16	1/16	1/16	1/16
4 High	1/4	0	0	0

$$H(Y) = -\sum_y p(y) \log p(y)$$

$$= -\sum_y \frac{1}{4} \log \frac{1}{4}$$

$$= -\log \frac{1}{4} = 2$$

- X: marginal $(1/2, 1/4, 1/8, 1/8)$

$$H(X) = -\sum_x p(x) \log p(x)$$

- Y: marginal $(1/4, 1/4, 1/4, 1/4)$

$$= -\frac{1}{2} \cdot \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \underbrace{\frac{1}{8} \log \frac{1}{8}}_{\cancel{-\frac{1}{8} \log \frac{1}{8}}} - \underbrace{\frac{1}{8} \log \frac{1}{8}}_{\cancel{-\frac{1}{8} \log \frac{1}{8}}}$$

- $H(X) = 7/4$ bits $H(Y) = 2$ bits

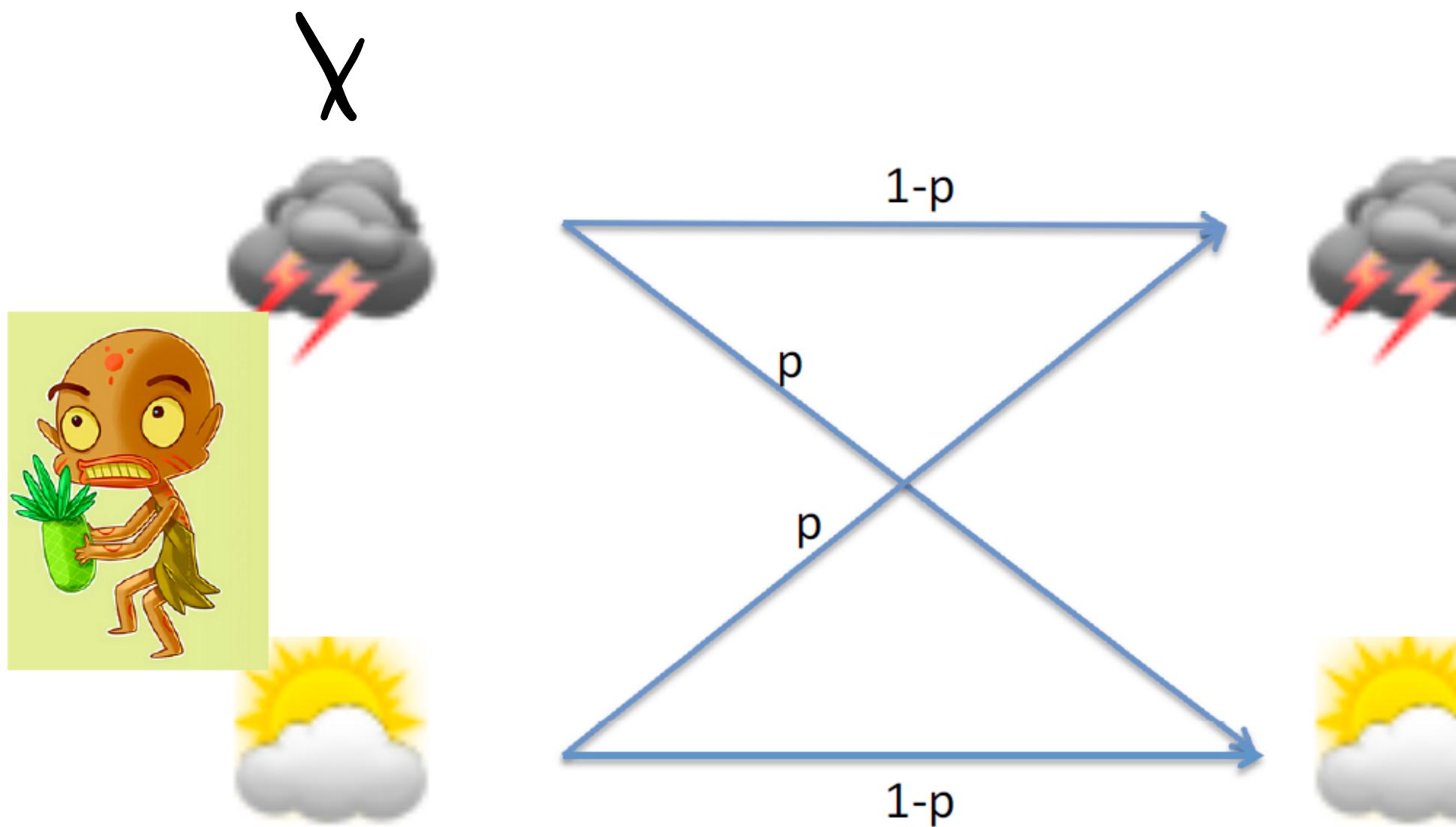
$$= \frac{1}{2} + \frac{1}{2} - \frac{3}{8} - \frac{3}{8} = 1 + \frac{3}{4} = 1.75 \text{ bits}$$

- Conditional entropy: $\underline{H(X|Y)} = 11/8$ bits, $H(Y|X) = 13/8$ bits

- $H(Y|X) \neq H(X|Y)$

- Mutual information: $I(X;Y) = H(X) - H(X|Y) = 0.375$ bit

Example: Binary Symmetric Channel



$$P(X, Y) = \begin{cases} \frac{1}{2}(1-p), & (x,y)=0,0 \\ \frac{1}{2}p, & (x,y)=0,1 \\ \frac{1}{2}p, & (x,y)=1,0 \\ \frac{1}{2}(1-p), & (x,y)=1,1 \end{cases}$$

$$\begin{aligned} & -\frac{1}{2}(1-p) \log \left[\frac{1}{2}(1-p) \right] \\ & -\frac{1}{2}p \log \left[\frac{1}{2}p \right] \\ & -\frac{1}{2}p \log \left[\frac{1}{2}p \right] \\ & -\frac{1}{2}(1-p) \log \left[\frac{1}{2}(1-p) \right] \end{aligned}$$

$$= -(1-p) \left[-1 + \log(1-p) \right]$$

$$- p \left[-1 + \log p \right]$$

$$P(X=0) = P(X=1) = \frac{1}{2}$$

~~$P(Y|X)$~~

$$P(Y|X) = \begin{pmatrix} 1-p & p \\ p & 1-p \end{pmatrix}$$

$$P(Y=0) = P(Y=0|X=0)P(X=0) + P(Y=0|X=1)P(X=1) = \frac{1}{2} = P(Y=1)$$

$$\begin{aligned} \underbrace{I(X; Y)}_{=} & H(X) + H(Y) - H(X, Y) \\ & = 1 + 1 - H(X; Y) \\ & = 1 - H(p) \end{aligned}$$

$$\begin{aligned} & = 1 - (1-p) \log_2(1-p) \\ & - p \log_2 p \end{aligned}$$

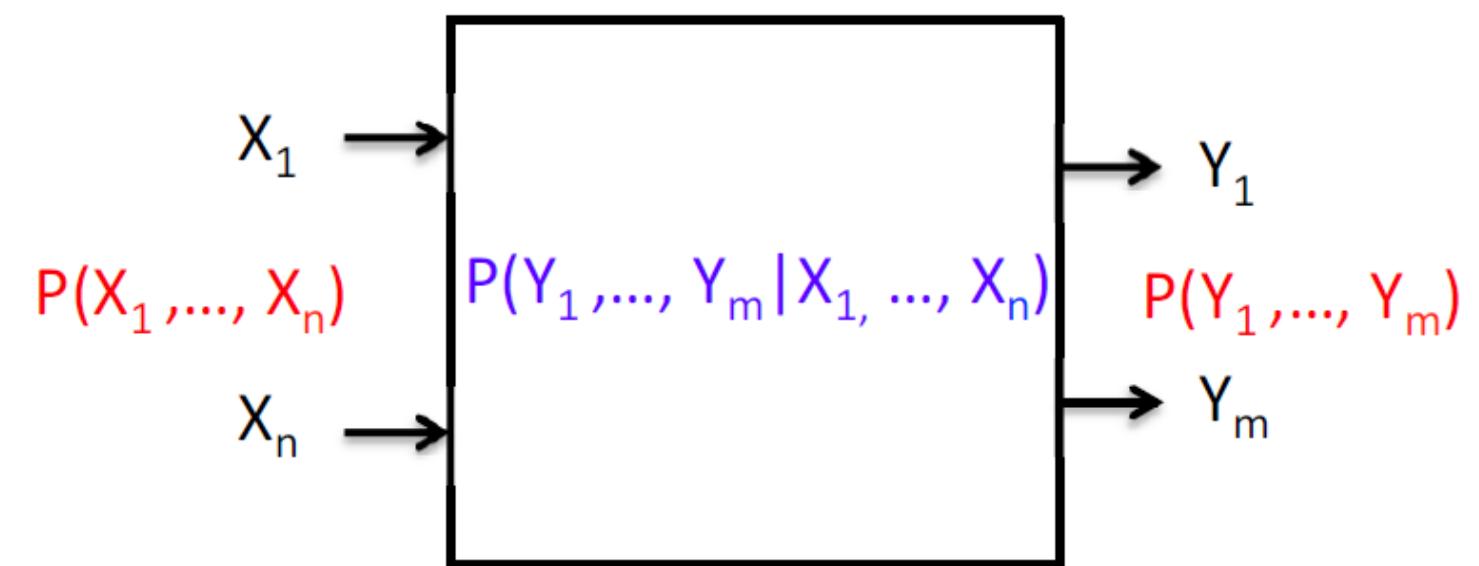
Summary

Entropy



$$H(X)$$

Mutual Information



$$I(X_1, \dots, X_n; Y_1, \dots, Y_m)$$