

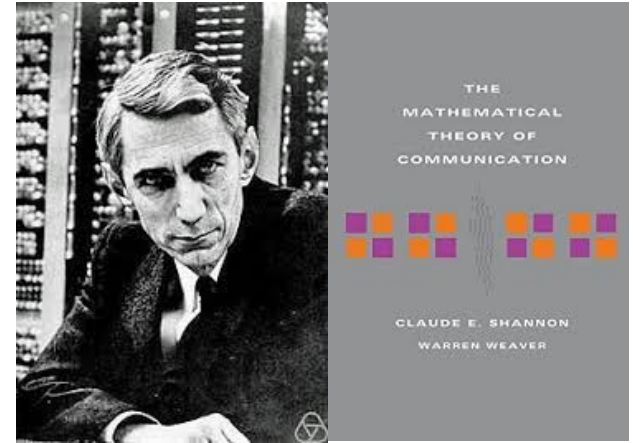
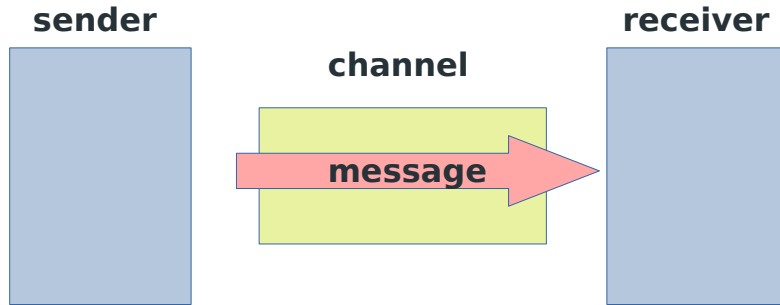
Michele Allegra
Information theory and Inference



Information theory

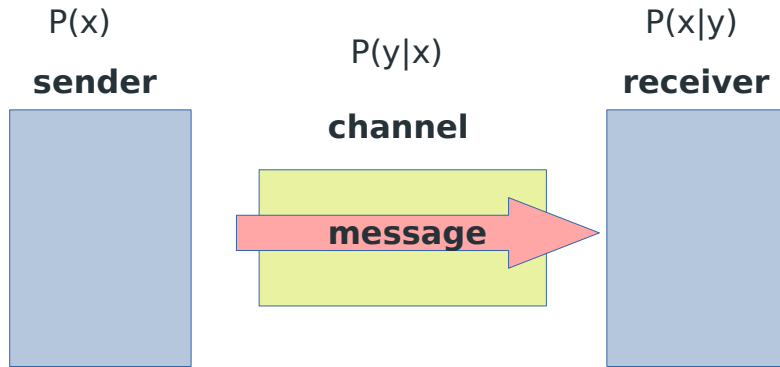
“the fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point”

“the significant aspect is that the actual message is one selected from a set of possible messages. The system must be designed to operate for each possible selection”



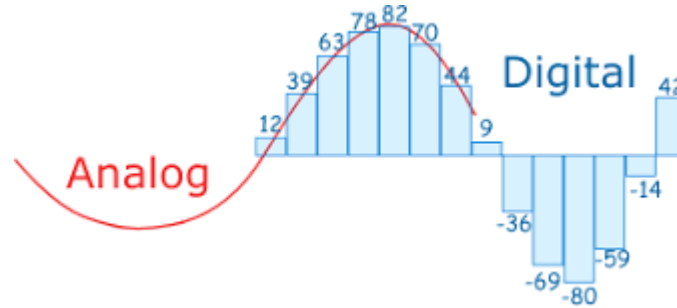
Information theory

- sender selects (input) message x in S with probability $P(x)$
- channel yields (output) y in R with probability $P(y|x)$
- Receiver reconstructs input message $P(x|y) = P(y|x)P(x)/P(y)$
(*signal processing*)

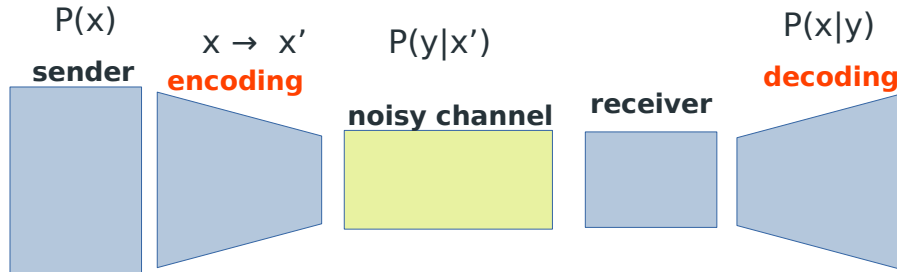


Information theory

encoding $x \rightarrow x'$



*Any type of information converted into digitized information
(string of 0's and 1's)*



Information theory

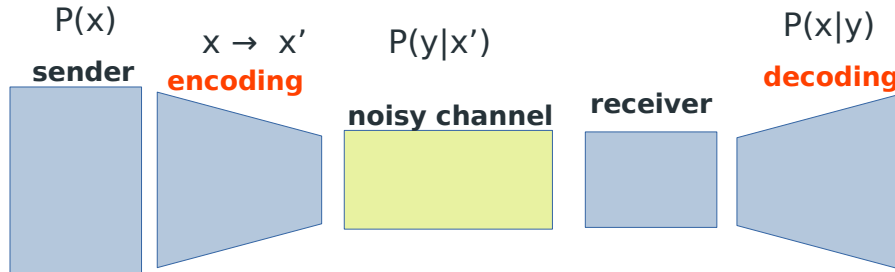
Reliable information transmission limited by mutual information: $I(X:Y)$

$I(X:Y) = \text{Mutual information} = H(X) + H(Y) - H(X,Y)$

$H(X) = \text{Shannon entropy} = -\sum p(x) \log[p(x)]$

$H(X)$ information needed to specify X

$I(X:Y)$ information Y provides about X



Information theory

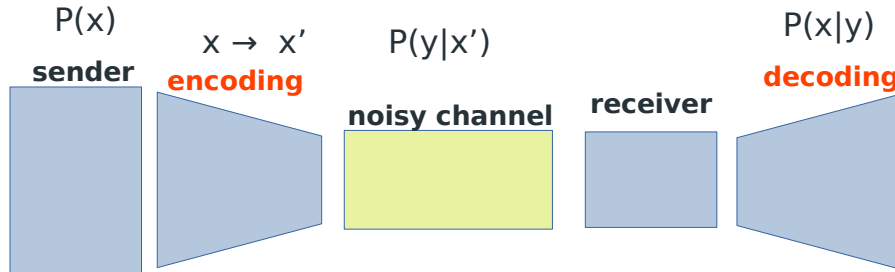
Reliable information transmission limited by mutual information: $I(X:Y)$

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

$H(X)$ information needed to specify X

$H(X|Y)$ information needed to specify X if Y is known

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



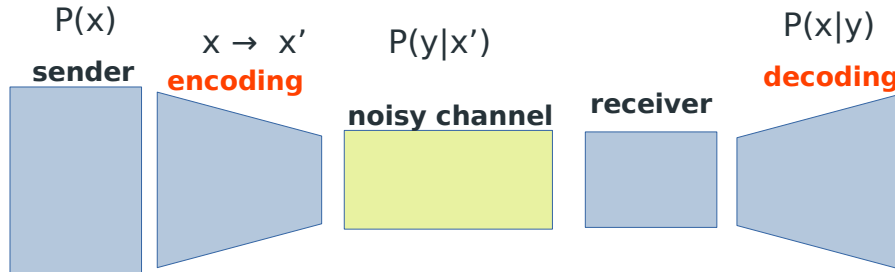
Information theory

Reliable information transmission limited by mutual information: $I(X:Y)$

$$I(X:Y) = D(P(X,Y)||P(X)P(Y))$$

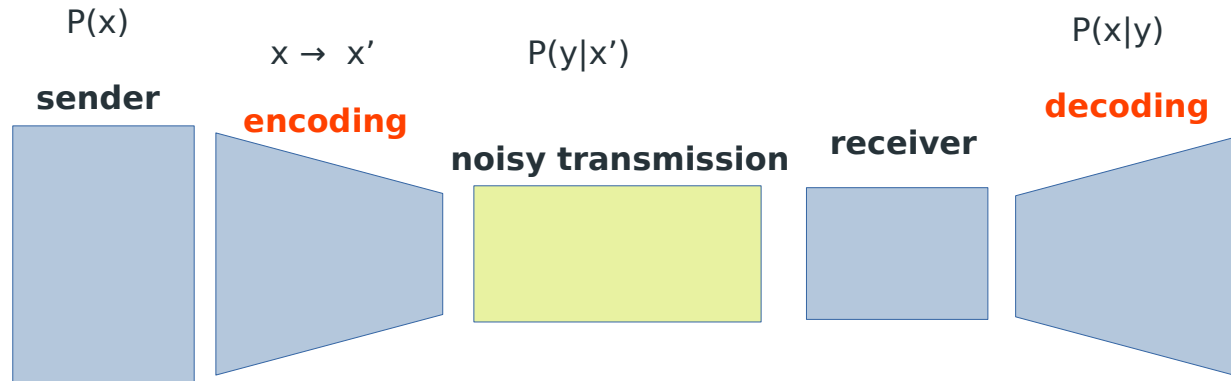
$$D(P||Q) = \text{relative entropy/K-L divergence} = -\sum p(x)\log[p(x)/q(x)]$$

$D(P||Q)$ = “distance between P and Q ”



Information theory

- **Quantify** the information that can be reliably transmitted
- **Design** channel to maximize information transmission rate



Information theory

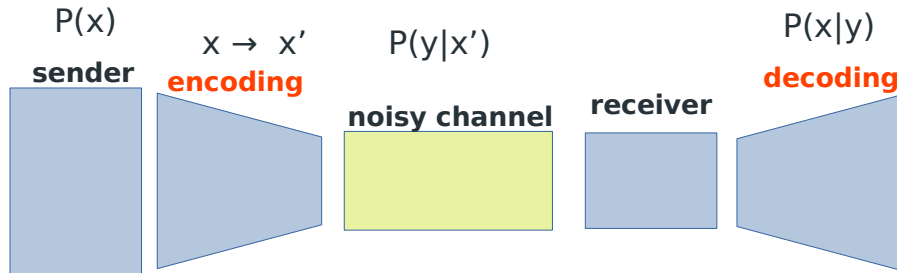
Achieving reliable information transmission with *minimal code length*

design smart encoding/decoding

e.g. send bit (0,1) though noisy channel with bit flip probability of ϵ

repetition code: 0 → 0000 1 → 1111

error probability ϵ^k



Inference

“the uncertainty of knowledge rests on events or their causes.”

“If we are sure an urn contains black and white papers in a given ratio, and we ask the chance that extracting random paper it will be white, the event is uncertain but the cause determining its probability (the white/black ratio), is known.”

“Hence the following problem: if an urn contains black and white papers in an unknown ratio, and we extract a white paper, determine the chance that the ratio is p/q ”

“the probability of each cause is equal to the probability of the event given the cause, divided by the sum of all the probabilities of the event given any cause”



Inference

Again, suppose that p_1, p_2, \dots, p_n are a number of mutually exclusive hypotheses such that one of them must be true.

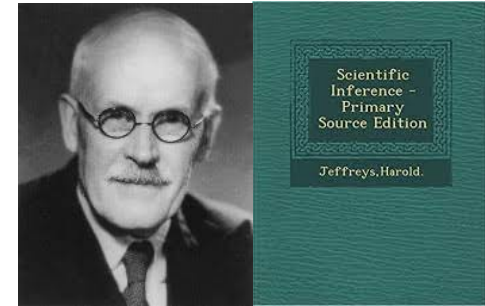
Therefore

$$P(p_r | q \cdot h) = \frac{P(q | p_r \cdot h) P(p_r | h)}{\sum_{r=1}^n P(q | p_r \cdot h) P(p_r | h)}. \quad (4)$$

This theorem* is to the theory of probability what Pythagoras's theorem is to geometry.

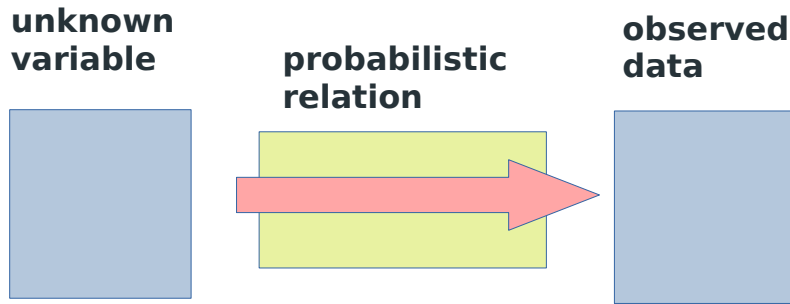
Then our result is that the posterior probability of p is the prior probability of p divided by the prior probability of the consequence.

Harold Jeffreys 1939



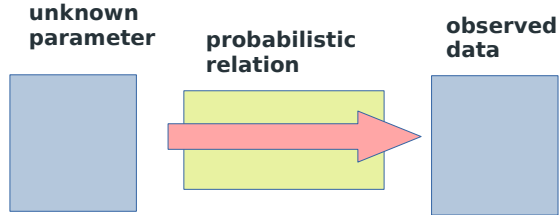
Inference

- unknown variable x in S with probability $P(x)$
- data y in R with probability $P(y|x)$
- reconstruct variable $P(x|y) = P(y|x)P(x)/P(y)$
(*data processing*)

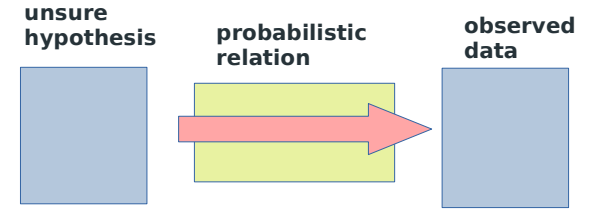


Inference

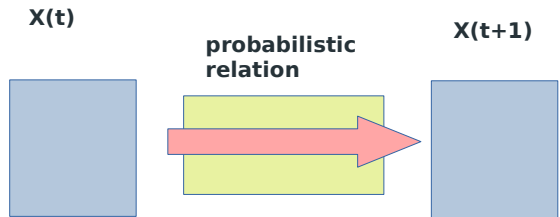
Parameter inference



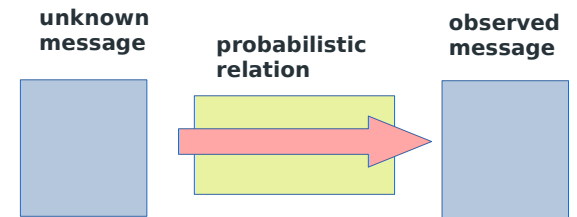
Hypothesis testing



Dynamical system analysis



Communication



Inference and sampling

- Inference requires to compute a posterior $P(x|y) = P(y|x)P(x)/P(y)$

- This allows to obtain *best estimates* and assess *uncertainty* about estimates

$$x^{est} = E[x|y]$$

$$\delta x^{est} = E[(x - x^{est})^2 | y]$$

- This involves *complicated high-dimensional sums/integrals*

- e.g. $E[x|y] = \int dx x P(y|x)P(x)/P(y) = \int dx x P(y|x)P(x) / \int dx P(y|x)P(x)$

sampling is needed!!

Inference and sampling

- Ising spins $\mathbf{s} = \{s_1, \dots, s_n\} \in \{-1, 1\}^N$
- observe spins:

$$P(\mathbf{s}|\mathbf{h}, J) = \frac{e^{-\mathbf{h} \cdot \mathbf{s} - \mathbf{s}^T J \mathbf{s}}}{Z(\mathbf{h}, J)}$$

- reconstruct \mathbf{h}, J :

$$P(\mathbf{h}, J | \mathbf{s}^{obs}) = \frac{P(\mathbf{s}^{obs} | \mathbf{h}, J)}{P(\mathbf{s}^{obs})} = \frac{P(\mathbf{s}^{obs} | \mathbf{h}, J)}{\int d\mathbf{h} dJ P(\mathbf{s}^{obs} | \mathbf{h}, J)}$$

$$\langle \mathbf{h} \rangle = \frac{\int d\mathbf{h} dJ \mathbf{h} P(\mathbf{s}^{obs} | \mathbf{h}, J)}{\int d\mathbf{h} dJ P(\mathbf{s}^{obs} | \mathbf{h}, J)}, \quad \langle J \rangle = \frac{\int d\mathbf{h} dJ J P(\mathbf{s}^{obs} | \mathbf{h}, J)}{\int d\mathbf{h} dJ P(\mathbf{s}^{obs} | \mathbf{h}, J)}$$

Inference and sampling

Computers were invented for numerical integration!

THE JOURNAL OF CHEMICAL PHYSICS VOLUME 21, NUMBER 6 JUNE, 1953

Equation of State Calculations by Fast Computing Machines

NICHOLAS METROPOLIS, ARIANNA W. ROSENBLUTH, MARSHALL N. ROSENBLUTH, AND AUGUSTA H. TELLER,
Los Alamos Scientific Laboratory, Los Alamos, New Mexico

AND

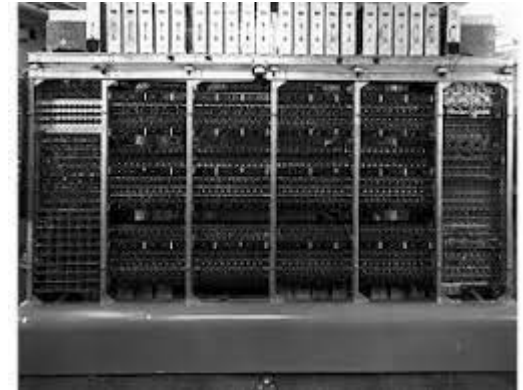
EDWARD TELLER, * *Department of Physics, University of Chicago, Chicago, Illinois*

(Received March 6, 1953)

A general method, suitable for fast computing machines, for investigating such properties as equations of state for substances consisting of interacting individual molecules is described. The method consists of a modified Monte Carlo integration over configuration space. Results for the two-dimensional rigid-sphere system have been obtained on the Los Alamos MANIAC and are presented here. These results are compared to the free volume equation of state and to a four-term virial coefficient expansion.



Metropolis



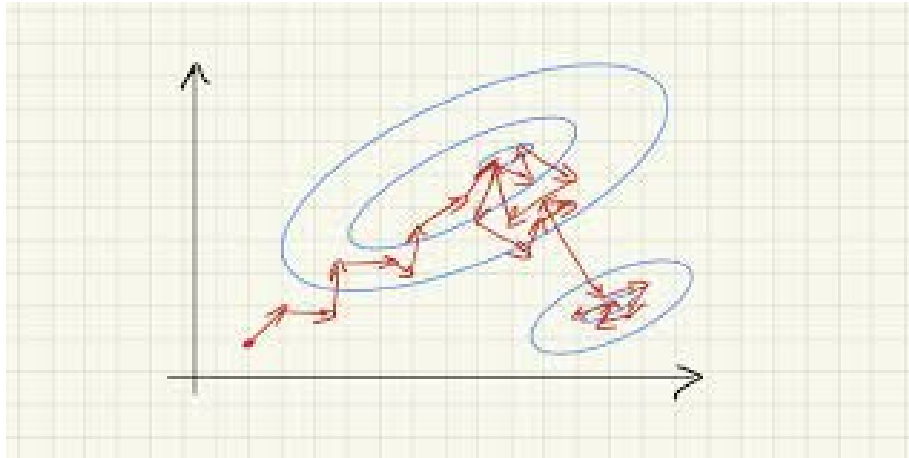
Maniac

Inference and sampling

Large computing power

Smart algorithms (Markov chain Monte Carlo)

Design stochastic dynamical system that automatically sample to required distribution



Information theory and inference

- a general inference process can be thought of as “noisy channel”
- Idea applicable to prediction, sensing, parameter estimation, hypothesis testing
- “message” to be reconstructed can be signal, parameter, hypothesis
- Information theory quantifies *how much information is contained in the data*
- Information theory establishes *fundamental limits* of inference
- Information theory suggests *recipes for data processing*

Example: Fano's inequality

- message $X \in \{1, \dots, M\}$
- “noisy” output Y
- reconstructed message $\hat{X} = f(Y)$
- probability of error: $P_e = \text{Prob}(\hat{X} \neq X)$
- Fano's inequality:

$$P_e \geq 1 - \frac{I(X : Y) + 1}{\log(M)}$$

- With n independent repetitions

$$P_e \geq 1 - \frac{nI(X : Y) + 1}{\log(M)}$$

- to achieve $P_e \leq \delta$,

$$n \geq \frac{(1 - \delta) \log(M) - 1}{I(X : Y)}$$

Example: Fano's inequality

- unknown variable $X \in \{1, \dots, M\}$
- known data Y
- estimator $\hat{X} = f(Y)$
- probability of inference error: $P_e = \text{Prob}(\hat{X} \neq X)$
- minimum data size to achieve $P_e \leq \delta$,

$$n \geq \frac{(1 - \delta) \log(M) - 1}{I(X : Y)}$$