

TAM 598 Lecture 11 :

Bayesian Inference - Selecting Your Prior

Announcements:

- HW 3 covers lectures 8-12 ; due on Mar 12

I. Selecting Prior Information

prior should reflect your beliefs about the variables before you see any data

approaches to selection:

(1) principle of insufficient reason) principle of indifference

(2) principle of maximum entropy

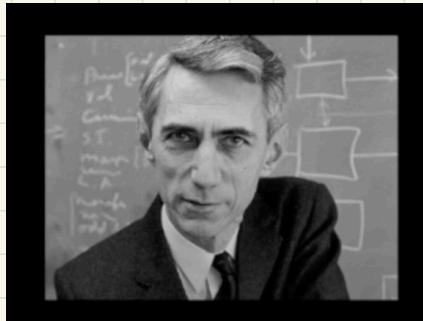
(3) principle of transformation groups

(1) principle of insufficient reason/ principle of indifference (Laplace)

The theory of chance consists in reducing all the events of the same kind to a certain number of cases equally possible, that is to say, to such as we may be equally undecided about in regard to their existence, and in determining the number of cases favorable to the event whose probability is sought. The ratio of this number to that of all the cases possible is the measure of this probability, which is thus simply a fraction whose numerator is the number of favorable cases and whose denominator is the number of all the cases possible. *Pierre-Simon Laplace*

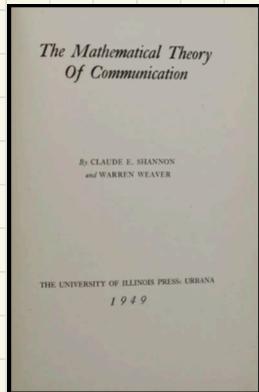
(2) principle of maximum entropy - pick most "uncertain" distribution

If you have some information about R.V. X , such as its expectation $E[X]$ or variance $V[X]$.

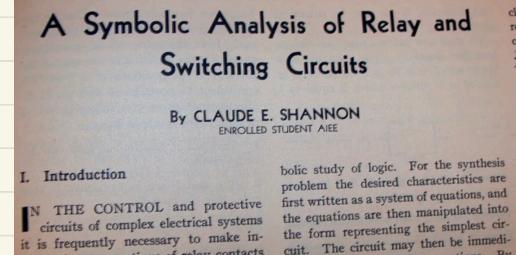


Shannon, Bell Labs

Shannon Information Content



Consider a discrete random variable X . What is the information content associated with a particular measurement $X = x$?



Why is this a good proposal for information content?

- (1) Deterministic outcomes contain no information.
- (2) Information content increases with decreasing probability.

Why is this a good proposal for information content?

- (3) Information content is additive for independent RVs

Example Submarine (D. Mackay)

?	?	?	?	?	?
?	?	?	?	?	?
?	?	?	?	?	?
?	?	?	?	?	?
?	?	?	?	?	?
?	?	?	?	?	?

Start: $P(X=x) = \frac{1}{36}$ for $x = 1, 2, \dots, 36$

Now we can think of the expected information content associated with a given random variable

$$H[X] = \sum_x p(X=x) h(X=x)$$

Example: Information entropy of a binary distribution.

We know that there are two possible outcomes, 0 and 1.
We know nothing else. What $P(x)$ should we choose?

Given some known information about a distribution, how do we select $p(x)$?

Given some known information about a distribution, how do we select $p(x)$?

Our problem now is to maximize $H[p(x)] = - \sum_i p_i \log p_i$

Subject to constraints:

$$\sum_i p_i = 1 \quad \text{normalization}$$

$$\mathbb{E}[f_k(x)] = F_k \quad \text{information } k=1, \dots, K$$

The general solution (**Karush-Kuhn-Tucker**) conditions

Example : Brandeis dice problem (E.T. Jaynes, 1962)

You toss a die N times. You measure $\mathbb{E}[x] = 4.5$ instead of the "fair dice" expected value of 3.5. Given this info and nothing else, what $p(x)$ should we assign to the next toss?

To find λ we need to solve:

Examples: discrete maximum entropy distributions:

• X takes N different values \Rightarrow Categorical $(\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N})$

• X is 0 with prob $(1-\theta)$, is 1 with prob θ
 \Rightarrow Bernoulli (θ)

• X takes values $0, 1, \dots, n$ with known expectation
 $E[X] = \theta n$ $\underline{\hspace{10em}}$ \Rightarrow Binomial $B(\theta, n)$

• X takes values $0, 1, 2, \dots, \infty$ with known $E[X] = \lambda$
 \Rightarrow Poisson (λ)

• thermodynamic ensembles: max entropy of
NVT canonical - states of a quantum system w/ known
expected energy

NVT grand canonical - states of a quant. system w/
varying # of particles, known expected #, Energy (15)

for continuous distributions, to a large extent:

$$H[X] = - \sum_x p(x) \log p(x) \longrightarrow - \int p(x) \ln p(x) dx$$

- X lies in $[a, b] \Rightarrow U[a, b]$
- X has $E[X] = \mu, V[X] = \sigma^2 \Rightarrow N(\mu, \sigma^2)$
- X has $E[X] = \mu, \text{ covariance } C[X, X] = \Sigma \Rightarrow N(\mu, \Sigma)$
- X lies in $[0, \infty)$ with $E[X] = \lambda^{-1} \Rightarrow Exp(\lambda)$

Finding a continuous distribution that maximizes the Shannon entropy is a little trickier. A common approach is the method of moments

For the discrete case we have:

$$p(x=x_i) = \frac{1}{Z} \exp \left\{ \sum_{k=1}^K \lambda_k f_k(x_i) \right\}$$

where $f_k(x_i)$ is some function of x_i and $E[f_k(x_i)] = F_k$

For the continuous case: we match expectations of moments $m=1, \dots, M$