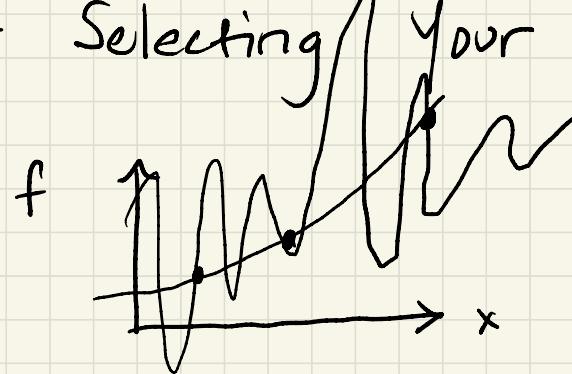


TAM 598

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:  
assign the same probability to all possible variable values
- (2) principle of maximum entropy - pick the most "uncertain" distribution that is consistent with what you know
- (3) principle of transformation groups - pick a distribution that is invariant under the problems' symmetries

# (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*

if R.V,  $X$  can take  $N$  values, then we assign

$$p(x) = \frac{1}{N} \quad \text{for} \quad x \in \{1, 2, \dots, N\}$$

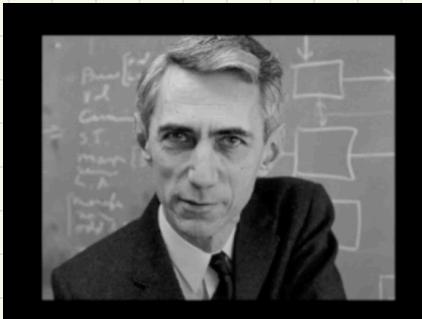
e.g) throwing a die  $p(x) = \frac{1}{6} \quad x \in \{1, 2, \dots, 6\}$

(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  $V[X]$ .

choose a maximally uncertain distribution that is consistent with the available information.

The uncertainty of a distribution is measured by its information entropy (Claude Shannon, 1948)



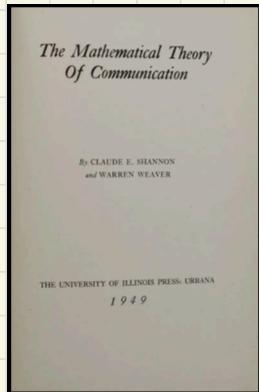
Shannon, Bell Labs

uncertainty  $H[H(p(x))]$

$$= - \sum_x p(x) \underbrace{\log_2 p(x)}_{\text{discrete}}$$

$$= - \int p(x) \underbrace{\ln p(x)}_{\text{continuous}} dx$$

# Shannon Information Content

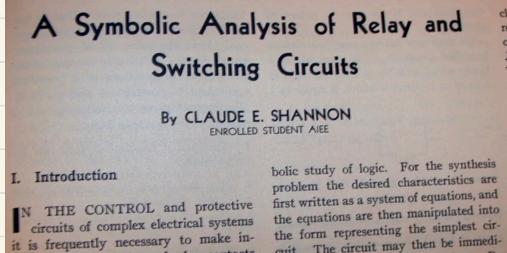


Consider a discrete random variable  $X$ . What is the information content associated with a particular measurement  $X = x$ ?

proposal :

$$h(X=x) = \underline{\log_2 \frac{1}{P(X=x)}} = -\log_2 P(X=x)$$

We can think of this as a measure of "surprise." We are more surprised when we measure an outcome that we thought was less probable.



Why is this a good proposal for information content?

(1) Deterministic outcomes contain no information.

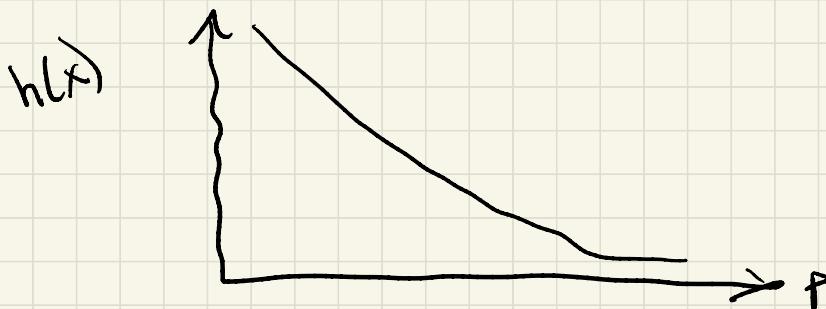
say  $P(X=x) = 1$

$$\text{Then } h(X=x) = \log_2 \frac{1}{P(X=x)} = \log_2 1 = 0$$

(2) Information content increases with decreasing probability.

If  $P(X=x) < P(X=x')$  then  $h(X=x) > h(X=x')$

consider  $\frac{d}{dp} h = \frac{d}{dp} \log_2 \frac{1}{p} = -\frac{1}{p \ln p} < 0$  for  $p > 0$



(6)

Why is this a good proposal for information content?

(3) Information content is additive for independent RVs

if  $P(X=x, Y=y) = P(X=x) P(Y=y)$

then  $h(X=x, Y=y) = h(X=x) + h(Y=y)$

$$h(X=x, Y=y) = \log_2 \frac{1}{P(X=x) P(Y=y)}$$

$$= \log_2 \frac{1}{P(X=x)} + \log_2 \frac{1}{P(Y=y)}$$

$$= h(X=x) + h(Y=y)$$

Example

Submarine (D. Mackay)

- gain information as we do measurements

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| ? | ? | ? | ? | ? | ? |
| ? | ? | ? | ? | X | ? |
| X | ? | ? | ? | ? | ? |
| ? | ? | ? | ? | ? | ? |
| ? | ? | ? | ? | ? | ? |
| ? | ? | ? | ? | ? | ? |

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

(1) First measurement is a miss

$$h(X=x_1) = \log_2 \frac{1}{P(x_1)} = \log_2 \left( \frac{36}{35} \right) = 0.0406 \text{ bits}$$

(2) Second measurement is a miss

$$h(X=x_1) + h(X=x_2) = \log_2 \left( \frac{36}{35} \right) + \log_2 \left( \frac{35}{34} \right) = 0.0824 \text{ bits}$$

(3) First 19 measurements are misses

$$\log_2 \left( \frac{36}{35} \right) + \log_2 \left( \frac{35}{34} \right) + \dots + \log_2 \left( \frac{19}{18} \right) = \log_2 \left( \frac{36}{18} \right) = 1$$

(4) Find the submarine on our 25th measurement

$$\log_2 \left( \frac{36}{35} \right) + \log_2 \left( \frac{35}{34} \right) + \dots + \log_2 \left( \frac{12}{11} \right) + \log_2 \left( \frac{11}{1} \right) = \log_2 36 = 5.17 \text{ bits}$$

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

$$\begin{aligned} \text{H}[X] &= \sum_x p(X=x) h(X=x) \\ \text{Shannon Entropy} &= \sum_x p(X=x) \log_2 \frac{1}{p(X=x)} \\ \text{Information Entropy} &= -\sum_x p(x) \log_2 p(x) \end{aligned}$$

measures how diffuse the range of outcomes is for a given R.V

eg) 5 possible outcomes

opt 1

|  |                       |                       |                       |                       |
|--|-----------------------|-----------------------|-----------------------|-----------------------|
|  | $\boxed{\phantom{0}}$ | $\boxed{\phantom{0}}$ | $\boxed{\phantom{0}}$ | $\boxed{\phantom{0}}$ |
|  | 0.2                   | 0.2                   | 0.2                   | 0.2                   |

$$\begin{aligned} H[X] &= -5(0.2) \log_2 (0.2) \\ &= 2.32 \end{aligned}$$

opt 2

|                       |                       |                       |                       |                       |
|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| $\boxed{\phantom{0}}$ | $\boxed{\phantom{0}}$ | $\boxed{\phantom{0}}$ | $\boxed{\phantom{0}}$ | $\boxed{\phantom{0}}$ |
| 0.35                  | 0.1                   | 0.35                  | 0.1                   | 0.1                   |

$$\begin{aligned} H[X] &= 2.05 \\ &\quad -2(0.35) \log_2 (0.35) \\ &\quad -3(0.1) \log_2 (0.1) \end{aligned}$$

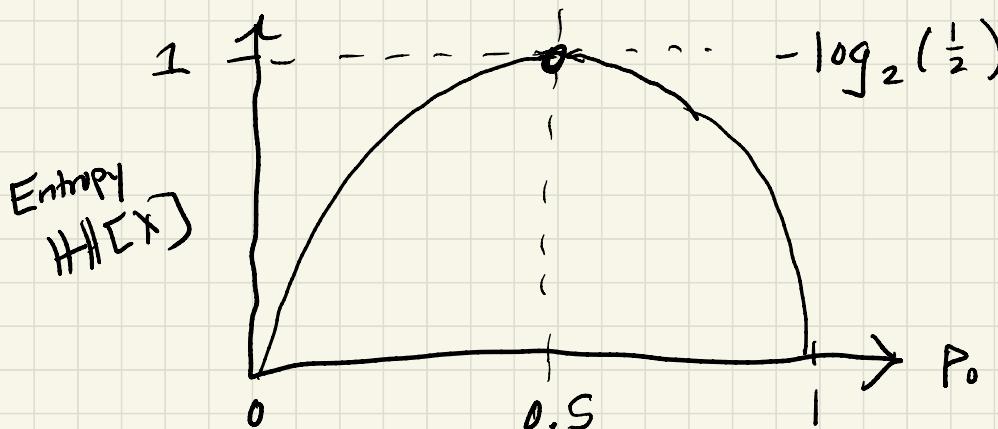
Example: Information entropy of a binary distribution.

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

$$\text{let } P_0 = P(X=0)$$

$$P_1 = 1 - P_0 = P(X=1)$$

$$H[X] = \underbrace{-P_0 \log_2 P_0 - (1-P_0) \log_2 (1-P_0)}$$



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

Our information about  $X$  comes in the form

$$\boxed{\mathbb{E}[f_k(x)] = F_k}$$

where  $f_k$  is a known function, and  $F_k$  are known constants

$$k = 1, 2, \dots, K$$

{ eg) info I = "the expected value of  $X$  is  $\mu$ "

$$\mathbb{E}[x] = \mu \quad k=1, f_1(x) = x, F_1 = \mu$$

eg) I = "the expectation is  $\mu$ ; variance is  $\sigma^2$ "

$$\mathbb{E}[x] = \mu$$

$$\mathbb{E}[x^2] = \sigma^2 + \mu^2 \quad \text{since } \mathbb{V}[x] = \mathbb{E}[x^2] - (\mathbb{E}[x])^2$$

$$K=2, f_1(x) = x, f_2(x) = x^2, F_1 = \mu, F_2 = \sigma^2 + \mu^2$$

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:

$$\left\{ \begin{array}{l} \sum_i p_i = 1 \quad \text{normalization} \\ \mathbb{E}[f_k(x)] = F_k \quad \text{information } k=1, \dots, K \end{array} \right.$$

The general solution (Karush-Kuhn-Tucker) conditions

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

System of  
nonlinear  
equations  
 $k=1, \dots, K$

$$\lambda_k = \text{constants obtained by } F_k = \frac{\partial}{\partial \lambda_k} \log Z$$
$$Z = \text{partition function, normalization} = \sum_i \exp \left\{ \sum_{k=1}^K \lambda_k f_k(x_i) \right\}$$

## 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?

$$\text{constraint: } \mathbb{E}[x] = \sum_{x=1}^6 x p(x) = 4.5$$

$$\text{general solution: } p(x) = \frac{\exp(\lambda x)}{Z(\lambda)}$$

$$\begin{aligned}\text{where } Z(\lambda) &= \sum_{i=1}^6 e^{\lambda i} \\ &= e^{\lambda} + e^{2\lambda} + \dots e^{6\lambda}\end{aligned}$$

To find  $\lambda$  we need to solve:

$$\frac{\partial}{\partial \lambda} \log Z = \frac{1}{Z} \left( \frac{\partial Z}{\partial \lambda} \right) = 4.5$$

$$\frac{1}{Z} \frac{\partial}{\partial \lambda} \left[ (e^\lambda)^1 + (e^\lambda)^2 + \dots + (e^\lambda)^b \right] = 4.5$$

$$\frac{1}{Z} \left[ e^\lambda + 2e^{2\lambda} + \dots + b e^{b\lambda} \right] = 4.5$$

$$\boxed{\frac{1}{Z} \sum_{i=1}^b i e^{i\lambda} = 4.5}$$

Need the root of  $f(\lambda) = \frac{1}{Z(\lambda)} \sum_{i=1}^b i e^{\lambda i} - 4.5$

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  $\theta$   $\Rightarrow$  Bernoulli  $(\theta)$
- $X$  takes values  $0, 1, \dots, n$  with known expectation  
 $E[X] = \theta n$   $\Rightarrow$  Binomial  $B(\theta, n)$
- $X$  takes values  $0, 1, 2, \dots, \infty$  with known  $E[X] = \lambda$   $\Rightarrow$  Poisson  $(\lambda)$

• thermodynamic ensembles: max entropy of

{ NVT  
NVT }

canonical - states of a ~~quantum~~ system w/ known mechanical expected energy

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

for continuous distributions, to a large extent:

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

- $X$  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  $m=1, \dots, M$