

Introduction to Bayesian Statistics - 3

Edoardo Milotti

Università di Trieste and INFN-Sezione di Trieste

Prior distributions

The choice of prior distribution is an important aspect of Bayesian inference

- prior distributions are one of the main targets of frequentists: how much do posteriors differ when we choose different priors?
- there are two main “objective” methods for the choice of priors
 1. Jeffreys' method
 2. The Maximum Entropy Method

Changing variables

We can reformulate a problem in terms of different random variables. But moving from, say x to x^2 , means changing the prior PDF. *This is obviously related to the fact that in many cases we do not know what the physically meaningful elementary, equiprobable events are.*

This means that result depends on the choice of the random variable. Can we obviate this problem?

Recall that the conservation of probability implies

$$p_y(y)dy = p_x(x)dx = p_x[x(y)] \left| \frac{dx}{dy} \right| dy$$

Therefore

$$p_y(y) = p_x[x(y)] \left| \frac{dx}{dy} \right|$$

Jeffreys started from the averaged distribution of data, which – in the context of a given model – is described by the likelihood function. Then he took a scale-invariant version of it, the log-derivative, and to correct for possible negativity, he took the RMS, i.e.,

$$p(\theta) \propto \sqrt{\mathbf{E} \left[\left(\frac{d}{d\theta} \ln L(x, \theta) \right)^2 \right]} = \sqrt{I(\theta)}$$

Indeed, taking the transformation

$$\varphi = \varphi(\theta)$$

we find

$$\begin{aligned} p(\varphi) &\propto \sqrt{\mathbf{E} \left[\left(\frac{d}{d\theta} \ln L(x, \theta) \right)^2 \right] \left| \frac{d\theta}{d\varphi} \right|} = \sqrt{\mathbf{E} \left[\left(\frac{d\theta}{d\varphi} \frac{d}{d\theta} \ln L(x, \theta(\varphi)) \right)^2 \right]} \\ &= \sqrt{\mathbf{E} \left[\left(\frac{d}{d\varphi} \ln L(x, \varphi) \right)^2 \right]} \end{aligned}$$

The rationale of Jeffreys' choice is that, although we do not know the physical constraints that determine the choice of the “least informative prior”, we can use the physical model provided by the Likelihood as a substitute.

Example: Gaussian case

$$L = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

$$\ln L = \ln \sqrt{2\pi} - \ln \sigma - \frac{(x - \mu)^2}{2\sigma^2}$$

1. with respect to mean

$$p(\mu) \sim \text{constant}$$

2. with respect to standard deviation (Jeffreys' prior, usually taken to mean “scale invariance”, additive in a log scale)

$$p(\sigma) \sim \frac{1}{\sigma}$$

A short refresher on (Boltzmann's) entropy in statistical mechanics

- consider a system where states n are occupied by N_n identical particles ($n, n=1, \dots, M$).
- the number of ways to fill these states is given by

$$\Omega = \frac{N!}{N_1!N_2!\dots N_M!}$$

- then Boltzmann's entropy is

$$\begin{aligned} S_B &= k_B \ln \Omega = k_B \ln \frac{N!}{N_1!N_2!\dots N_M!} \approx k_B \left((N \ln N - N) - \sum_n (N_n \ln N_n - N_n) \right) \\ &= k_B \left(N \ln N - \sum_n N p_n (\ln p_n + \ln N) \right) = k_B \sum_n p_n \ln \frac{1}{p_n} \end{aligned}$$

$$S_B = k_B \sum_i p_i \ln \frac{1}{p_i}$$

*Boltzmann's entropy is just like
Shannon's entropy*

probability of physical
states

$$S_I = \sum_i p_i \log_2 \frac{1}{p_i}$$

*Shannon's entropy is the average
information output by a source of
symbols*

this logarithmic function is
the information carried by
the i -th symbol


probability of source
symbols

Examples:


- just two symbols, 0 and 1, same source probability

$$S_I = -2 \left(\frac{1}{2} \log_2 \frac{1}{2} \right) = 1 \text{ bit}$$


there are 2
equal terms



average information
conveyed by each
symbol



the result is given in
pseudounit “bits” (for
natural logarithms this is
“nats”)



- just two symbols, 0 and 1, probabilities $\frac{1}{4}$ and $\frac{3}{4}$, respectively

$$S_I = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} \approx 0.81 \text{ bit}$$

- 8 symbols, equal probabilities

$$S_I = -\sum_1^8 \frac{1}{8} \log_2 \frac{1}{8} = \log_2 8 = 3 \text{ bit}$$

The Shannon entropy is additive for independent sources.

If symbols are emitted simultaneously and independently by two sources, the joint probability distribution is

$$p(j, k) = p_1(j)p_2(k)$$

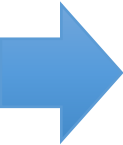
and therefore the joint entropy is

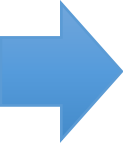
$$\begin{aligned} S &= - \sum_{j,k} p(j, k) \log_2 p(j, k) = - \sum_{j,k} p_1(j)p_2(k) \log_2 [p_1(j)p_2(k)] \\ &= - \sum_j p_1(j) \log_2 p_1(j) - \sum_k p_2(k) \log_2 p_2(k) \\ &= S_1 + S_2 \end{aligned}$$

The Shannon entropy is the highest for the uniform distribution.

This is an easy result that follows using one Lagrange multiplier to keep probability normalization into account

$$\begin{aligned} S + \lambda \sum_{k=1}^N p_k &= - \sum_{k=1}^N p_k \log_2 p_k + \lambda \sum_{k=1}^N p_k \\ &= - \frac{1}{\ln 2} \sum_{k=1}^N p_k \ln p_k + \lambda \sum_{k=1}^N p_k \end{aligned}$$


$$\frac{\partial}{\partial p_j} \left(S + \lambda \sum_{k=1}^N p_k \right) = - \frac{1}{\ln 2} (\ln p_j + 1) + \lambda = 0$$


$$p_j = \exp(\lambda \ln 2 - 1) = 1/N$$

all probabilities have the same value

Edwin T. Jaynes (1922-1998), introduced the method of maximum entropy in statistical mechanics: when we start from the informational entropy (Shannon's entropy) and we use it introduce Boltzmann's entropy we reobtain the whole of statistical mechanics by maximizing entropy.

In a sense, statistical mechanics also arises from a comprehensive "principle of maximum entropy".

<http://bayes.wustl.edu/etj/etj.html>



Information Theory and Statistical Mechanics

E. T. JAYNES

Department of Physics, Stanford University, Stanford, California

(Received September 4, 1956; revised manuscript received March 4, 1957)

Information theory provides a constructive criterion for setting up probability distributions on the basis of partial knowledge, and leads to a type of statistical inference which is called the maximum-entropy estimate. It is the least biased estimate possible on the given information; i.e., it is maximally noncommittal with regard to missing information. If one considers statistical mechanics as a form of statistical inference rather than as a physical theory, it is found that the usual computational rules, starting with the determination of the partition function, are an immediate consequence of the maximum-entropy principle. In the resulting "subjective statistical mechanics," the usual rules are thus justified independently of any physical argument, and in particular independently of experimental verification; whether

or not the results agree with experiment, they still represent the best estimates that could have been made on the basis of the information available.

It is concluded that statistical mechanics need not be regarded as a physical theory dependent for its validity on the truth of additional assumptions not contained in the laws of mechanics (such as ergodicity, metric transitivity, equal *a priori* probabilities, etc.). Furthermore, it is possible to maintain a sharp distinction between its physical and statistical aspects. The former consists only of the correct enumeration of the states of a system and their properties; the latter is a straightforward example of statistical inference.

In these papers Jaynes argues that information theory provides a constructive criterion for setting up probability distributions on the basis of partial knowledge, and leads to a type of statistical inference which is called the maximum-entropy estimate.

It is the least biased estimate possible on the given information; i.e., it is maximally noncommittal with regard to missing information.

If one considers statistical mechanics as a form of statistical inference rather than as a physical theory, it is found that the usual computational rules, starting with the determination of the partition function, are an immediate consequence of the maximum-entropy principle.

In the resulting "subjective statistical mechanics," the usual rules are justified independently of any physical argument, and in particular independently of experimental verification; whether or not the results agree with experiment, they still represent the best estimates that could have been made on the basis of the information available.

Jaynes concludes that statistical mechanics need not be regarded as a physical theory dependent for its validity on additional assumptions not contained in the laws of mechanics (such as ergodicity, metric transitivity, equal a priori probabilities, etc.).

Furthermore, it is possible to maintain a sharp distinction between physical and statistical aspects. The former consists only of the correct enumeration of the states of a system; the latter is a straightforward example of statistical inference.

Now let's move on and maximize entropy in order to solve problems and find prior distributions ...



The kangaroo problem (Jaynes)

- *Basic information:* one third of all kangaroos has blue eyes, and one third is left-handed.
- *Question:* which fraction of kangaroos has both blue eyes and is left-handed?

	left	~left
blue	1/9	2/9
~blue	2/9	4/9

no correlation

	left	~left
blue	0	1/3
~blue	1/3	1/3

maximum negative correlation

	left	~left
blue	1/3	0
~blue	0	2/3

maximum positive correlation

probabilities p_{bl} $p_{\bar{b}l}$ $p_{b\bar{l}}$ $p_{\bar{b}\bar{l}}$

entropy (proportional to Shannon's entropy)

$$S = p_{bl} \ln \frac{1}{p_{bl}} + p_{\bar{b}l} \ln \frac{1}{p_{\bar{b}l}} + p_{b\bar{l}} \ln \frac{1}{p_{b\bar{l}}} + p_{\bar{b}\bar{l}} \ln \frac{1}{p_{\bar{b}\bar{l}}}$$

constraints (3 constraints, 4 unknowns)

$$p_{bl} + p_{\bar{b}l} + p_{b\bar{l}} + p_{\bar{b}\bar{l}} = 1$$

$$p_{bl} + p_{b\bar{l}} = 1/3$$

$$p_{\bar{b}l} + p_{\bar{b}\bar{l}} = 1/3$$

entropy maximization with constraints

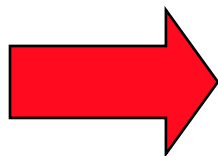
$$S_V = \left(p_{bl} \ln \frac{1}{p_{bl}} + p_{\bar{bl}} \ln \frac{1}{p_{\bar{bl}}} + p_{b\bar{l}} \ln \frac{1}{p_{b\bar{l}}} + p_{\bar{b}\bar{l}} \ln \frac{1}{p_{\bar{b}\bar{l}}} \right) \\ + \lambda_1 (p_{bl} + p_{\bar{bl}} + p_{b\bar{l}} + p_{\bar{b}\bar{l}} - 1) + \lambda_2 (p_{bl} + p_{b\bar{l}} - 1/3) + \lambda_3 (p_{\bar{bl}} + p_{\bar{b}\bar{l}} - 1/3)$$

$$\frac{\partial S_V}{\partial p_{bl}} = -\ln p_{bl} - 1 + \lambda_1 + \lambda_2 + \lambda_3 = 0$$

$$\frac{\partial S_V}{\partial p_{\bar{bl}}} = -\ln p_{\bar{bl}} - 1 + \lambda_1 + \lambda_3 = 0$$

$$\frac{\partial S_V}{\partial p_{b\bar{l}}} = -\ln p_{b\bar{l}} - 1 + \lambda_1 + \lambda_2 = 0$$

$$\frac{\partial S_V}{\partial p_{\bar{b}\bar{l}}} = -\ln p_{\bar{b}\bar{l}} - 1 + \lambda_1 = 0$$



$$p_{bl} = \exp(-1 + \lambda_1 + \lambda_2 + \lambda_3)$$

$$p_{\bar{bl}} = \exp(-1 + \lambda_1 + \lambda_3)$$

$$p_{b\bar{l}} = \exp(-1 + \lambda_1 + \lambda_2)$$

$$p_{\bar{b}\bar{l}} = \exp(-1 + \lambda_1)$$

$$\begin{cases} p_{\bar{b}l} = p_{\bar{b}l} \exp(\lambda_3) \\ p_{b\bar{l}} = p_{\bar{b}l} \exp(\lambda_2) \\ p_{bl} = p_{\bar{b}l} \exp(\lambda_2 + \lambda_3) \end{cases} \Rightarrow p_{\bar{b}l} p_{b\bar{l}} = p_{bl} p_{\bar{b}l}$$

$$\begin{cases} p_{bl} + p_{\bar{b}l} + p_{b\bar{l}} + p_{\bar{b}\bar{l}} = 1 \\ p_{bl} + p_{b\bar{l}} = 1/3 \\ p_{bl} + p_{\bar{b}l} = 1/3 \\ p_{\bar{b}l} p_{b\bar{l}} = p_{bl} p_{\bar{b}\bar{l}} \end{cases} \Rightarrow \begin{cases} p_{b\bar{l}} = p_{\bar{b}\bar{l}} = 1/3 - p_{bl} \\ p_{\bar{b}l} = 1/3 + p_{bl} \\ (1/3 - p_{bl})^2 = p_{bl}/3 + p_{bl}^2 \\ 1/9 - 2p_{bl}/3 + p_{bl}^2 = p_{bl}/3 + p_{bl}^2 \end{cases}$$

$$\Rightarrow p_{bl} = \frac{1}{9}; \quad p_{b\bar{l}} = p_{\bar{b}l} = \frac{2}{9}; \quad p_{\bar{b}\bar{l}} = \frac{4}{9}$$

this solution coincides with the least informative distribution (no correlation)

Solution of underdetermined systems of equations

In this problem there are fewer equations than unknowns; the system of equations is underdetermined, and in general there is no unique solution.

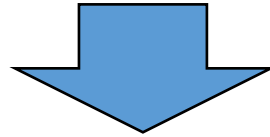
The maximum entropy method helps us find a reasonable solution, the least informative one (least correlations between variables)

Example:

$$\begin{aligned} 3x + 5y + 1.1z &= 10 \\ -2.1x + 4.4y - 10z &= 1 \end{aligned} \quad (x, y, z > 0)$$

$$\begin{aligned}
 3x + 5y + 1.1z &= 10 \\
 -2.1x + 4.4y - 10z &= 1
 \end{aligned}
 \quad (x, y, z > 0)$$

this ratio can be taken to be a "probability"



$$\begin{aligned}
 S &= - \left(\frac{x}{x+y+z} \ln \frac{x}{x+y+z} + \frac{y}{x+y+z} \ln \frac{y}{x+y+z} + \frac{z}{x+y+z} \ln \frac{z}{x+y+z} \right) \\
 &= - \frac{1}{x+y+z} \left[x \ln x + y \ln y + z \ln z - (x+y+z) \ln(x+y+z) \right]
 \end{aligned}$$

$$Q = S + \lambda(3x + 5y + 1.1z - 10) + \mu(-2.1x + 4.4y - 10z - 1)$$

$$\begin{aligned}
 \frac{\partial Q}{\partial x} &= - \frac{\ln x - \ln(x+y+z)}{x+y+z} + \frac{x \ln x + y \ln y + z \ln z - (x+y+z) \ln(x+y+z)}{(x+y+z)^2} + 3\lambda - 2.1\mu \\
 &= \frac{(y+z) \ln x + y \ln y + z \ln z}{(x+y+z)^2} + 3\lambda - 2.1\mu = 0
 \end{aligned}$$

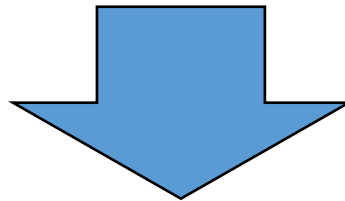
$$\frac{\partial Q}{\partial x} = \frac{(y+z)\ln x + y\ln y + z\ln z}{(x+y+z)^2} + 3\lambda - 2.1\mu = 0$$

$$\frac{\partial Q}{\partial y} = \frac{x\ln x + (x+z)\ln y + z\ln z}{(x+y+z)^2} + 5\lambda + 4.4\mu = 0$$

$$\frac{\partial Q}{\partial z} = \frac{x\ln x + y\ln y + (x+y)\ln z}{(x+y+z)^2} + 1.1\lambda - 10\mu = 0$$

$$10 = 3x + 5y + 1.1z$$

$$1 = -2.1x + 4.4y - 10z$$



$$x = 0.606275; \quad y = 1.53742; \quad z = 0.449148;$$
$$\lambda = 0.0218739; \quad \mu = -0.017793$$

this is an example of an “ill-posed” problem

the solution that we found is a kind of **regularization**
of the ill-posed problem

Finding priors with the maximum entropy method

$$S = \sum_k p_k \ln \frac{1}{p_k} = -\sum_k p_k \ln p_k \quad \text{Shannon entropy}$$

entropy maximization when all information is missing and normalization is the only constraint:


$$\frac{\partial}{\partial p_k} \left[-\sum_k p_k \ln p_k + \lambda \left(\sum_k p_k - 1 \right) \right] = -(\ln p_k + 1) + \lambda = 0$$

$$p_k = e^{\lambda-1}; \quad \sum_k p_k = \sum_k e^{\lambda-1} = N e^{\lambda-1} = 1 \quad \Rightarrow \quad p_k = 1/N$$

entropy maximization when the mean is known μ

$$\frac{\partial}{\partial p_k} \left[-\sum_k p_k \ln p_k + \lambda_0 \left(\sum_k p_k - 1 \right) + \lambda_1 \left(\sum_k x_k p_k - \mu \right) \right]$$
$$= -(\ln p_k + 1) + \lambda_0 + \lambda_1 x_k = 0$$

incomplete
solution...


$$p_k = e^{\lambda_0 + \lambda_1 x_k - 1};$$

We must satisfy two constraints now ...

$$p_k = e^{\lambda_0 + \lambda_1 x_k - 1}$$

$$\sum_k p_k = \sum_k e^{\lambda_0 + \lambda_1 x_k - 1} = e^{\lambda_0 - 1} \sum_k e^{\lambda_1 x_k} = 1$$

$$\sum_k x_k p_k = \sum_k x_k e^{\lambda_0 + \lambda_1 x_k - 1} = e^{\lambda_0 - 1} \sum_k x_k e^{\lambda_1 x_k} = \mu$$

$$e^{\lambda_0 - 1} = \frac{1}{\sum_k e^{\lambda_1 x_k}}; \quad \frac{\sum_k x_k e^{\lambda_1 x_k}}{\sum_k e^{\lambda_1 x_k}} = \mu$$

no analytic solution,
only numerical



Example : the biased die

(E. T. Jaynes: *Where do we stand on Maximum Entropy?* In *The Maximum Entropy Formalism*; Levine, R. D. and Tribus, M., Eds.; MIT Press, Cambridge, MA, 1978)

mean value of throws for an unbiased die

$$\frac{1}{6}(1 + 2 + 3 + 4 + 5 + 6) = \frac{21}{6} = 3.5$$

mean value for a biased die

$$3.5(1 + \varepsilon)$$

Problem: for a given mean value of the biased die, what is the probability distribution of each value?

The mean value is insufficient information, and we use the maximum entropy method to find the most likely distribution (the least informative one).

entropy maximization with the biased die:

$$\frac{\partial}{\partial p_k} \left[-\sum_{k=1}^6 p_k \ln p_k + \lambda_0 \left(\sum_{k=1}^6 p_k - 1 \right) + \lambda_1 \left(\sum_{k=1}^6 k p_k - \frac{7}{2}(1 + \varepsilon) \right) \right]$$
$$= -(\ln p_k + 1) + \lambda_0 + k\lambda_1 = 0$$

$$p_k = e^{\lambda_0 + \lambda_1 k - 1}$$

$$\sum_{k=1,6} p_k = e^{\lambda_0 - 1} \sum_{k=1,6} e^{\lambda_1 k} = 1$$

$$\sum_{k=1,6} k p_k = e^{\lambda_0 - 1} \sum_{k=1,6} k e^{\lambda_1 k} = \frac{7}{2}(1 + \varepsilon)$$

$$e^{\lambda_0 - 1} = \frac{1}{\sum_{k=1,6} e^{\lambda_1 k}}; \quad \frac{\sum_{k=1,6} k p_k}{\sum_{k=1,6} p_k} = \frac{7}{2}(1 + \varepsilon)$$

we still have to satisfy the constraints ...

$$e^{\lambda_0 - 1} \sum_{k=1,6} e^{\lambda_1 k} = e^{\lambda_0 - 1} \left(\sum_{k=0,6} e^{\lambda_1 k} - 1 \right) = e^{\lambda_0 - 1} \left(\frac{1 - e^{7\lambda_1}}{1 - e^{\lambda_1}} - 1 \right) = 1$$

$$\begin{aligned} \frac{\sum_{k=1,6} k e^{\lambda_1 k}}{\sum_{k=1,6} e^{\lambda_1 k}} &= \frac{\partial}{\partial \lambda_1} \ln \sum_{k=1,6} e^{\lambda_1 k} = \frac{\partial}{\partial \lambda_1} \ln \left(e^{\lambda_1} \sum_{k=0,5} e^{\lambda_1 k} \right) \\ &= \frac{\partial}{\partial \lambda_1} \left[\lambda_1 + \ln(1 - e^{6\lambda_1}) - \ln(1 - e^{\lambda_1}) \right] \\ &= 1 - \frac{6e^{6\lambda_1}}{1 - e^{6\lambda_1}} + \frac{e^{\lambda_1}}{1 - e^{\lambda_1}} = \frac{7}{2}(1 + \varepsilon) \end{aligned}$$

the Lagrange multipliers are obtained from nonlinear equations and we must use numerical methods

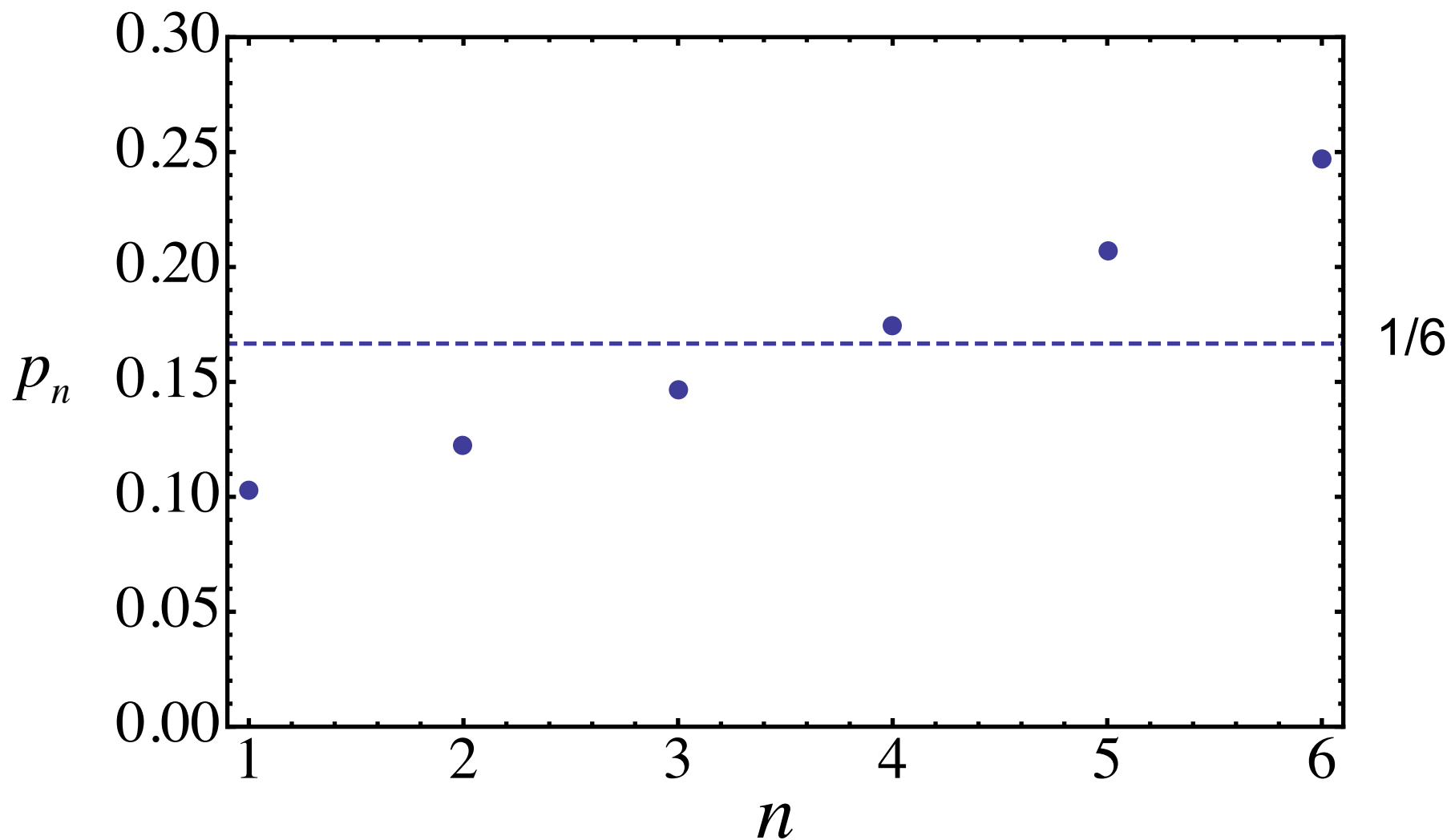
numerical solution

media	p_1	p_2	p_3	p_4	p_5	p_6
3.0	0.246782	0.20724	0.174034	0.146148	0.122731	0.103065
3.1	0.22929	0.199582	0.173723	0.151214	0.131622	0.114568
3.2	0.212566	0.191659	0.172808	0.155811	0.140487	0.126669
3.3	0.196574	0.183509	0.171313	0.159928	0.149299	0.139377
3.4	0.181282	0.175168	0.16926	0.163551	0.158035	0.152704
3.5	0.166667	0.166667	0.166667	0.166667	0.166666	0.166666
3.6	0.152704	0.158035	0.163551	0.16926	0.175168	0.181282
3.7	0.139377	0.149299	0.159928	0.171313	0.183509	0.196574
3.8	0.126669	0.140487	0.155811	0.172808	0.191659	0.212566
3.9	0.114568	0.131622	0.151214	0.173723	0.199582	0.22929
4.0	0.103065	0.122731	0.146148	0.174034	0.20724	0.246782

with a biased die we obtain skewed distributions.

These are examples of UNINFORMATIVE PRIORS

Example: mean = 4



Entropy with continuous probability distributions

(relative entropy, Kullback-Leibler divergence)

$$S \rightarrow - \int_a^b [p(x) dx] \ln [p(x) dx]$$

this diverges!

$$S_{p|m} = - \sum_k p_k \ln \frac{p_k}{m_k}$$

relative entropy

$$S_{p|m} = - \int_a^b p(x) \ln \frac{p(x)}{m(x)} dx$$

this does not diverge!

Mathematical aside on the Kullback-Leibler divergence

The obvious extension of the Shannon entropy to continuous distributions

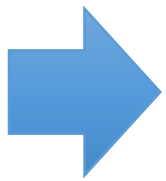
$$S = \int_{-\infty}^{+\infty} p(x) dx \log_2 \frac{1}{p(x) dx}$$

does not work, because it diverges.

A solution is suggested again by statistical mechanics ...

Boltzmann entropy with degeneracy number attached to each level

$$\Omega = \frac{N!}{N_1! N_2! \dots N_M!} g_1^{N_1} g_2^{N_2} \dots g_M^{N_M}$$



$$\ln \Omega = \ln N! - \sum_{k=1}^M \ln N_k! + \sum_{k=1}^M N_k \ln g_k$$

$$= -N \sum_{k=1}^M (N_k/N) \ln \frac{(N_k/N)}{g_k}$$

Kullback-Leibler
divergence

$$= -N \sum_{k=1}^M p_k \ln \frac{p_k}{g_k}$$



$$I_{KL} = \sum_{k=1}^M p_k \ln \frac{p_k}{g_k}$$

Properties of the Kullback-Leibler divergence

- extremal value when $p_k = g_k$.

Indeed, using again a Lagrange multiplier we must consider the auxiliary function

$$I_{KL} + \lambda \sum_k p_k$$

and we find the extremum at

$$p_k = g_k e^{\lambda-1} = g_k$$

(homework!)

normalization

- the KL divergence is a measure of the number of excess bits that we must use when we take a distribution of symbols which is different from the true distribution

$$\begin{aligned} I_{KL} &= \sum_{k=1}^M p_k \ln \frac{p_k}{g_k} \\ &= \sum_{k=1}^M p_k \ln \frac{1}{g_k} - \sum_{k=1}^M p_k \ln \frac{1}{p_k} \end{aligned}$$

- the KL divergence for continuous distributions does not diverge

$$\begin{aligned} I_{KL} &= \sum_k p_k \ln \frac{p_k}{g_k} \\ &\rightarrow \int_{-\infty}^{+\infty} p(x) dx \ln \frac{p(x) dx}{g(x) dx} \\ &= \int_{-\infty}^{+\infty} p(x) \ln \frac{p(x)}{g(x)} dx \end{aligned}$$

- the KL divergence is non-negative

Notice first that when we define $\phi(t) = t \ln t$ we find

$$\phi(t) = \phi(1) + \phi'(1)(t - 1) + \frac{1}{2}\phi''(h)(t - 1)^2 = (t - 1) + \frac{1}{2h}(t - 1)^2$$

where $t < h < 1$ and therefore

$$\begin{aligned} I_{KL} &= \int_{-\infty}^{+\infty} p(x) \ln \frac{p(x)}{g(x)} dx = - \int_{-\infty}^{+\infty} \frac{p(x)}{g(x)} \ln \frac{p(x)}{g(x)} g(x) dx = \int_{-\infty}^{+\infty} \phi\left(\frac{p(x)}{g(x)}\right) g(x) dx \\ &= \int_{-\infty}^{+\infty} \left[\left(\frac{p(x)}{g(x)} - 1\right) + \frac{1}{2h} \left(\frac{p(x)}{g(x)} - 1\right)^2 \right] g(x) dx = \int_{-\infty}^{+\infty} \frac{1}{2h} \left(\frac{p(x)}{g(x)} - 1\right)^2 g(x) dx \\ &= \int_{-\infty}^{+\infty} \frac{1}{2h} \frac{(p(x) - g(x))^2}{g(x)} dx \geq 0 \end{aligned}$$

The KL divergence is a quasi-metric (however a local version of the KL divergence is the Fisher information, which is a true metric)

The KL divergence can be used to measure the “distance” between two distributions.

Example: the KL divergence

$$I_{KL}(p, q) = \int_{-\infty}^{+\infty} p(x) \ln \frac{p(x)}{q(x)} dx$$

for the distributions

$$\begin{aligned} p(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x^2}{2\sigma^2}\right) \\ q(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \end{aligned} \quad \Rightarrow \quad I_{KL}(p, q) = \frac{\mu^2}{2\sigma^2}$$

Now consider a family of parametric distributions and evaluate the KL divergence between two close elements of the family

$$\begin{aligned} I_{KL} (p(x, \theta), p(x, \theta + \epsilon)) &= \int_{-\infty}^{+\infty} p(x, \theta) \ln \frac{p(x, \theta)}{p(x, \theta + \epsilon)} dx \\ &= \mathbf{E} (\ln p(x, \theta) - \ln p(x, \theta + \epsilon)) \end{aligned}$$

Since

$$\ln p(x, \theta + \epsilon) \approx \ln p(x, \theta) + \frac{\partial \ln p(x, \theta)}{\partial \theta} \epsilon + \frac{1}{2} \frac{\partial^2 \ln p(x, \theta)}{\partial \theta^2} \epsilon^2$$

we find, using the first Bartlett identity,

$$\begin{aligned} I_{KL} (p(x, \theta), p(x, \theta + \epsilon)) &= -\mathbf{E} \left(\frac{\partial \ln p(x, \theta)}{\partial \theta} \epsilon + \frac{1}{2} \frac{\partial^2 \ln p(x, \theta)}{\partial \theta^2} \epsilon^2 \right) \\ &= -\frac{1}{2} \mathbf{E} \left[\frac{\partial^2 \ln p(x, \theta)}{\partial \theta^2} \right] \epsilon^2 = \frac{1}{2} I(\theta) \epsilon^2 \end{aligned}$$

i.e., locally the KL divergence is just the Fisher information

Homework: go back to the estimate of the parameter of the binomial distribution and find the KL divergence of successive estimates

End of mathematical aside

Entropy extremization with additional conditions (partial knowledge of moments of the prior distribution)

$$\langle x^k \rangle = \int_a^b x^k p(x) dx$$

function (functional) that must be extremized

$$Q[p] = - \int_a^b p(x) \ln \frac{p(x)}{m(x)} dx + \sum_k \lambda_k \left\{ \int_a^b x^k p(x) dx - M_k \right\}$$

variation

$$\delta Q = - \int_a^b \delta p \left\{ \ln \frac{p(x)}{m(x)} + 1 - \sum_k \lambda_k x^k \right\} dx = 0$$

$$\ln \frac{p(x)}{m(x)} + 1 - \sum_k \lambda_k x^k = 0$$

$$p(x) = m(x) \exp \left(\sum_k \lambda_k x^k - 1 \right)$$

$$p(x) = m(x) \exp\left(\sum_n \lambda_n x^n - 1\right)$$

$p(x)$ is determined by the choice of $m(x)$ and by the constraints

The constraints can be the moments themselves:

$$M_k = \int_a^b x^k m(x) \exp\left(\sum_n \lambda_n x^n - 1\right) dx$$

1. no moment is known, normalization is the only constraint, and $p(x)$ is defined in the interval (a,b)

$$M_0 = \int_a^b m(x) \exp(\lambda_0 - 1) dx = 1$$

we take a reference distribution which is uniform on (a,b) ,
i.e.,

$$m(x) = \frac{1}{b-a}$$

$$M_0 = \frac{1}{b-a} \int_a^b \exp(\lambda_0 - 1) dx = \exp(\lambda_0 - 1) = 1$$

$$\Rightarrow \lambda_0 = 1; \quad p(x) = m(x) \exp\left(\sum_{n=0}^0 \lambda_n x^n - 1\right) = \frac{1}{b-a}$$

2. only the first moment is known, i.e, the mean, and $p(x)$ is defined on (a,b)

$$M_0 = \frac{1}{b-a} \int_a^b \exp(\lambda_0 + \lambda_1 x - 1) dx = 1$$

$$M_1 = \frac{1}{b-a} \int_a^b x \exp(\lambda_0 + \lambda_1 x - 1) dx$$

$$M_0 = 1 = \frac{\exp(\lambda_0 - 1)}{b-a} \int_a^b \exp(\lambda_1 x) dx = \frac{\exp(\lambda_0 - 1)}{b-a} \cdot \frac{\exp(\lambda_1 b) - \exp(\lambda_1 a)}{\lambda_1}$$

$$M_1 = \frac{\exp(\lambda_0 - 1)}{b-a} \int_a^b x \exp(\lambda_1 x) dx = \frac{\exp(\lambda_0 - 1)}{b-a} \left[\frac{1}{\lambda_1} (b \exp(\lambda_1 b) - a \exp(\lambda_1 a)) - \frac{1}{\lambda_1^2} (\exp(\lambda_1 b) - \exp(\lambda_1 a)) \right]$$

in general these equations can only be solved numerically...

special case:

$$a \rightarrow -\frac{L}{2}; \quad b \rightarrow \frac{L}{2}; \quad M_1 = 0$$

$$\frac{\exp(\lambda_0 - 1) \exp(\lambda_1 L/2) - \exp(-\lambda_1 L/2)}{L \lambda_1} = 1$$

$$\frac{\exp(\lambda_0 - 1)}{L} \left[\frac{1}{\lambda_1} \left(\frac{L}{2} \exp(\lambda_1 L/2) + \frac{L}{2} \exp(-\lambda_1 L/2) \right) - \frac{1}{\lambda_1^2} (\exp(\lambda_1 L/2) - \exp(-\lambda_1 L/2)) \right] = 0$$

$$\frac{\exp(\lambda_0 - 1) \exp(\lambda_1 L/2) - \exp(-\lambda_1 L/2)}{L \lambda_1} = 1$$

$$\frac{L}{2} (\exp(\lambda_1 L/2) + \exp(-\lambda_1 L/2)) - \frac{1}{\lambda_1} (\exp(\lambda_1 L/2) - \exp(-\lambda_1 L/2)) = 0$$

$$\exp(\lambda_0 - 1) \frac{\sinh(\lambda_1 L/2)}{\lambda_1 L/2} = 1$$

$$L \cosh(\lambda_1 L/2) - \frac{2}{\lambda_1} \sinh(\lambda_1 L/2) = 0$$

$$\Rightarrow (\lambda_1 L/2) = \tanh(\lambda_1 L/2) \Rightarrow \lambda_1 = 0; \quad \lambda_0 = 1$$

$$p(x) = m(x) \exp\left(\sum_{k=0}^1 \lambda_k x^k - 1\right) = \frac{1}{L}$$

nonzero mean

$$a \rightarrow -\frac{L}{2}; \quad b \rightarrow \frac{L}{2}; \quad M_1 = \varepsilon$$

$$\frac{\exp(\lambda_0 - 1)}{L} \cdot \frac{\exp(\lambda_1 L/2) - \exp(-\lambda_1 L/2)}{\lambda_1} = 1$$

$$\frac{\exp(\lambda_0 - 1)}{\lambda_1 L} \left[\frac{L}{2} (\exp(\lambda_1 L/2) + \exp(-\lambda_1 L/2)) - \frac{1}{\lambda_1} (\exp(\lambda_1 L/2) - \exp(-\lambda_1 L/2)) \right] = \varepsilon$$

$$\frac{\exp(\lambda_0 - 1)}{(\lambda_1 L/2)} \cdot \sinh(\lambda_1 L/2) = 1$$

$$\frac{L}{2} \frac{1}{\tanh(\lambda_1 L/2)} - \frac{1}{\lambda_1} = \varepsilon$$

$$\tanh(\lambda_1 L/2) = \left(\frac{1}{\lambda_1 L/2} + \frac{2\varepsilon}{L} \right)^{-1} \quad \tanh(z) = \left(\frac{1}{z} + \frac{2\varepsilon}{L} \right)^{-1}$$

this is similar to the equations of ferromagnetism

$$z - \frac{z^3}{3} \approx \left(\frac{1}{z} + \frac{2\varepsilon}{L} \right)^{-1} \Rightarrow \left(z - \frac{z^3}{3} \right) \left(\frac{1}{z} + \frac{2\varepsilon}{L} \right) \approx 1 + \frac{2\varepsilon}{L} z - \frac{z^2}{3} = 1$$

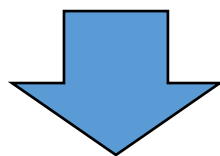
$$\Rightarrow \frac{2\varepsilon}{L} - \frac{z}{3} \approx 0 \Rightarrow z \approx \frac{6\varepsilon}{L}$$

$$\frac{\lambda_1 L}{2} \approx \frac{6\varepsilon}{L} \Rightarrow p(x) \approx \frac{1}{L} \exp(\lambda_1 x) \approx \frac{1}{L} \left(1 - \frac{12\varepsilon}{L} x \right)$$

another special case $a = 0; b \rightarrow \infty$

$$M_0 = \frac{1}{b-a} \int_a^b \exp(\lambda_0 + \lambda_1 x - 1) dx = 1$$

$$M_1 = \frac{1}{b-a} \int_a^b x \exp(\lambda_0 + \lambda_1 x - 1) dx$$



$$M_0 = 1 = m_0 \exp(\lambda_0 - 1) \cdot \frac{1}{(-\lambda_1)}$$

$$M_1 = m_0 \exp(\lambda_0 - 1) \left[\frac{1}{\lambda_1^2} \right] = (-\lambda_1) \left[\frac{1}{\lambda_1^2} \right] = -\frac{1}{\lambda_1} = \langle x \rangle$$

then

$$m_0 \exp(\lambda_0 - 1) = -\lambda_1 = \frac{1}{\langle x \rangle}$$

and we obtain the exponential distribution

$$\begin{aligned} p(x) &= m(x) \exp\left(\sum_n \lambda_n x^n - 1\right) \\ &= m_0 \exp(\lambda_0 - 1) \exp(\lambda_1 x) = \frac{1}{\langle x \rangle} \exp\left(-\frac{x}{\langle x \rangle}\right) \end{aligned}$$

3. both mean and variance are known, and the interval is the whole real axis

$$M_0 = m_0 \int_a^b \exp(\lambda_0 + \lambda_1 x + \lambda_2 x^2 - 1) dx = 1$$

$$M_1 = m_0 \int_a^b x \exp(\lambda_0 + \lambda_1 x + \lambda_2 x^2 - 1) dx$$

$$M_2 = m_0 \int_a^b x^2 \exp(\lambda_0 + \lambda_1 x + \lambda_2 x^2 - 1) dx$$

$$\begin{aligned} \exp(\lambda_0 + \lambda_1 x + \lambda_2 x^2 - 1) &= \exp \left[\lambda_2 \left(x^2 + 2 \frac{\lambda_1}{\lambda_2} x + \frac{\lambda_1^2}{\lambda_2^2} \right) + \left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2} \right) \right] \\ &= \exp \left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2} \right) \exp \left[\lambda_2 \left(x + \frac{\lambda_1}{\lambda_2} \right)^2 \right] \end{aligned}$$

$$M_0 = m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) \int_{-\infty}^{+\infty} \exp\left[-\frac{1}{2(-1/2\lambda_2)}\left(x + \frac{\lambda_1}{\lambda_2}\right)^2\right] dx = m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) \sqrt{-\frac{\pi}{\lambda_2}} = 1$$

$$M_1 = m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) \int_{-\infty}^{+\infty} x \exp\left[-\frac{1}{2(-1/2\lambda_2)}\left(x + \frac{\lambda_1}{\lambda_2}\right)^2\right] dx = m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) \sqrt{-\frac{\pi}{\lambda_2}} \left(-\frac{\lambda_1}{\lambda_2}\right) = -\mu$$

$$M_2 = m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) \int_{-\infty}^{+\infty} x^2 \exp\left[-\frac{1}{2(-1/2\lambda_2)}\left(x + \frac{\lambda_1}{\lambda_2}\right)^2\right] dx = m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) \sqrt{-\frac{\pi}{\lambda_2}} \left(-\frac{1}{2\lambda_2} + \frac{\lambda_1^2}{\lambda_2^2}\right) = \sigma^2 + \mu^2$$

$$M_0 = m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) \sqrt{-\frac{\pi}{\lambda_2}} = 1$$

$$M_1 = \frac{\lambda_1}{\lambda_2} = \mu$$

$$M_2 = \left(-\frac{1}{2\lambda_2} + \frac{\lambda_1^2}{\lambda_2^2}\right) = \sigma^2 + \mu^2$$

$$\Rightarrow \lambda_1 = -\frac{\mu}{2\sigma^2}; \quad \lambda_2 = -\frac{1}{2\sigma^2}; \quad m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) = \frac{1}{\sqrt{2\pi\sigma^2}}$$

$$\begin{aligned}
p(x) &= m_0 \exp(\lambda_0 + \lambda_1 x + \lambda_2 x^2 - 1) \\
&= m_0 \exp\left(\lambda_0 - 1 - \frac{\lambda_1^2}{\lambda_2}\right) \exp\left[-\frac{1}{2(-1/2\lambda_2)}\left(x + \frac{\lambda_1}{\lambda_2}\right)^2\right] \\
&= \frac{1}{\sqrt{2\sigma^2\pi}} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right]
\end{aligned}$$

... in this case where mean and variance are known, the entropic prior is Gaussian

An alternative form of entropy that incorporates the normalization constraint

$$\begin{aligned} Q[p; m] &= - \int_{\mathcal{X}} dx p(x) \ln \frac{p(x)}{m(x)} + \lambda \left(\int_{\mathcal{X}} dx p(x) - \int_{\mathcal{X}} dx m(x) \right) \\ &= \int_{\mathcal{X}} dx \left(-p(x) \ln \frac{p(x)}{m(x)} + \lambda p(x) - \lambda m(x) \right) \end{aligned}$$

$$\delta Q = \int_{\mathcal{X}} \delta p dx \left(-\ln \frac{p(x)}{m(x)} - 1 + \lambda \right) = 0$$

$$p(x) = m(x) \exp(\lambda - 1)$$

$$\int_{\mathcal{X}} dx p(x) = \int_{\mathcal{X}} dx m(x) \exp(\lambda - 1) = \exp(\lambda - 1) \int_{\mathcal{X}} dx m(x) = \exp(\lambda - 1) = 1$$

$$\Rightarrow \lambda = 1$$

$$Q[p; m] = \int_{\mathcal{X}} dx \left(-p(x) \ln \frac{p(x)}{m(x)} + p(x) - m(x) \right)$$

Until now we have emphasized the role of the momenta of the distribution, however other information can be incorporated in the same way in the entropic prior.

A “crystallographic” example (Jaynes, 1968)

Consider a simple version of a crystallographic problem, where a 1-D crystal has atoms at the positions

$$x_j = jL \quad (L = 1, \dots, n)$$

and such that these positions may be occupied by impurities.

From X-ray experiments it has been determined that impurity atoms prefer sites where

$$\cos(kx_j) > 0$$

so that

$$\langle \cos(kx_j) \rangle = 0.3$$

which means that we have the constraint

$$\langle \cos(kx_j) \rangle = \sum_{j=1}^n p_j \cos(kx_j) = 0.3$$

where p_j is the probability that an impurity atom is at site j .

Then the constrained entropy that must be maximized is

$$Q = -\sum_{j=1}^n p_j \ln p_j + \lambda_0 \left(\sum_{j=1}^n p_j - 1 \right) + \lambda_1 \left(\sum_{j=1}^n p_j \cos(kx_j) - 0.3 \right)$$

from which we find the maximization condition

$$\frac{\partial Q}{\partial p_j} = -(\ln p_j + 1) + \lambda_0 + \lambda_1 \cos(kx_j) = 0$$

i.e.,

$$p_j = \exp \left[1 - \lambda_0 - \lambda_1 \cos(kx_j) \right]$$

The rest of the solution proceeds either by approximation or by numerical calculation.

References:

- G. D' Agostini, Rep. Prog. Phys. **66** (2003) 1383
- V. Dose: “Bayes in five days”, lecture notes, Max-Planck Research School on bounded plasmas, Greifswald, may 14-18 2002
- V. Dose, Rep. Prog. Phys. **66** (2003) 1421
- E. T. Jaynes, “Monkeys, Kangaroos and N”, in Maximum-Entropy and Bayesian Methods in Applied Statistics, edited by J. H. Justice, Cambridge Univ. Press, Cambridge, UK, 1986, updated (1996) version at <http://bayes.wustl.edu>
- E. T. Jaynes, “Prior probabilities”, IEEE Transactions On Systems Science and Cybernetics, vol. sec-4, (1968) 227

Information theory

- N. Abramson: “Information Theory and Coding”, McGraw-Hill 1963
- R. M. Gray: “Entropy and Information Theory”, Springer-Verlag 1990