

Introduction to Bayesian Statistics - 4

Edoardo Milotti

Università di Trieste and INFN-Sezione di Trieste

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.

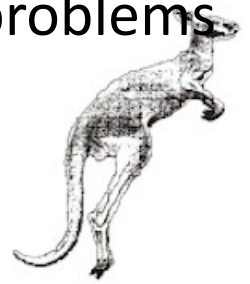
If one considers statistical mechanics as a form of statistical inference rather than as a physical theory, it is found that the usual computational rules, starting with the determination of the partition function, are an immediate consequence of the maximum-entropy principle.

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

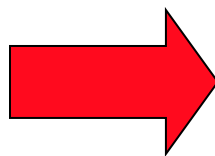
$$S_V = \left(p_{bl} \ln \frac{1}{p_{bl}} + p_{\bar{bl}} \ln \frac{1}{p_{\bar{bl}}} + p_{b\bar{l}} \ln \frac{1}{p_{b\bar{l}}} + p_{\bar{b}\bar{l}} \ln \frac{1}{p_{\bar{b}\bar{l}}} \right) \\ + \lambda_1 (p_{bl} + p_{\bar{bl}} + p_{b\bar{l}} + p_{\bar{b}\bar{l}} - 1) + \lambda_2 (p_{bl} + p_{b\bar{l}}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

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

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

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

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



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

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

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

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

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

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

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

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.

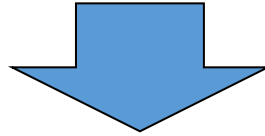
The maximum entropy method helps us find a reasonable solution, the least informative one (least correlations between variables)

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"



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

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

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

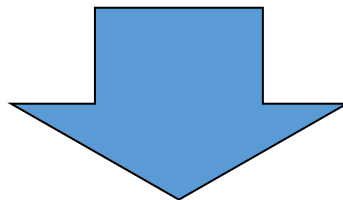
$$\frac{\partial Q}{\partial x} = \frac{(y+z)\ln x + y\ln y + z\ln z}{(x+y+z)^2} + 3\lambda - 2.1\mu = 0$$

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

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

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

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



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

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

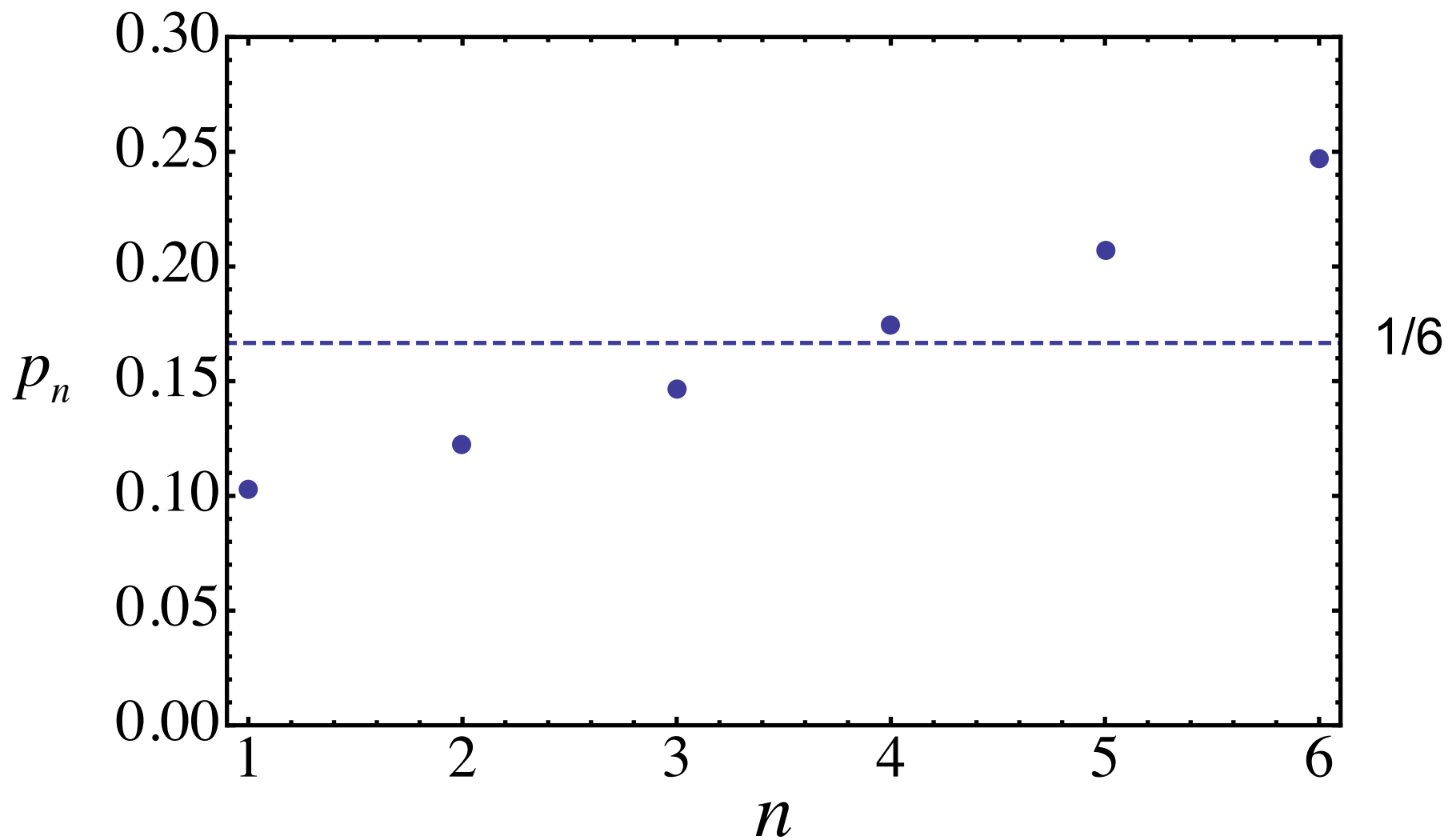
numerical solution

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

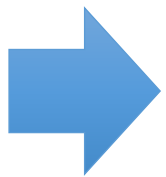
$$S = \int_{-\infty}^{+\infty} p(x) dx \log_2 \frac{1}{p(x) dx}$$

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 reference 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

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

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



$$I_{KL}(p, q) = \frac{\mu^2}{2\sigma^2}$$

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

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

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

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

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

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

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

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

we find an approximate solution

$$\begin{aligned} z - \frac{z^3}{3} &\approx \left(\frac{1}{z} + \frac{2\varepsilon}{L} \right)^{-1} \Rightarrow \left(z - \frac{z^3}{3} \right) \left(\frac{1}{z} + \frac{2\varepsilon}{L} \right) \approx 1 + \frac{2\varepsilon}{L} z - \frac{z^2}{3} = 1 \\ \Rightarrow \frac{2\varepsilon}{L} z - \frac{z^2}{3} &\approx 0 \Rightarrow z \approx \frac{6\varepsilon}{L} \end{aligned}$$

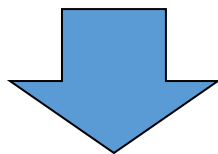
$$\frac{\lambda_1 L}{2} \approx \frac{6\varepsilon}{L} \Rightarrow p(x) \approx \frac{1}{L} \exp(\lambda_1 x) \approx \frac{1}{L} \left(1 - \frac{12\varepsilon}{L} x \right)$$

another special case

$$a = 0; \quad b \rightarrow \infty$$

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

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



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

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

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.

Example of MaxEnt in action: unconstrained problem in image restoration



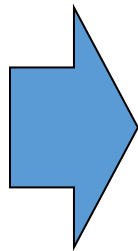
J. Skilling, Nature 309 (1984) 748

Car movement introduces linear correlations among pixels. The model of linear corrections does not allow direct inversion to find the corrected image because the number of variables is larger than the number of equations. The MaxEnt methods regularizes the problem and finds a reasonable solution.



J. Skilling, Nature 309 (1984) 748

Reconstruction of missing data (from <http://www.maxent.co.uk>)



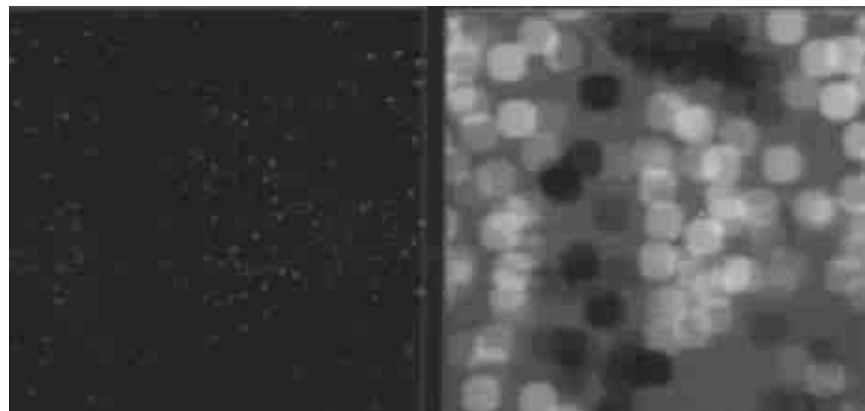
50%

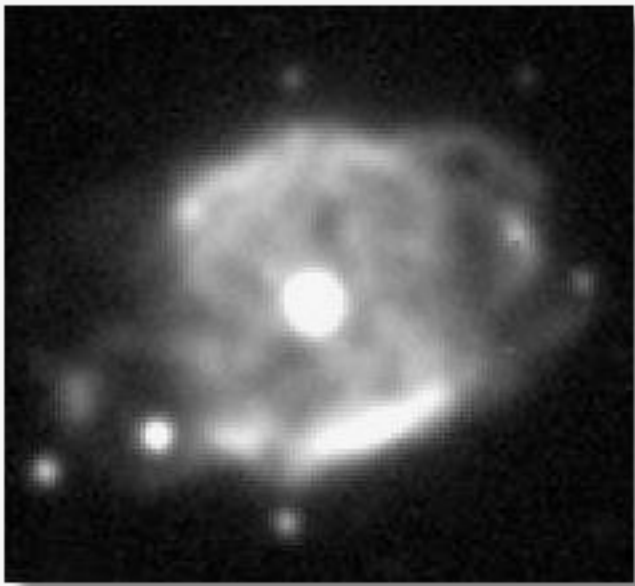


95%

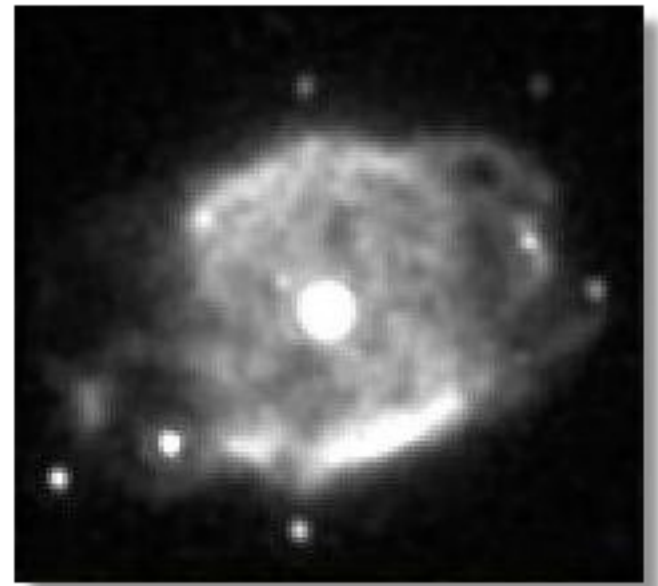


99%





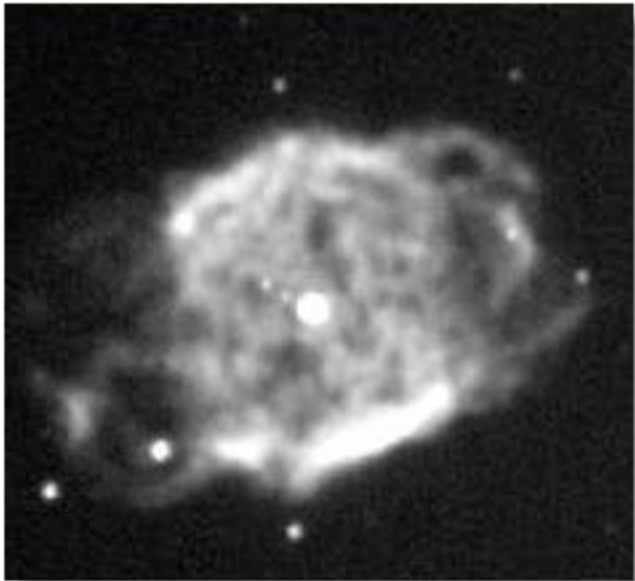
NGC 40



low resolution (MEM enhanced)

low resolution

high resolution



John Skilling: Biographical information

John is Scientific Director of MEDC. He did his Ph.D. (on cosmic rays) in the Department of Physics at Cambridge University, and went on to become a Lecturer in the Department of Applied Mathematics and Theoretical Physics, and a Fellow of St Johns College.

In the late 1970s, another radio astronomer, [Steve Gull](#), introduced him to the power of the Maximum Entropy Method. John wrote what was to become the first MemSys kernel system, and helped lay the Bayesian foundations for MEM. In 1981 he and Steve founded MEDC to exploit opportunities to apply MEM in other fields.

John resigned his Lectureship in 1990 in order to go fulltime with MSL and MEDC. Thanks to the wonders of modern technology John is able to telecommute from his new home in the West of Ireland, and he makes regular visits to clients both in the UK and further afield.



[Home](#) | [Applications](#) | [Products](#) | [Prices](#) | [Documents](#) | [About MEDC](#) | [Contact Us](#) | [Full search](#)

[Home](#)

[About MEDC](#)

[Applications](#)

[Examples](#)

[Products](#)

[Prices](#)

[Documents](#)

[Contact us](#)

[Search MEDC](#)

Quick Search:

Search

©MEDC Ltd. Last revised Wed Sep 19 22:19:39 2007

<http://www.maxent.co.uk/>

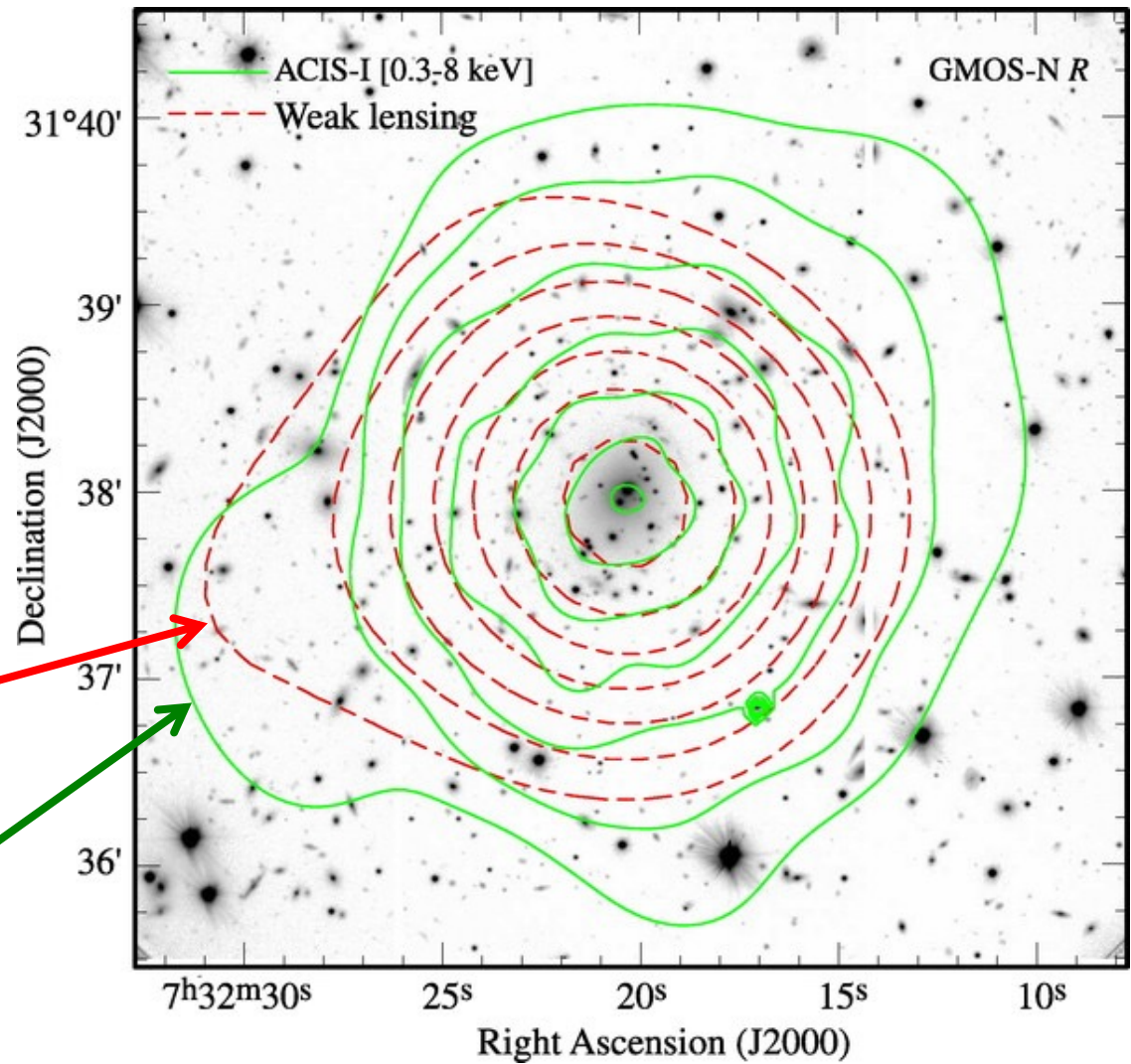
(the company still exists but the website has disappeared from the web)

Example of LensEnt usage
(Bridle et al, 1998):

reconstruction of mass
density from lensing data,
using Max Ent

reconstructed mass
density

X-ray emission data



GMOS image of the central region of Abell 586 with **logarithmically spaced X-ray isophotes (solid lines) and weak-lensing reconstructed mass density (dashed lines) superposed**. The X-ray point source near the southwest corner is the Seyfert 1 galaxy C171_3650.

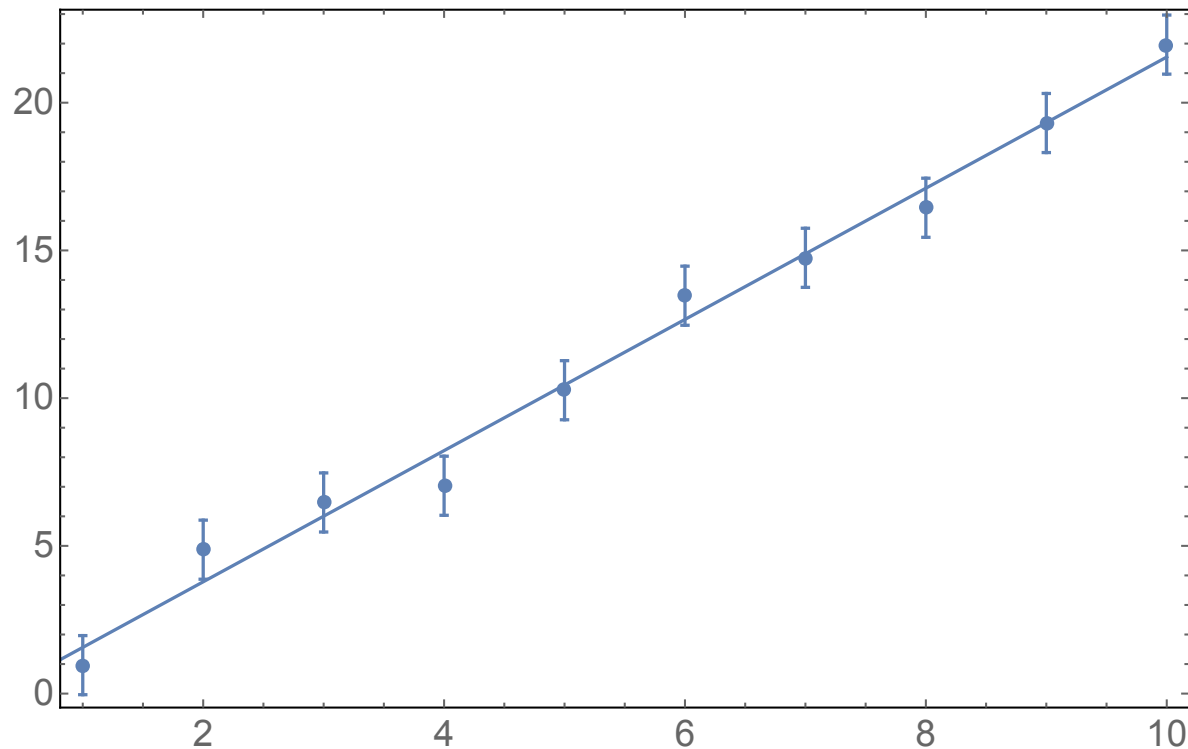
(from Cypriano et al., ApJ, **630** (2005) 38-49)

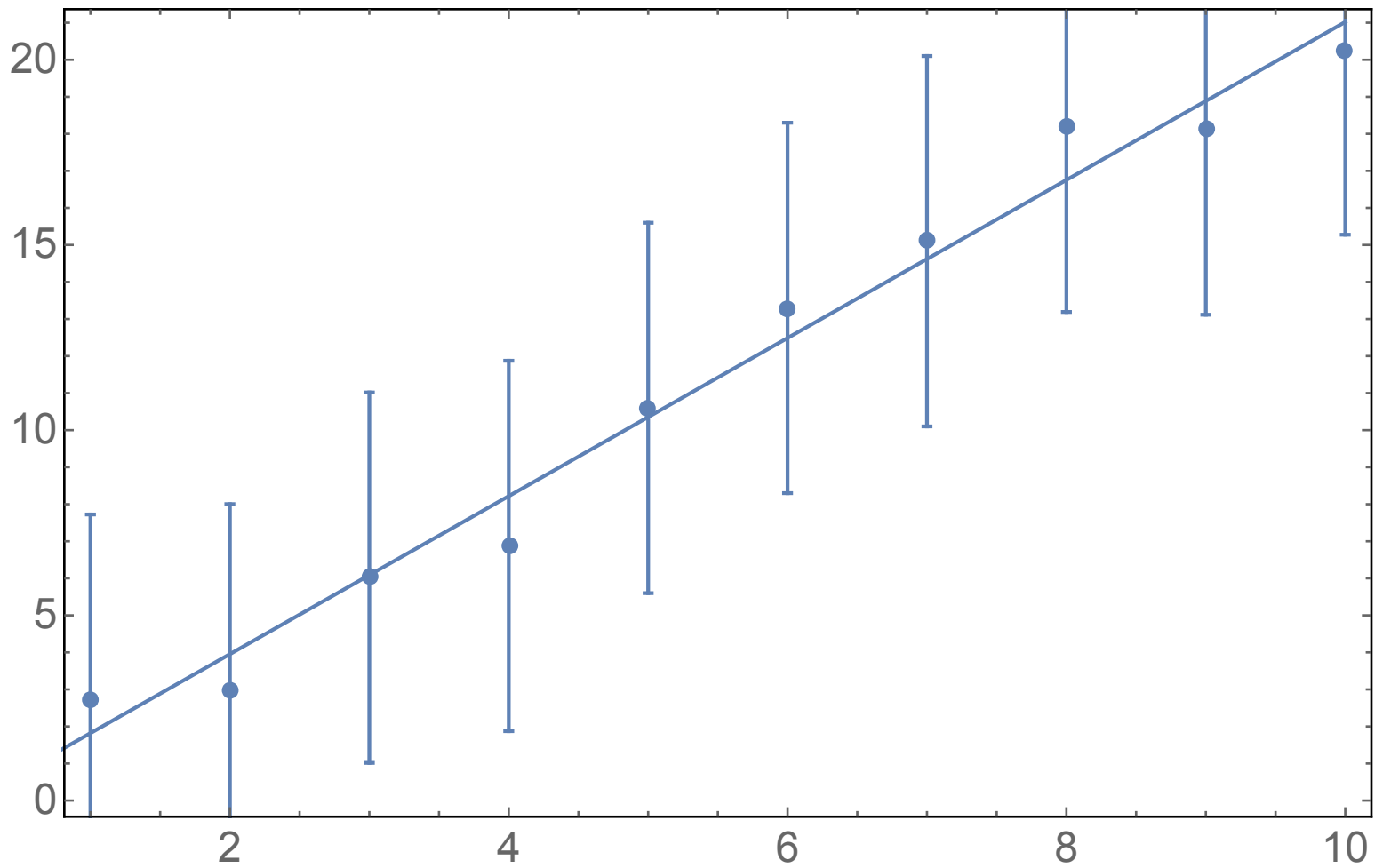
A few Bayesian examples

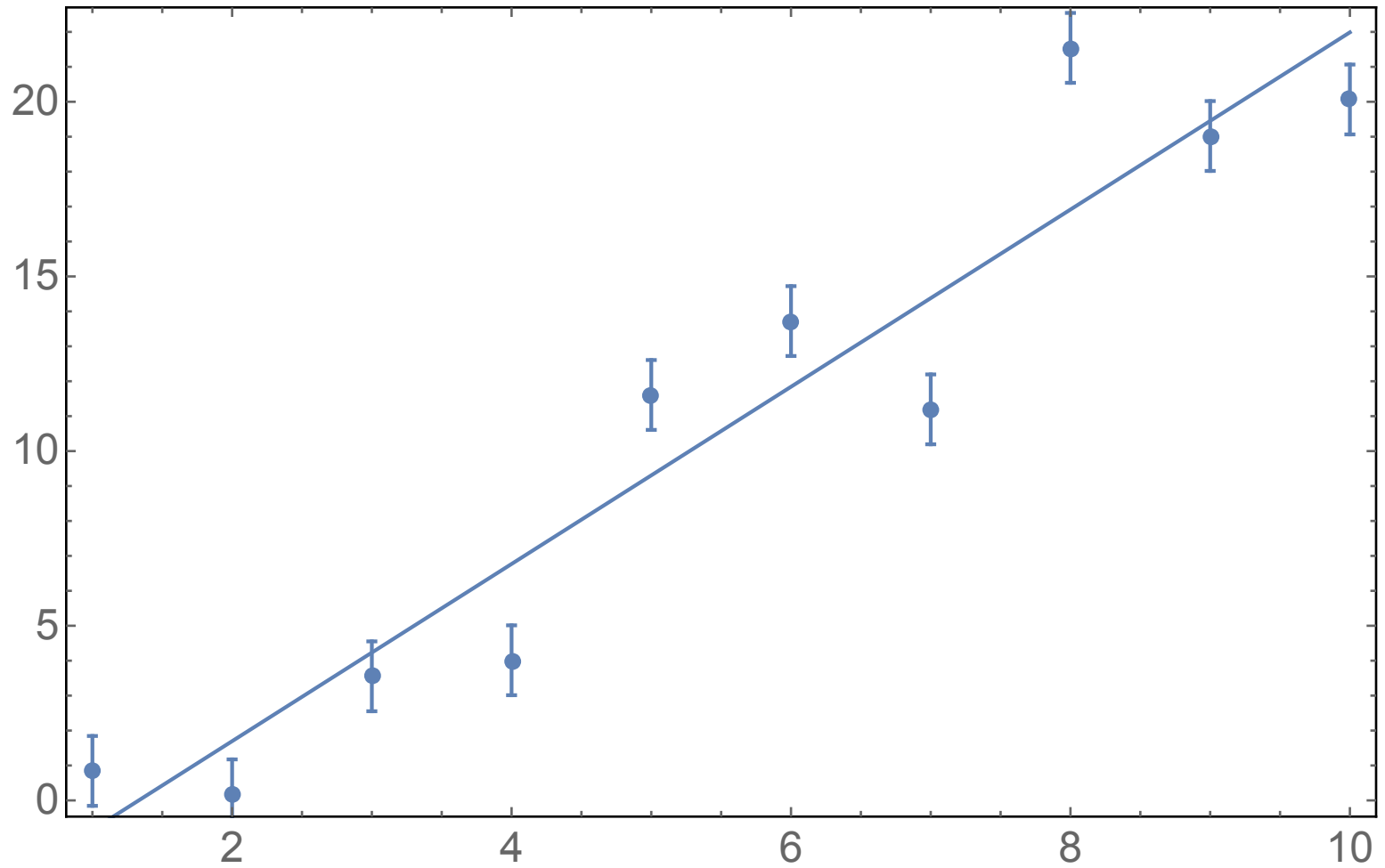
1. Miscalibrated Gaussian measurement errors
2. Search for weak signals in spectra
3. The statistical link between smoking and lung cancer

2. Miscalibrated Gaussian measurement errors, a Bayesian estimate using objective priors

Here, we consider the case where we must find the mean value with given measurement errors, and where the errors are Gaussian and they are systematically multiplied by an unknown scale factor.







The likelihood has a Gaussian structure

$$\begin{aligned} P(\mathbf{d} \mid \mu, \boldsymbol{\sigma}, \alpha) &= \prod_{k=1}^N \frac{1}{\sqrt{2\pi\alpha^2\sigma_k^2}} \exp\left[-\frac{(d_k - \mu)^2}{2\alpha^2\sigma_k^2}\right] \\ &= \frac{1}{(2\pi)^{N/2} \alpha^N} \left(\prod_{k=1}^N \frac{1}{\sigma_k} \right) \exp\left[-\frac{1}{2\alpha^2} \sum_{k=1}^N \frac{(d_k - \mu)^2}{\sigma_k^2}\right] \end{aligned}$$

we must rearrange the exponent as usual ...

$$\begin{aligned}\sum_{k=1}^N \frac{(d_k - \mu)^2}{\sigma_k^2} &= \sum_{k=1}^N \frac{d_k^2}{\sigma_k^2} - 2\mu \sum_{k=1}^N \frac{d_k}{\sigma_k^2} + \mu^2 \sum_{k=1}^N \frac{1}{\sigma_k^2} = \frac{ND}{\sigma_M^2} - 2\mu \frac{NM}{\sigma_M^2} + \mu^2 \frac{N}{\sigma_M^2} \\ &= \frac{N}{\sigma_M^2} (D - 2\mu M + \mu^2)\end{aligned}$$

$$\text{dove } \frac{1}{\sigma_M^2} = \frac{1}{N} \sum_{k=1}^N \frac{1}{\sigma_k^2}; \quad M = \sum_{k=1}^N \frac{d_k}{\sigma_k^2} / \sum_{k=1}^N \frac{1}{\sigma_k^2}; \quad D = \sum_{k=1}^N \frac{d_k^2}{\sigma_k^2} / \sum_{k=1}^N \frac{1}{\sigma_k^2}$$

therefore the likelihood is

$$P(\mathbf{d} \mid \mu, \boldsymbol{\sigma}, \alpha) = \frac{1}{(2\pi)^{N/2} \alpha^N} \left(\prod_{k=1}^N \frac{1}{\sigma_k} \right) \exp \left[-\frac{N}{2\alpha^2 \sigma_M^2} (D - 2\mu M + \mu^2) \right]$$

Now we estimate the scale factor from Bayes' theorem

$$p(\alpha|\mathbf{d}, \boldsymbol{\sigma}) = \frac{p(\mathbf{d}|\alpha, \boldsymbol{\sigma})}{\int_{\alpha} p(\mathbf{d}|\alpha', \boldsymbol{\sigma})p(\alpha')d\alpha'}p(\alpha)$$

however we need first to marginalize the likelihood with respect to the mean, which in this case is a *nuisance parameter*

we take a uniform prior for the mean

$$\begin{aligned} P(\mathbf{d}|\boldsymbol{\sigma}, \alpha) &= \int_{\mu} P(\mathbf{d}|\mu, \boldsymbol{\sigma}, \alpha)P(\mu|\boldsymbol{\sigma}, \alpha)d\mu \\ &= \frac{1}{W} \int_{\mu_{\min}}^{\mu_{\max}} P(\mathbf{d}|\mu, \boldsymbol{\sigma}, \alpha)d\mu \\ &\approx \frac{1}{W} \frac{1}{(2\pi)^{N/2} \alpha^N} \left(\prod_{k=1}^N \frac{1}{\sigma_k} \right) \int_{-\infty}^{+\infty} \exp \left[-\frac{N}{2\alpha^2 \sigma_M^2} (D - 2\mu M + \mu^2) \right] d\mu \end{aligned}$$

$$(W = \mu_{\max} - \mu_{\min})$$

as usual ...

$$\begin{aligned} D - 2\mu M + \mu^2 &= \mu^2 - 2\mu M + M^2 + D - M^2 \\ &= (\mu - M)^2 + D - M^2 \end{aligned}$$

... therefore the likelihood is:

$$\begin{aligned} P(\mathbf{d} | \boldsymbol{\sigma}, \alpha) &\approx \frac{1}{W} \frac{1}{(2\pi)^{N/2} \alpha^N} \left(\prod_{k=1}^N \frac{1}{\sigma_k} \right) \int_{-\infty}^{+\infty} \exp \left\{ -\frac{N}{2\alpha^2 \sigma_M^2} \left[(\mu - M)^2 + D - M^2 \right] \right\} d\mu \\ &= \frac{1}{W} \frac{1}{(2\pi)^{N/2} \alpha^N} \left(\prod_{k=1}^N \frac{1}{\sigma_k} \right) \exp \left(-\frac{N(D - M^2)}{2\alpha^2 \sigma_M^2} \right) \sqrt{\frac{2\pi\alpha^2 \sigma_M^2}{N}} \end{aligned}$$

$$\begin{aligned}
p(\alpha|\mathbf{d}, \boldsymbol{\sigma}) &= \frac{p(\mathbf{d}|\alpha, \boldsymbol{\sigma})}{\int_{\alpha} p(\mathbf{d}|\alpha', \boldsymbol{\sigma})p(\alpha')d\alpha'}p(\alpha) \\
&= \frac{\frac{1}{\alpha^{N-1}} \exp\left(-\frac{N(D - M^2)}{2\alpha^2\sigma_M^2}\right)}{\int_{\alpha} \frac{1}{\alpha'^{N-1}} \exp\left(-\frac{N(D - M^2)}{2\alpha'^2\sigma_M^2}\right) p(\alpha')d\alpha'}p(\alpha)
\end{aligned}$$

$P(\alpha) \propto \frac{1}{\alpha}$ for the standard deviation we take again the Jeffreys prior

$$p(\alpha|\mathbf{d}, \boldsymbol{\sigma}) = \frac{\frac{1}{\alpha^{N-1}} \exp\left(-\frac{N(D - M^2)}{2\alpha^2\sigma_M^2}\right) \frac{1}{\alpha}}{\int_{\alpha} \frac{1}{\alpha'^{N-1}} \exp\left(-\frac{N(D - M^2)}{2\alpha'^2\sigma_M^2}\right) \frac{1}{\alpha'} d\alpha'}; \quad A^2 = \frac{N(D - M^2)}{2\sigma_M^2}$$

$$p(\alpha|\mathbf{d}, \boldsymbol{\sigma}) \rightarrow \frac{\frac{1}{\alpha^N} \exp\left(-\frac{N(D - M^2)}{2\alpha^2\sigma_M^2}\right)}{\int_0^{\infty} \frac{1}{\alpha'^N} \exp\left(-\frac{N(D - M^2)}{2\alpha'^2\sigma_M^2}\right) d\alpha'}$$

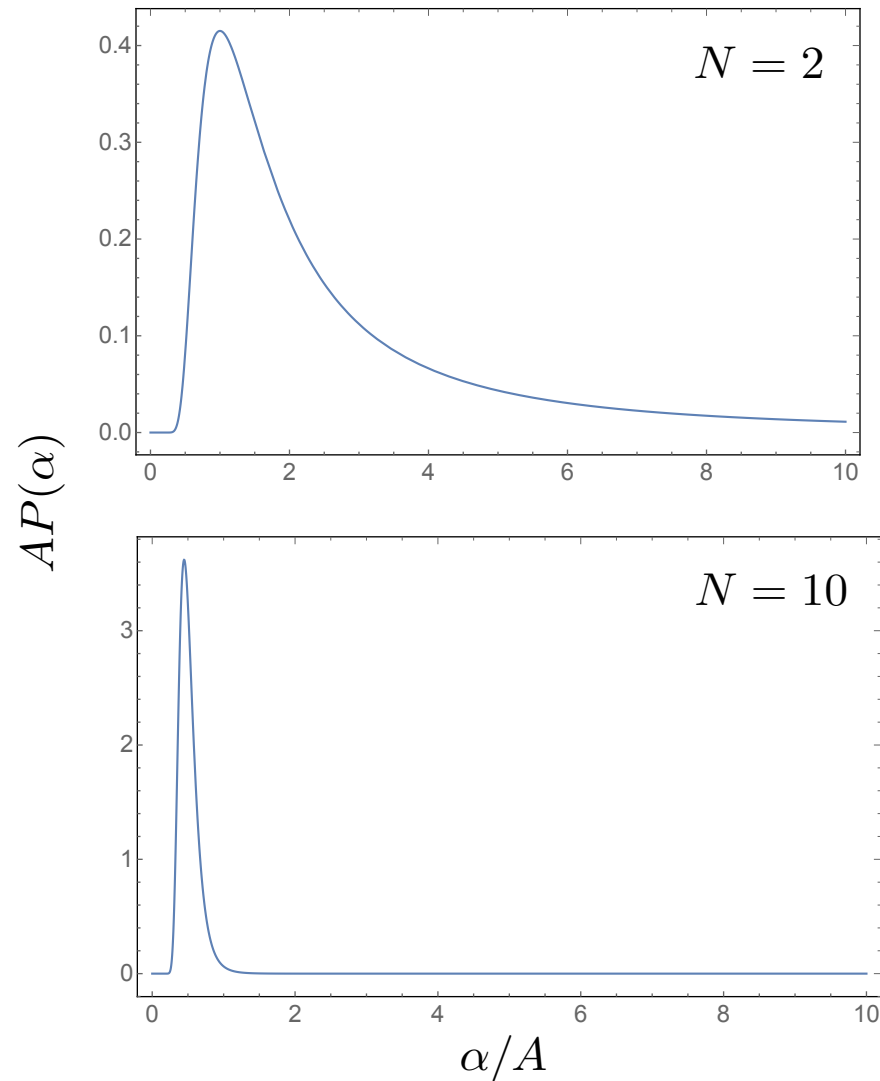
evaluation of $\int_0^\infty \frac{1}{\alpha'^N} \exp\left(-\frac{N(D - M^2)}{2\alpha'^2 \sigma_M^2}\right) d\alpha'$

$$\frac{A^2}{\alpha^2} = x; \quad \alpha = \frac{A}{\sqrt{x}}; \quad d\alpha = -\frac{A}{2x^{3/2}} dx$$

$$\int_0^\infty \frac{x^{N/2}}{A^N} \exp(-x) \frac{A}{2x^{3/2}} dx = \frac{1}{2A^{N-1}} \int_0^\infty x^{\frac{N-1}{2}-1} \exp(-x) dx = \frac{1}{2A^{N-1}} \Gamma\left(\frac{N-1}{2}\right)$$

$$p(\alpha | \mathbf{d}, \boldsymbol{\sigma}) \rightarrow \frac{\frac{2A^{N-1}}{\alpha^N} \exp\left(-\frac{A^2}{\alpha^2}\right)}{\Gamma\left(\frac{N-1}{2}\right)}$$

$$P(\alpha|\mathbf{d}, \boldsymbol{\sigma}) = \frac{(2A^{N-1}/\alpha^N) \exp(-A^2/\alpha^2)}{\Gamma[(N-1)/2]}$$



we take the MAP estimate of the scale parameter from the pdf

$$P(\alpha | \mathbf{d}, \boldsymbol{\sigma}) \rightarrow \frac{2A^{N-1} \exp\left(-\frac{A^2}{\alpha^2}\right)}{\Gamma\left(\frac{N-1}{2}\right)}$$

$$\frac{d}{d\alpha} P(\alpha | \mathbf{d}, \boldsymbol{\sigma}) \propto -\frac{N}{\alpha^{N+1}} \exp\left(-\frac{A^2}{\alpha^2}\right) + \frac{2A^2}{\alpha^{N+3}} \exp\left(-\frac{A^2}{\alpha^2}\right) = 0$$

$$\begin{array}{ccc} \rightarrow & N\alpha^2 = 2A^2 & \rightarrow \alpha_{MAP} = \sqrt{\frac{2}{N}}A = \sqrt{\frac{(D-M^2)}{\sigma_M^2}} \end{array}$$

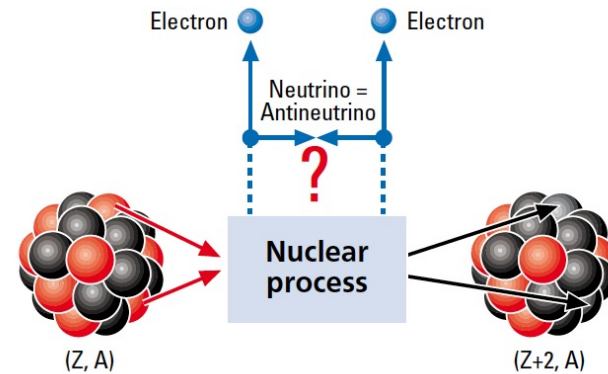
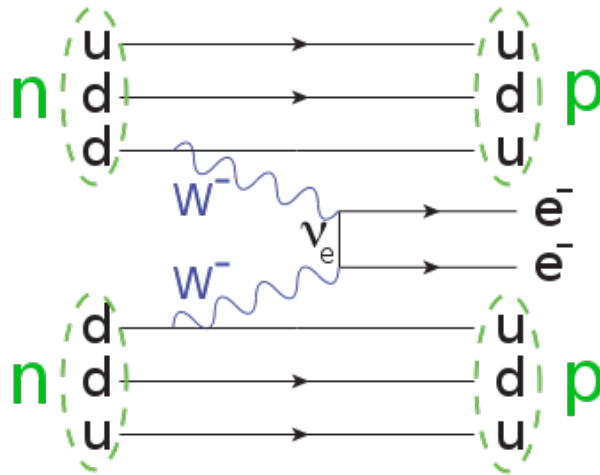
2. Search of signals in binned spectra: Bayesian analysis in the GERDA experiment

(Caldwell and Kröninger, PRD 74 (2006) 092003)

Consider the search for sparse signals in a spectrum where

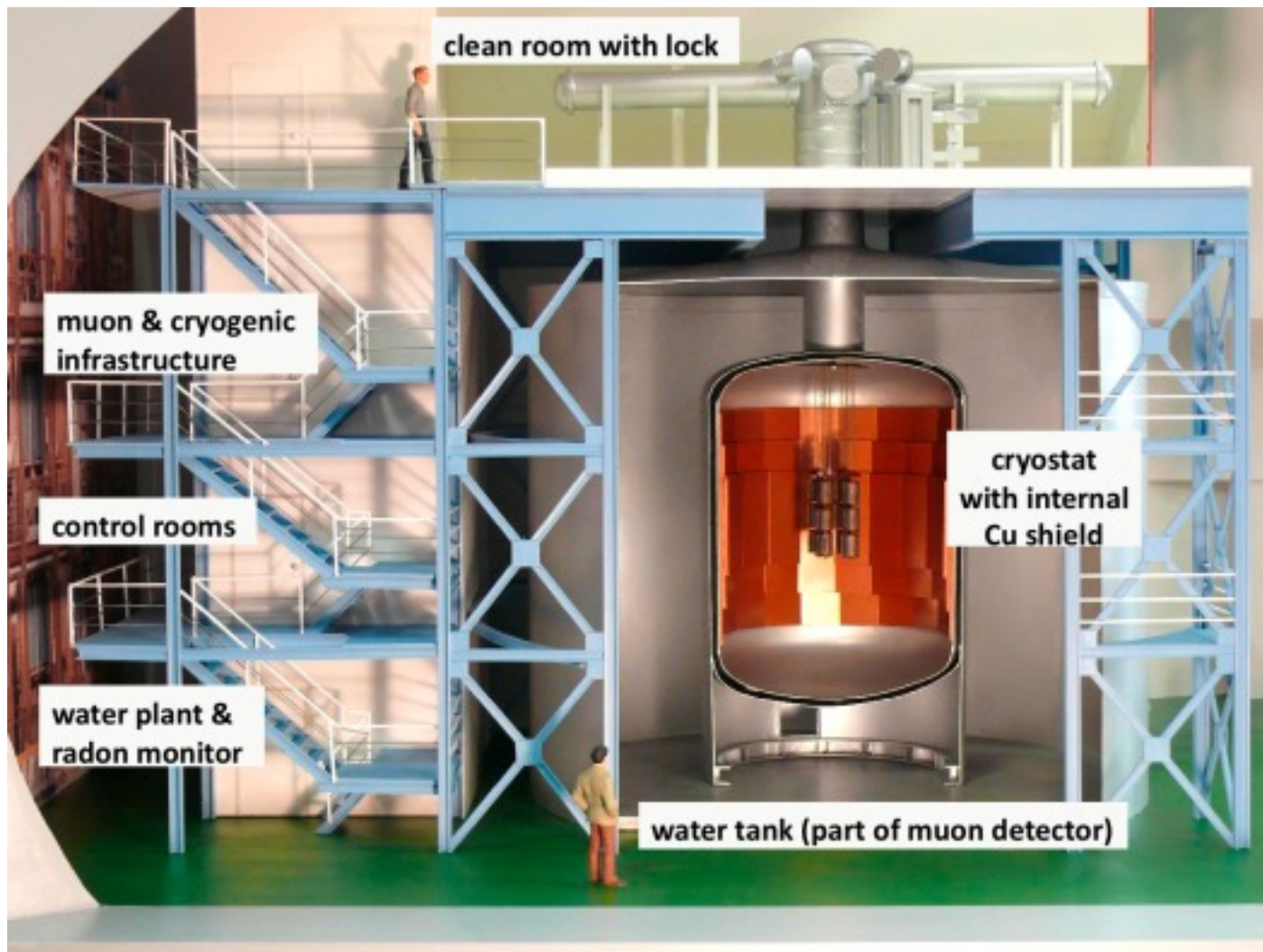
- The spectrum is confined to a certain region of interest.
- The spectral shape of a possible signal is known.
- The spectral shape of the background is known.
- The spectrum is divided into bins and the event numbers in the bins follow Poisson distributions.

The work of Caldwell and Kröniger has been carried out in the context of GERDA (GERmanium Detector Array), an experiment that aims to detect weak signals from neutrinoless beta decay in germanium detectors kept in a very low background environment.



*Feynman diagram of neutrinoless double-beta decay, with two neutrons decaying to two protons. **The only emitted products in this process are two electrons, which can occur if the neutrino and antineutrino are the same particle** (i.e. Majorana neutrinos) so the same neutrino can be emitted and absorbed within the nucleus.*

In conventional double-beta decay, two antineutrinos - one arising from each W vertex - are emitted from the nucleus, in addition to the two electrons. The detection of neutrinoless double-beta decay is thus a sensitive test of whether neutrinos are Majorana particles. (from Wikipedia)



Artist's view of GERDA
<http://www.mpi-hd.mpg.de/gerda/>

We introduce the normalized spectral shapes of background and signal

$$f_B(E); \quad f_S(E);$$

(flat spectrum, known signal shape). Then we find the average number of events in each bin

$$v_i(B, S) = v_i(E_i, \Delta E_i, B, S) = B \int_{E_i}^{E_i + \Delta E_i} f_B(E) dE + S \int_{E_i}^{E_i + \Delta E_i} f_S(E) dE$$

An observed spectrum is defined by the numbers of counts in each bin: $\{n_i\}_{i=1, n}$ and since we assume a Poisson statistics in each bin, we find the following likelihoods for a given spectral observation

$$p(\text{spectrum} | B, I) = \prod_{i=1}^N \frac{[v_i(B, 0)]^{n_i}}{n_i!} \exp[-v_i(B, 0)]$$

$$p(\text{spectrum} | B, S, I) = \prod_{i=1}^N \frac{[v_i(B, S)]^{n_i}}{n_i!} \exp[-v_i(B, S)]$$

A specific spectral shape depends on the **average number of background (B) and signal (S) events**, and we can write

$$p(\text{spectrum} | H_{\text{bkg}}, I) = \int_B p(\text{spectrum} | B, I) p_B(B) dB$$

$$p(\text{spectrum} | H_{\text{bs}}, I) = \int_{B,S} p(\text{spectrum} | B, S, I) p_B(B) p_S(S) dB dS$$

distribution for the average B



distribution for the average S



the possible spectra are the results of many possible choices of the background and of the signal rates, and therefore of the average number of background and signal events; here we marginalize over these dependencies

The competing hypotheses (observation of binned energy spectra) are

- H_{bkg} = background only
- H_{bs} = background + signal

then

$$p(H_{bkg} | spectrum, I) = \frac{p(spectrum | H_{bkg}, I)}{p(spectrum | I)} p(H_{bkg} | I)$$

$$p(H_{bs} | spectrum, I) = \frac{p(spectrum | H_{bs}, I)}{p(spectrum | I)} p(H_{bs} | I)$$

Bayes

$$p(spectrum | I) = p(spectrum | H_{bkg}, I) p(H_{bkg} | I) + p(spectrum | H_{bs}, I) p(H_{bs} | I)$$

Evidence

Then we find the complete likelihood functions:

$$\begin{aligned} p(\text{spectrum} | H_{bkg}, I) &= \int_B p(\text{spectrum} | B, I) p_B(B) dB \\ &= \int_B \prod_{i=1}^N \frac{[v_i(B, 0)]^{n_i}}{n_i!} \exp[-v_i(B, 0)] p_B(B) dB \end{aligned}$$

$$\begin{aligned} p(\text{spectrum} | H_{bs}, I) &= \int_B p(\text{spectrum} | B, S, I) p_B(B) p_S(S) dB \\ &= \int_B \prod_{i=1}^N \frac{[v_i(B, S)]^{n_i}}{n_i!} \exp[-v_i(B, S)] p_B(B) p_S(S) dB \end{aligned}$$

The final, complete expressions are:

$$\begin{aligned}
 p(H_{bkg} | spectrum, I) &= \frac{p(spectrum | H_{bkg}, I)}{p(spectrum | I)} p(H_{bkg} | I) \\
 &= \frac{\int \prod_{B, i=1}^N \frac{[v_i(B, 0)]^{n_i}}{n_i!} \exp[-v_i(B, 0)] p_B(B) dB}{\int \prod_{B, i=1}^N \frac{[v_i(B, 0)]^{n_i}}{n_i!} \exp[-v_i(B, 0)] p_B(B) dB p(H_{bkg} | I) + \int \prod_{B, i=1}^N \frac{[v_i(B, S)]^{n_i}}{n_i!} \exp[-v_i(B, S)] p_B(B) p_S(S) dB p(H_{bs} | I)} p(H_{bkg} | I)
 \end{aligned}$$

$$\begin{aligned}
 p(H_{bs} | spectrum, I) &= \frac{p(spectrum | H_{bs}, I)}{p(spectrum | I)} p(H_{bs} | I) \\
 &= \frac{\int \prod_{B, i=1}^N \frac{[v_i(B, S)]^{n_i}}{n_i!} \exp[-v_i(B, S)] p_B(B) p_S(S) dB}{\int \prod_{B, i=1}^N \frac{[v_i(B, 0)]^{n_i}}{n_i!} \exp[-v_i(B, 0)] p_B(B) dB p(H_{bkg} | I) + \int \prod_{B, i=1}^N \frac{[v_i(B, S)]^{n_i}}{n_i!} \exp[-v_i(B, S)] p_B(B) p_S(S) dB p(H_{bs} | I)} p(H_{bs} | I)
 \end{aligned}$$

One can use these expressions to test hypotheses (by means of Bayes factors), and find values for B and S.

3. The statistical link between smoking and lung cancer

Cornfield, Jerome

Born: October 30, 1912, in New York City, New York.

Died: September 17, 1979, in Herndon, Virginia.



A METHOD OF ESTIMATING COMPARATIVE RATES FROM CLINICAL DATA. APPLICATIONS TO CANCER OF THE LUNG, BREAST, AND CERVIX ¹

JEROME CORNFIELD, *National Cancer Institute, National Institutes of Health, U. S. Public Health Service, Bethesda, Md.*

¹ Received for publication February 23, 1961.



FIGURE 1. Passport photograph of Ronald Aylmer Fisher at age 34. Reprinted from Box JF. *RA Fisher: the life of a scientist*. New York: John Wiley & Sons, Inc., 1978.

Fisher developed four lines of argument in questioning the causal relation of lung cancer to smoking.

- 1) If A is associated with B, then not only is it possible that A causes B, but it is also possible that B is the cause of A. In other words, smoking may cause lung cancer, but it is a logical possibility that lung cancer causes smoking.
- 2) There may be a genetic predisposition to smoke (and that genetic predisposition is presumably also linked to lung cancer).
- 3) Smoking is unlikely to cause lung cancer because secular trend and other ecologic data do not support this relation.
- 4) Smoking does not cause lung cancer because inhalers are less likely to develop lung cancer than are noninhalers

Lung cancer and cigarette smoking

Consider the following data for fractions of the population (Cornfield, 1951)

	Having cancer of the lung	Healthy	Total
Smokers	$0.119 \cdot 10^{-3}$	0.579910	0.580025
Nonsmokers	$0.036 \cdot 10^{-3}$	0.419935	0.419971
Total	$0.155 \cdot 10^{-3}$	0.999845	1.000000

what is the proportion having cancer of the lung in each population?

Smokers: $0.119 \cdot 10^{-3} / 0.580025 = 2.05164 \cdot 10^{-4}$

Nonsmokers: $0.036 \cdot 10^{-3} / 0.419971 = 8.57202 \cdot 10^{-5}$

And the prevalence of lung cancer in smokers with respect to nonsmokers is

$$\text{Smokers/Nonsmokers} \approx 2.4$$

We can also write an easy Bayesian equation that leads to some information as to the **causation** of cancer of the lung

$$P(\text{Cancer}|\text{Smoker}) = \frac{P(\text{Smoker}|\text{Cancer})P(\text{Smoker})}{P(\text{Cancer})}$$

$$P(\text{Cancer}|\text{Nonsmoker}) = \frac{P(\text{Nonsmoker}|\text{Cancer})P(\text{Nonsmoker})}{P(\text{Cancer})}$$

Therefore

$$\frac{P(\text{Cancer}|\text{Smoker})}{P(\text{Cancer}|\text{Nonsmoker})} = \frac{P(\text{Smoker}|\text{Cancer})P(\text{Smoker})}{P(\text{Nonsmoker}|\text{Cancer})P(\text{Nonsmoker})}$$

and with the numbers in the table one finds that this ratio is about 3.5.

According to Jeffreys, a Bayes ratio of 3.5 is already substantial support in favor of the hypothesis that smoking does cause lung cancer.

$\log_{10}(B)$	B	Evidence support
0 to $1/2$	1 to 3.2	Not worth more than a bare mention
$1/2$ to 1	3.2 to 10	Substantial
1 to 2	10 to 100	Strong
> 2	> 100	Decisive

In 1954 Richard Doll and Bradford Hill published evidence in the British Medical Journal showing a strong link between smoking and lung cancer. They published further evidence in 1956.

Fisher was a paid tobacco industry consultant and a devoted pipe smoker. He did not think the statistical evidence for a link was convincing.

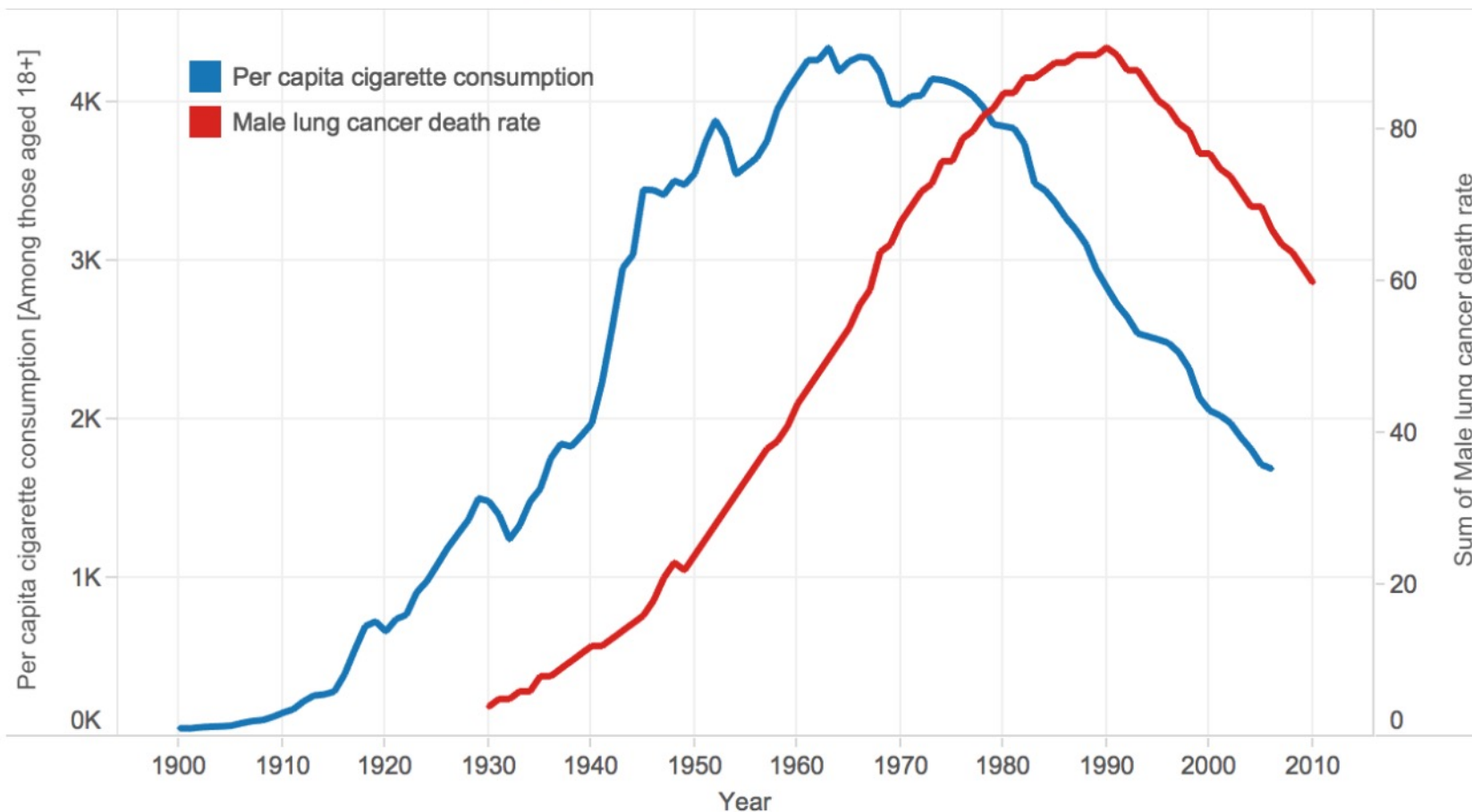
Ronald Fisher died aged 72 on July 29, 1962, in Adelaide, Australia following an operation for colon cancer.

With bitter irony, we now know that the likelihood of getting this disease increases in smokers.

Ronald Fisher was cremated and his ashes interred in St. Peter's Cathedral, Adelaide.

(from "Ronald Fisher." Famous Scientists. famousscientists.org. 17 Sep. 2015. Web. 5/30/2017 <www.famousscientists.org/ronald-fisher/>.)

Trends in Tobacco Use and Lung Cancer Death Rates in the U.S.



Death rates source: US Mortality Data, 1960-2010, US Mortality Volumes, 1930-1959, National Center for Health Statistics, Centers for Disease Control and Prevention.

Cigarette consumption source: US Department of Agriculture, 1900-2007.

From 1948 to his death 31 years later, Cornfield devoted the major portion of his career to the development and application of statistical theory and methods to the biomedical sciences. His contributions were diverse both in the nature of his statistical interests and in the areas of biostatistical applications. He was involved in and touched upon every major public health issue that arose in that period – the polio vaccines [23], smoking and lung cancer (*see **Smoking and Health***) [22, 29], risk factors for cardiovascular disease [5, 30], and the difficult statistical issues of estimating the low-dose carcinogenic effects in humans (*see **Extrapolation, Low Dose***) of a food additive that becomes suspect because it produces cancer in animals at much higher doses [14, 20].

Encyclopedia of Biostatistics, Online © 2005 John Wiley & Sons, Ltd.

This article is © 2005 John Wiley & Sons, Ltd.

This article was published in the *Encyclopedia of Biostatistics* in 2005 by John Wiley & Sons, Ltd.

DOI: 10.1002/0470011815.b2a17032

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