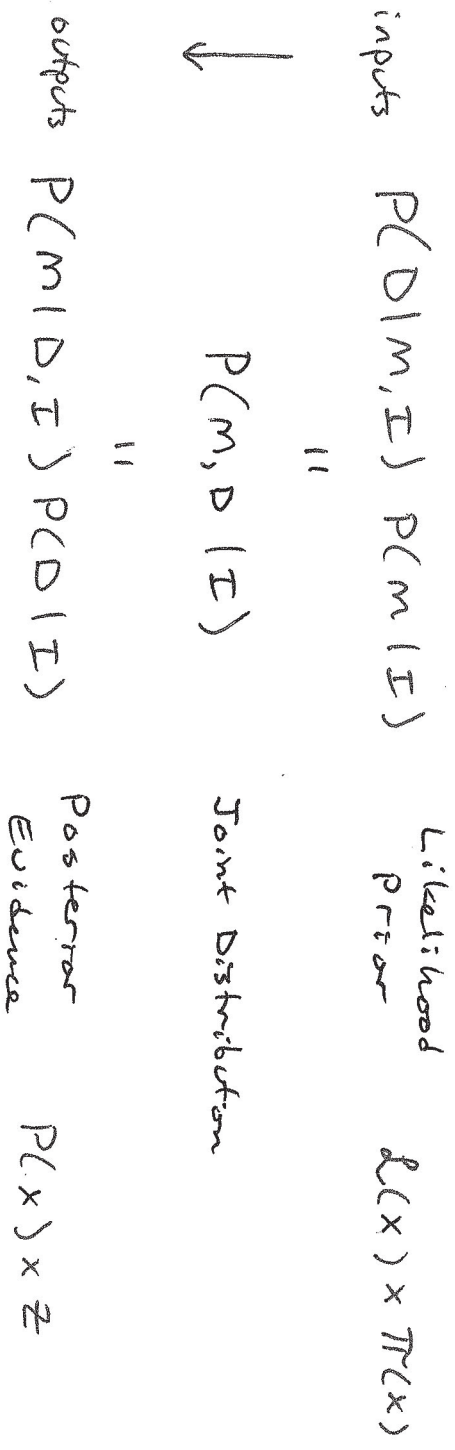


Nested Sampling

Skilling 2005

So far we have focused on the posterior prob.

This method derives from another perspective.



Normalization implies that

$$\int P(M|I) dM = 1 \quad \text{and} \quad \int P(M|D, I) dM = 1$$

Together, this implies that

$$Z = \int P(D|M, I) P(M|I) dM = 1$$

$$Z = \int \mathcal{L}(x) \pi(x) dx$$

$$P(x) = \frac{\mathcal{L}(x) \pi(x)}{Z}$$

Nested Sampling: Sorted Likelihood

Consider a discrete 2D hypothesis space.

We assume a uniform prior and plot a map of the likelihood below

$$\mathcal{L} =$$

0	8	15	3
11	24	22	10
19	30	26	16
9	23	18	6

and sort the values into a vector

$$\mathcal{L} = (30, 26, 24, 23, 22, 19, 18, 16, 15, 11, 10, 9, 8, 6, 3, 0)$$

with

$$\pi = \frac{1}{16} \times (1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1)$$

We now define the function

$$\xi(\lambda) = \text{proportion of the prior probability mass with likelihood greater than } \lambda$$

Formally, this is

$$\xi(\lambda) = \int \pi(x) dx$$

$$\mathcal{L}(x) \geq \lambda$$

The element of prior mass is $d\mathcal{L} = \pi(x) dx$

Nested Sampling: Density of States

Select $\lambda = 15$

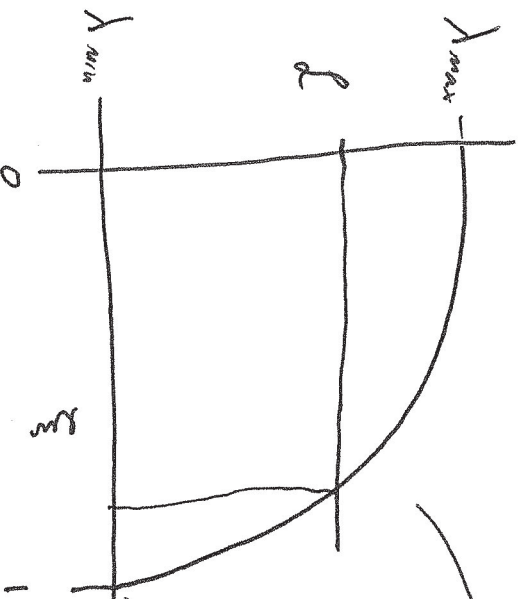
Then $\xi(\lambda=15) = 8/16 = 0.5$

We can look at the inverse of this function

$$\lambda = \xi^{-1}(\xi(x))$$

and since this function returns a likelihood, Skilling calls $\xi^{-1} \equiv \mathcal{L}$

where $\mathcal{L}(\xi)$ is distinct from $\mathcal{L}(x)$



This is a "density of states"
It describes how the states of
the system are distributed with
respect to their probability

The entire volume $\xi = 1$
includes all λ values $> \lambda_{min}$

$$\text{Now } Z = \int \mathcal{L}(x) \pi(x) dx$$

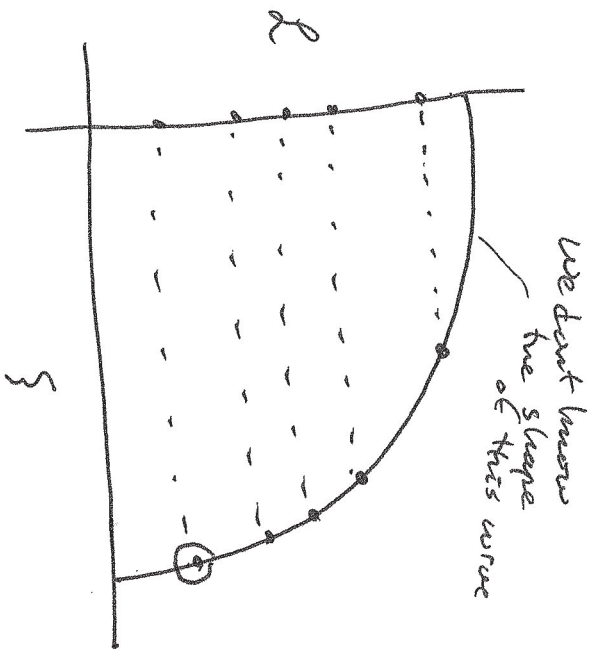
$$\text{let } \xi = \int \pi(x) dx$$
$$\Rightarrow d\xi = \pi(x) dx$$

$$\Rightarrow Z = \int_0^1 \mathcal{L}(\xi) d\xi$$

$$\mathcal{L}(x) = \mathcal{L}(\xi)$$

How do we perform this integral?

Nested Sampling: Stochastic Integration



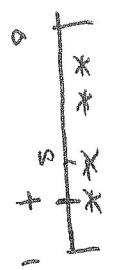
model space

Choose N samples in X
 Compute likelihoods
 Each sample represents $1/N$ of the total volume.
 The worst sample defines $l_k = l^*$ which is the likelihood value that we now use as an implicit boundary.

Order Statistics

When the worst sample is rejected, the interval implicitly defined by l^* shrinks.
 By how much?

Consider the unit interval



with n samples.

One sample is the rightmost sample.

What is the probability that S

will be to the left of t ?

$$P(S < t | I) = t$$

For $n-1$ samples

$$P(\text{all } n-1 \text{ to the left of } t) = t^{n-1} \Rightarrow P(t) = n t^{n-1}$$

$$\Rightarrow \text{Prob}(S/S^* = t) = t^{n-1} \Rightarrow C = n$$

Order Statistics

Consider the unit interval covered with n samples...



Let t be the greatest sample

What is the probability that another sample $s < t$?

$$P(S < t | I) = t$$

For $n-1$ samples, what is the probability that $n-1$ samples are less than t ?

$$P(s_1 < t, s_2 < t, \dots, s_{n-1} < t | I) \propto \prod_{i=1}^{n-1} P(s_i < t | I)$$

Same as the Probability that t is the largest of n samples $\propto \prod_{i=1}^{n-1} t$

Normalize it

$$P(t | I) = C t^{n-1}$$

$$\int P(t | I) dt = \int_0^1 C t^{n-1} dt = 1$$

$$\Rightarrow C \frac{1}{n} = 1$$

$$\Rightarrow C = n \quad P(t) = n t^{n-1}$$

$$\Rightarrow C = n$$

Order Statistics

Find mean value of $\log t$...

$t = \frac{3}{3}^x$ Shrinkage ratio

$$P(t) = n t^{n-1}$$

$$\text{let } s = \log t \Rightarrow \frac{ds}{dt} = \frac{1}{t}$$

$$P(s) = P(t) \left| \frac{dt}{ds} \right|$$

$$= t P(t)$$

$$= n t^n$$

$$= n e^{ns} \quad (\text{since } s = \log t \Rightarrow t = e^s)$$

$$\langle \log t \rangle = \langle s \rangle$$

$$= \int_{-\infty}^0 s n e^{ns} ds$$

$$= \frac{e^{ns} (ns - 1)}{n} \Big|_{-\infty}^0$$

$$= -\frac{1}{n} - 0$$

Similarly for variance s.t. $\langle (\log t)^2 \rangle = -\frac{1}{n} + \frac{1}{n^2}$

Nested Sampling

At the k^{th} iteration

$$f_k = f^*$$

$$\xi_k = \xi^* = \frac{k}{n} + \epsilon_j$$

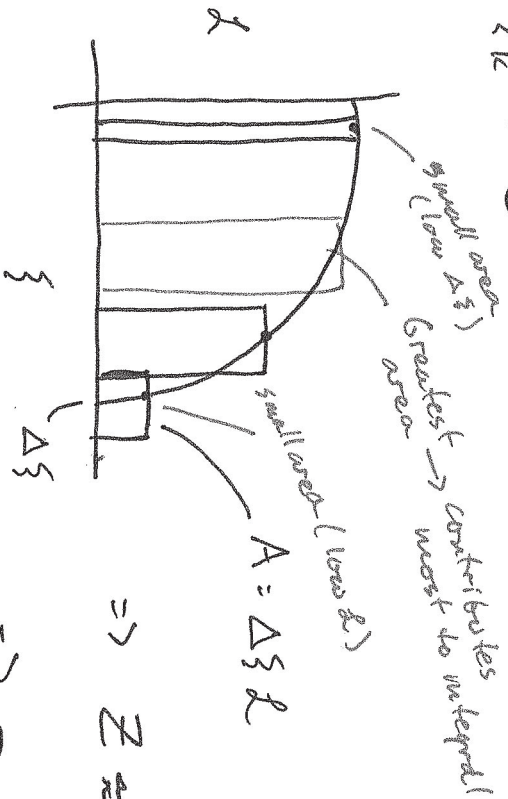
$$\Rightarrow \log \xi_k = (-k \pm \sqrt{k})/n = -\frac{k}{n} \pm \frac{\sqrt{k}}{n}$$

Ignoring our uncertainty about $\log \xi_k$

$$\Rightarrow \log \xi_k = -\frac{k}{n}$$

$$\Rightarrow \xi_k = e^{-k/n}$$

Better to sample $\log \xi_k$!



$$\Rightarrow Z \approx \sum_k f_k \Delta \xi_k$$

$$\Rightarrow Z \approx \sum_k A_k$$

The entropy can also be found:

where $\Delta \xi_k \approx \xi_{k-1} - \xi_k$

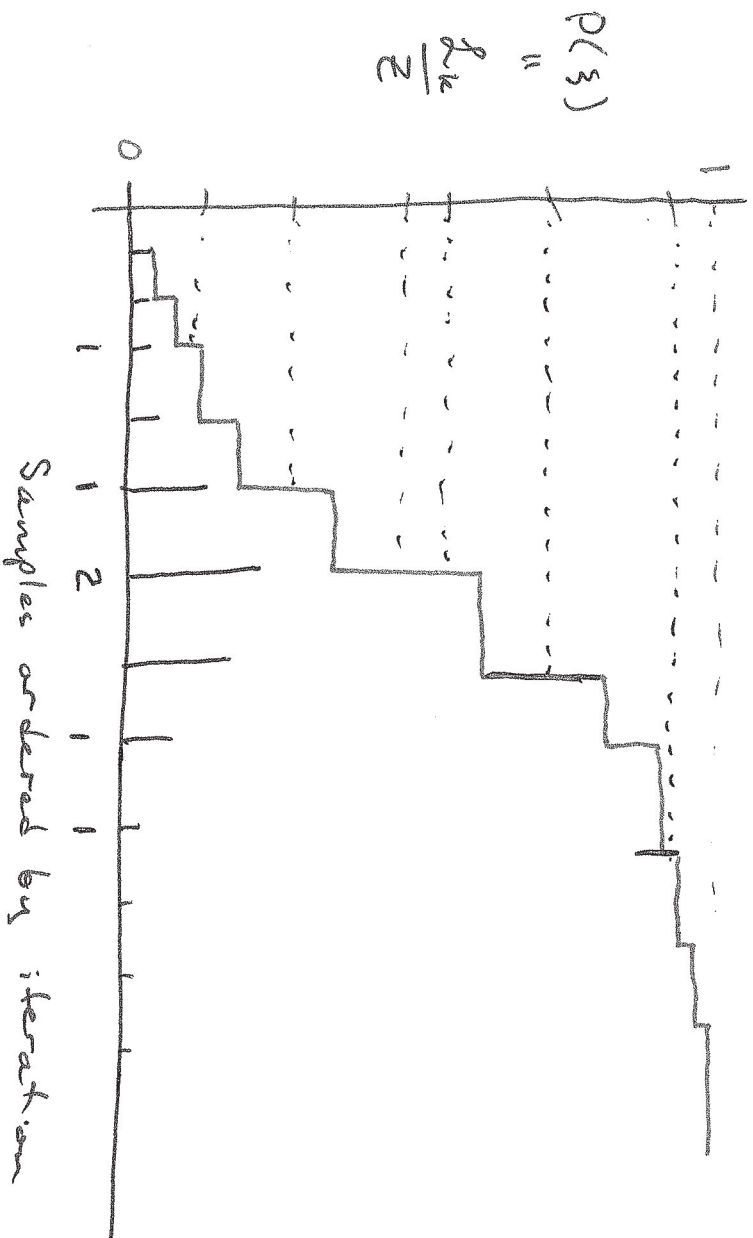
$$\mathcal{H} = \int P(\xi) \log [P(\xi)] d\xi \approx \sum_k \frac{A_k}{Z} \log \left[\frac{A_k}{Z} \right]$$

$$\Rightarrow \log Z \approx \log \left(\sum_k A_k \right) \pm \sqrt{\mathcal{H}}$$

Nested Sampling: Key Points

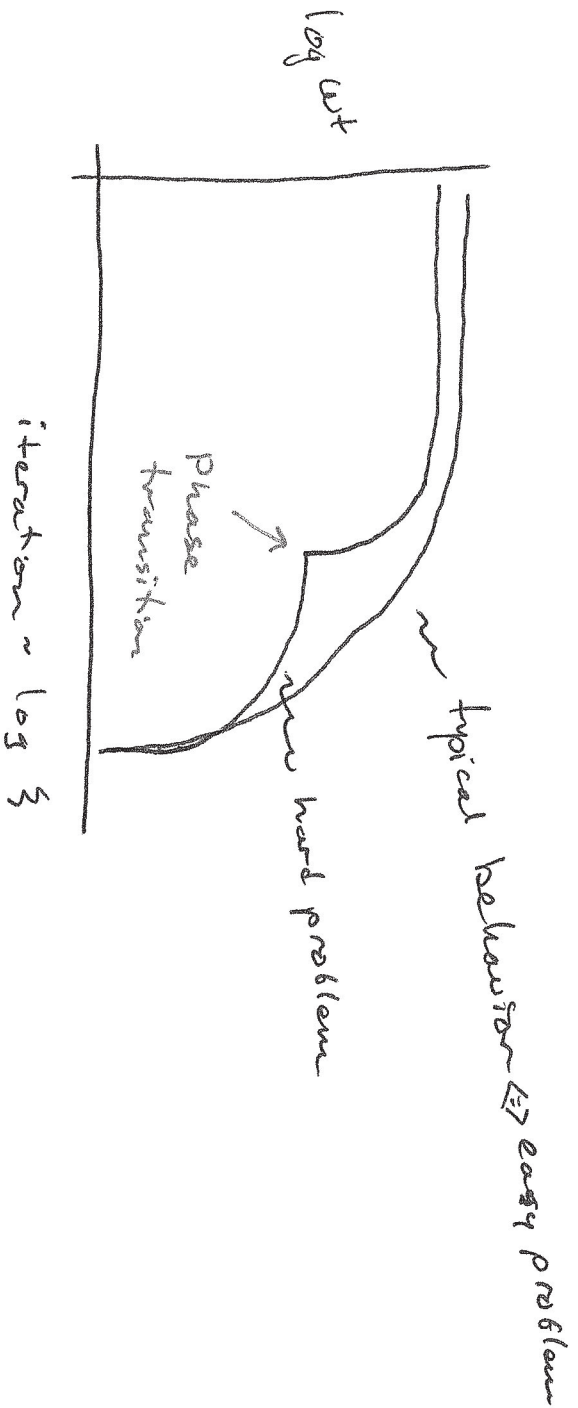
- The samples must be selected from the Prior probability. In this implementation, it is assumed that the prior was uniform.
- Often One can perform a change of variables so that your priors are uniform.
- The algorithm is focused on integrating to get the evidence
- The set of thrown-away samples do not constitute a set of samples distributed according to the posterior probability.

Sampling from the Posterior using Rejected Sampling

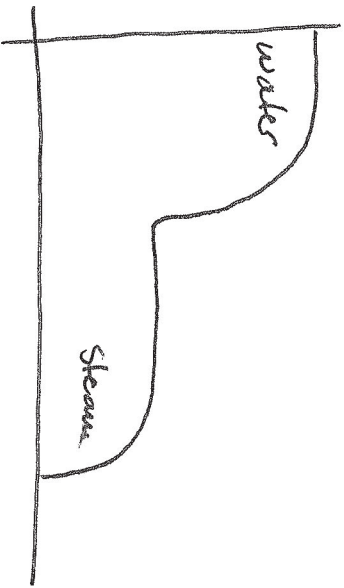


1. Imagine the cumulative distribution function, which is the integral of the pdf from $-\infty$ to x . Note that the pdf is represented by discrete samples.
2. If you want N samples, select N samples from $[0, 1]$. Mark them on the y-axis.
3. Project N samples onto the x-axis to select samples along the x-axis. The resulting set of samples (including N replications) is distributed according to the posterior.

Density of States and Phase Transitions.



}}}



If you do not iterate long enough, it looks like you have converged! But no!

The water phase is a highly probable region with small volume sitting on top of the steam phase.

Phase Transitions \Leftrightarrow AHAs!!

Inspiration

Interestingly, one can often look at the behavior of the algorithm and identify which specific concepts are being learned at each phase transition.

Simulated Annealing

Turn data on slowly.

Priors are smoother and easy to navigate.

Likelihoods can be rough.

Z^λ

$$P(M, D, I; \lambda) \propto \frac{P(M|I) P(D|M, I)^\lambda}{P(D|I)}$$

let λ go from 0 to 1

This "turns on" the data term.

λ is the "temperature"

$\lambda = 0$ cold

$\lambda = 1$ hot

$$\langle \log R \rangle_\lambda = \int \log R dP_\lambda = \frac{\int R^\lambda \log R dZ}{\int R^\lambda dZ} = \frac{d}{d\lambda} \left[\log \left(\int R^\lambda dZ \right) \right]$$

If we integrate from $\lambda = 0$ to 1

$$\int_0^1 \langle \log R \rangle_\lambda d\lambda = \log \left(\int R^\lambda dZ \right) \Big|_0^1 = \log \left(\int R dZ \right) - \log \left(\int dZ \right)$$

$$= \log \frac{\left(\int R dZ \right)}{\left(\int dZ \right)}$$

This methodology is another way to find Z .

$$= \log Z \quad \dots \quad \text{the slope}$$

This methodology is Integration

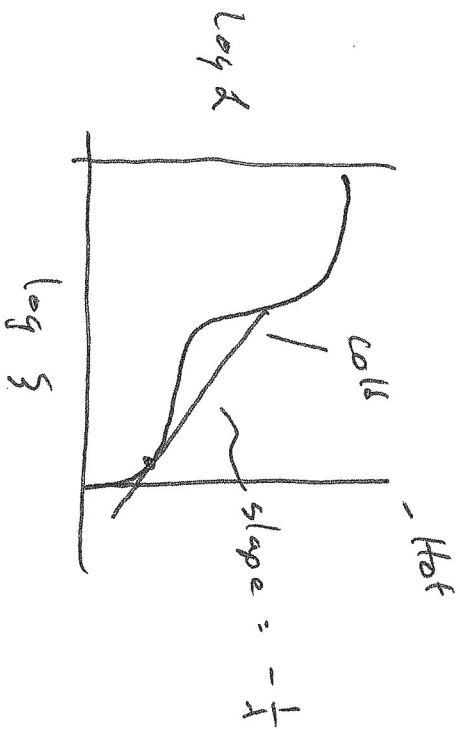
has problems at phase transitions since if $\frac{d \log R}{d \log Z} = \dots$

Simulated Annealing Tracker Slope

$$dP \propto L^{\lambda} dS$$

The greatest probability mass is most likely to be found under the maximum at ~~$L^{\lambda} dS$~~ $L^{\lambda} dS$.

This maximum occurs at $\frac{d \log L}{d \log S} = -\frac{1}{\lambda}$



As one cools from $\lambda = 0$ (infinite slope) to $\lambda = 1$ (45°), simulated annealing will eventually require a jump to get over the phase transition.

The algorithm gets stuck here.