# Nested sampling

- Original method due to Skilling (2004)

- Actually a method for computing the <u>Bayesian evidence</u>, $Z$

- Useful by-product: We get $\bar{\theta}$ samples distributed according to $p(\bar{\theta}|D)$

- Want to compute $Z = \int L(\bar{\theta}) \Pi(\bar{\theta}) d\bar{\theta}$     (∗)

- High-dim integrals are hard! One-dim. are easy!

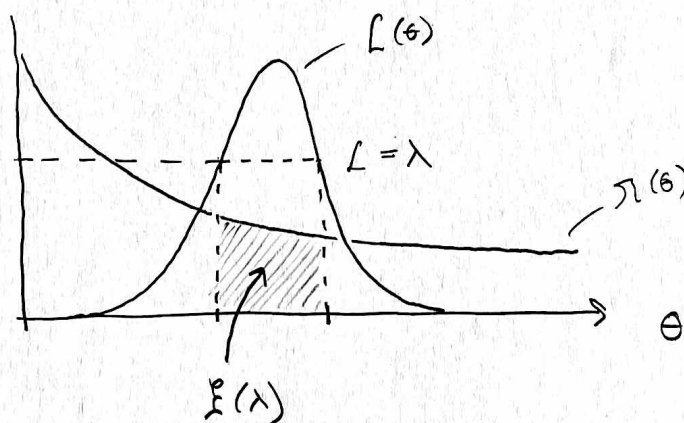- Can we turn (∗) into a one-dim. integral?

- Introduce variable: "prior mass"

$$\begin{array}{|l|}
\hline
\text{Terminology:} \\
\text{prob. density} = \dfrac{d(\text{prob. mass})}{d(\text{volume})} \\
\hline
\end{array}$$

$$\xi(\lambda) = \int_{L(\theta) > \lambda} \Pi(\theta) d\theta$$

$$d\xi = \Pi(\theta) d\theta$$

$$\begin{array}{|l|}
\hline
\text{Transf. of random variable:} \\
p_x(x) dx = p_y(y) dy \\
\text{Here: } \Pi_\xi(\xi) d\xi = \Pi_\theta(\theta) d\theta \\
\text{with } \Pi_\xi(\xi) = 1 \\
\text{and } \xi \in [0,1] \\
\hline
\end{array}$$



$$\begin{array}{|l|}
\hline
d\xi: \text{ The small additional} \\
\text{prior mass included} \\
\text{by lowering the} \\
\text{likelihood threshold} \\
\text{by } dL \\
\hline
\end{array}$$

$\xi(\lambda)$: The amount of <u>prior probability</u> contained within the regions of parameter space where the likelihood $L(\theta)$ is greater than some value $\lambda$
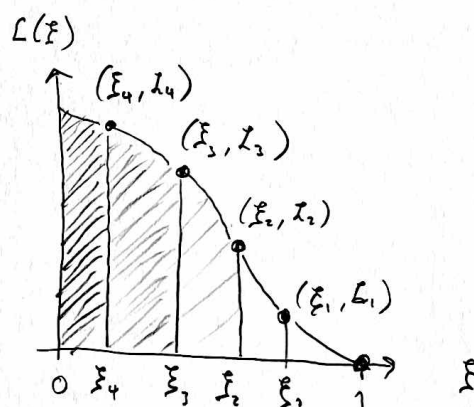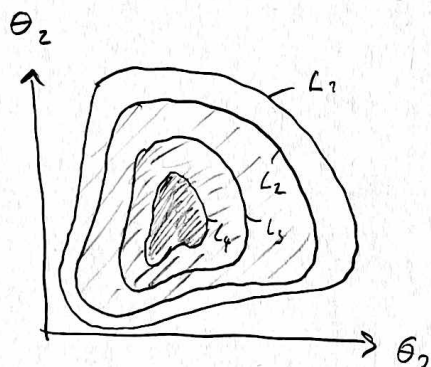
Examples: $\xi(0) = 1$ , $\xi(\lambda = L_{max}) = 0$

○ Note that $\xi(\lambda)$ is a one-dim, decreasing function of $\lambda$

○ <u>Inverse function</u>, denoted as $L(\xi)$ is simply

$$L(\xi(\lambda)) \equiv \lambda = \text{the value for the likelihood contour that contains a given prior mass } \xi$$

○ Can now express $Z$ as <u>one-dim, integral</u> over $\xi$ :

$$Z = \int_0^1 L(\xi)\, d\xi$$



○ If we can get a set of <u>ordered pairs of values</u> $(\xi_i, L_i)$ we can evaluate $Z$ integral using standard methods, (e.g trapezoidal rule )

○ <u>Nested sampling</u> is algorithm to get these samples

# Algorithm   [Show slides]

1) Draw $N$ "live points" $\bar{\theta}$ according to prior $\pi(\bar{\theta})$

2) Evaluate $L(\bar{\theta})$ at each live point

3) Discard (but record) point with lowest likelihood

4) Draw a new point from _prior_, but with additional req. that $L(\bar{\theta}_{new}) > L(\bar{\theta}_{disc.})$   ← Main challenge for algo. efficiency!

5) Repeat from step 3

○ The _discarded_ points form ordered set of likelihood samples

$$0 < L_1 < L_2 < \dots$$

○ For each likelihood sample, _can estimate_ corresponding prior mass $\xi_i$ to obtain   (will show this later)

$$1 > \xi_1 > \xi_2 > \dots$$

○ Result:

Evidence estimate:

$$Z \approx \sum_{i=1}^{M} L_i w_i = \sum_{i=1}^{M} L_i \frac{1}{2}\left[\xi_{i-1} - \xi_{i+1}\right]$$

$w_i$ : The slice of prior mass associated with the likelihood value $L_i$.
(Here chosen according to trapezoidal rule)

Posterior samples:

Assign each discarded parameter sample $\bar{\theta}_i$ its share of the posterior prob.

$$p_i = \frac{L_i w_i}{Z}$$

Main challenge: How to efficiently draw replacement samples from the "likelihood-constrained prior"?
MultiNest + friends solve this!

• How can we estimate the prior mass $\xi_i$ corresponding to a likelihood value $L_i$ ?

$$0 < L_1 < L_2 < \ldots$$

$$1 > \xi_1 > \xi_2 > \ldots$$

- From $d\xi = \pi(\bar{\theta})d\bar{\theta}$, we know that sampling $\bar{\theta}$ according to $\pi(\bar{\theta})$ corresponds to sampling $\xi$ from uniform distribution $U(0,1)$

$$\begin{bmatrix} \text{Recall relation :} \\ \\ \bar{\theta} \longrightarrow L(\bar{\theta}) \rightarrow \lambda \longrightarrow \begin{array}{c} \text{Integration} \\ \text{contour for} \\ \int\limits_{L(\bar{\theta}) > \lambda} \pi(\bar{\theta})d\bar{\theta} \end{array} \longrightarrow \xi \end{bmatrix}$$

- Sampling constraint $L(\bar{\theta}_{new}) > L(\bar{\theta}_{disc.})$ ensures that the prior mass associated with $\bar{\theta}_{new}$ is __smaller__ than that for $\bar{\theta}_{disc.}$

$$\xi_{new} < \xi_{disc.}$$

- We use $N$ live points. At start of iteration $i$ we have $N$ $\bar{\theta}$-samples that should correspond to $N$ $\xi$-samples from uniform distr. on $(0, \xi_{i-1})$

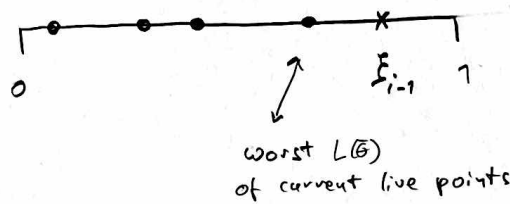- The prior mass $\xi_i$ of next point to be discorded is an unknow/random variable

$$\xi_i = t_i \, \xi_{i-1}$$

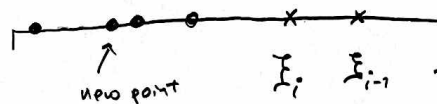where the shrinkage factor $t_i = \dfrac{\xi_i}{\xi_{i-1}}$ has a pdf

$$p(t) = N t^{N-1} \qquad \begin{bmatrix} \text{Pdf for the largest value } t \text{ of} \\ N \text{ samples drawn from } U(0,1) \end{bmatrix}$$
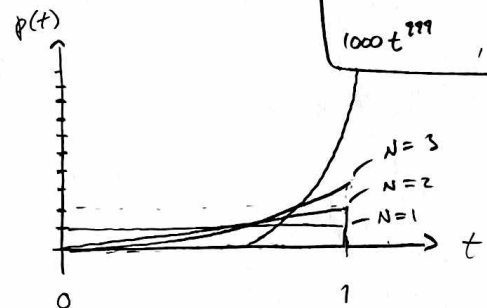
<u>Iteration $i$</u> $(N=4)$



worst $L(\vec{\theta})$
of current live points

$\Downarrow$     Discard and sample
a new point under constraint that $L(\vec{\theta})$



new point    $\xi_i$   $\xi_{i-1}$   1

$$\left[ \begin{array}{l} \text{We don't try to compute the exact } \xi_i \text{ corr} \\ \text{for the iteration } i \text{ , we just estimate it} \end{array} \right]$$

$$p(t) = N t^{N-1} \quad = \quad \begin{cases} 1 & , \ N=1 \\ 2t & , \ N=2 \\ 3t^2 & , \ N=3 \\ \\ 1000 \, t^{999} & , \ N=1000 \end{cases}$$



- Since we start from $\xi_0 = 1$, we can express $\xi_i$ as the random variable

$$\xi_i = t_i \, t_{i-1} \dots t_1 \qquad \left( \text{since } \xi_i = t_i \, \xi_{i-1} = t_i \, t_{i-1} \, \xi_{i-2} = \dots \right)$$

or $\quad \ln \xi_i = \ln t_i + \ln t_{i-1} + \dots$

- All the $t_i$ have pdf $p(t) = N t^{N-1}$ which give

$$E[\ln t] = -\frac{1}{N} \quad , \quad \text{Var}[\ln t] = \frac{1}{N^2}$$

- This means that the sum $\ln \xi_i = \ln t_i + \ln t_{i-1} + \ldots = \Sigma$ has expectation am variance

$$E\left[\ln \xi_i\right] = E\left[\ln t_i\right] + E\left[\ln t_{i-1}\right] + \ldots = \frac{i}{N^2}$$

$$= -\frac{1}{N} - \frac{1}{N} - \ldots$$

$$= -\frac{i}{N}$$

$$\mathrm{Var}\left[\ln \xi_i\right] = \mathrm{Var}\left[\ln t_i\right] + \mathrm{Var}\left[\ln t_{i-1}\right] + \ldots \quad \left(\begin{array}{l} \text{Since the} \\ t_i \text{ are uncorrelated} \end{array}\right)$$

$$= \frac{i}{N^2}$$

- In short: $\quad \ln \xi_i \approx -\frac{i}{N} \pm \frac{\sqrt{i}}{N}$

- So we approximate the prior mass $\xi_i$ associated with the likelihood value $L_i$ of the discarded point $\bar{\theta}$ at iteration $i$ as

$$\boxed{\xi_i \approx e^{-\frac{i}{N}}}$$

- Then we have what we need to compute a new approx. for $Z$ after each new iteration $i$

- The sampling stops when the largest possible contribution $\Delta Z$ from the current live points is much smaller than the current estimate for $Z$

  (But this will fail if the sampling has missed some region of high likelihood)


- Uncertainty on final evidence estimate is dominated by uncertainty in $\xi_i$ estimates

  (assuming the sampling has found all relevant parameter regions.)

- Efficiency challenge:

  — Naively sampling $\theta$ points from entire $\Pi(\theta)$ at every iteration will lead to ever decreasing efficiency, due to constraint $L(\theta_{new}) > L(\theta_{disc.})$

  — One appr. used to alleviate problem:

    — Draw samples from ellipsoids containing current live points

    — Use clustering algos. to assign sep. ellipsoids to sep. clusters of live points

- Much used packages: MultiNest, PolyChord

  (pymultinest)