

Bayesian model selection and parameter estimation

Bayesian inference

Aim: use available data to

- Construct probability density distributions for parameters associated with these hypotheses
 - **Parameter estimation**
- Evaluate which out of several hypotheses is the most likely
 - **Model selection**

Do this while making explicit all extraneous assumptions

Inductive logic

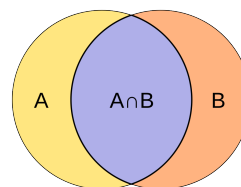
- Propositions (i.e. statements, events) denoted by uppercase letters, e.g. A , B , C , ..., X
- Boolean algebra:
 - Conjunction: *A and B are both true*
 AB or $A \wedge B$
 - Disjunction: *At least one of A or B is true*
 $A + B$ or $A \vee B$
 - Negation: *A is false*
 \bar{A} or $\neg A$
 - Implication: *From A follows B*
 $A \rightarrow B$ or $A \Rightarrow B$

Probabilities for propositions

- Useful to view statements as sets which are subsets of a “Universe”

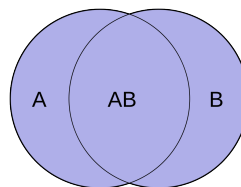
- Conjunction: intersection of sets

$$AB \text{ or } A \wedge B$$



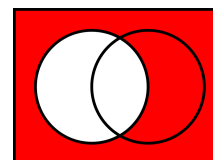
- Disjunction: union of sets

$$A + B \text{ or } A \vee B$$



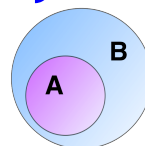
- Negation: complement within Universe

$$\bar{A} \text{ or } \neg A$$

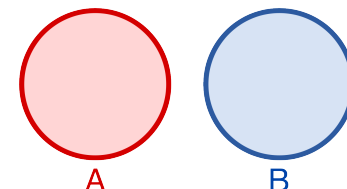


- Each of these sets have a probability associated with them

- If $A \subset B$ then $p(A) \leq p(B)$



- If A and B disjoint then $p(A \cup B) = p(A) + p(B)$



- The Universe has probability 1, so that e.g.

$$p(A) + p(\bar{A}) = 1$$

Conditional probability

- *Conditional probability:* $p(A|B) \equiv \frac{p(A \cap B)}{p(B)}$

- *Product rule:*

$$p(A, B) = p(A \cap B) = p(A|B) p(B)$$

- It is customary to explicitly denote probabilities being conditional on “all background information we have”: $p(A|I)$, $p(B|I)$, ...

- All essential formulae unaffected, e.g. product rule:

$$p(A, B|I) = p(A|B, I) p(B|I)$$

- From the product rule follows Bayes' theorem:

$$p(A|B, I) = \frac{p(B|A, I) p(A|I)}{p(B|I)}$$

Marginalization

- Note that for any A and B,

$$A \cap B \quad \text{and} \quad A \cap \bar{B}$$

are disjoint sets whose union is A, so

$$p(A|I) = p(A, B|I) + p(A, \bar{B}|I)$$

- Consider sets $\{B_k\}$ such that

- They are disjoint: $B_k \cap B_l = \emptyset, \quad k \neq l$

- They are exhaustive: $\cup_k B_k$ is the Universe, or $\sum_k p(B_k|I) = 1$

Then

$$p(A|I) = \sum_k p(A, B_k|I)$$

Marginalization rule

Marginalization over a continuous variable

- Consider the proposition

“The continuous variable x has the value α ”

Then not necessarily a well-defined meaning of probability $p(x = \alpha|I)$

- Instead assign probabilities to finite intervals:

$$p(x_1 \leq x \leq x_2|I) = \int_{x_1}^{x_2} \text{pdf}(x)dx$$

where “pdf” is the *probability density function*

- Exhaustiveness written as

$$\int_{x_{\min}}^{x_{\max}} \text{pdf}(x)dx = 1$$

- Marginalization for continuous variables:

$$p(A) = \int_{x_{\min}}^{x_{\max}} \text{pdf}(A, x)dx$$

Parameter estimation

- Experiment performed, data d collected
- Parameter θ being measured
- Consider a model H that allows to calculate probability of getting data d if parameter θ is known (“*generative model*”)

– Can calculate the *likelihood* $p(d|\theta, H, I)$

- What is wanted is instead posterior probability of θ , $p(\theta|d, H, I)$

- Use Bayes' theorem:

$$p(\theta|d, H, I) = \frac{p(d|\theta, H, I)p(\theta|H, I)}{p(d|H, I)}$$

- “*Prior*” $p(\theta|H, I)$ is our knowledge of θ before experiment
- “*Evidence*” $p(d|H, I)$ doesn't depend on θ , ignore for now

$$p(\theta|d, H, I) \propto p(d|\theta, H, I)p(\theta|H, I)$$

Parameter estimation

$$p(\theta|d, H, I) \propto p(d|\theta, H, I)p(\theta|H, I)$$

- Posterior is likelihood weighted by prior

Conclusions drawn based on:

- Information available before experiment
- Experimental data obtained

- Can extend to more parameters: joint posterior $p(\theta_1, \dots, \theta_N|d, H, I)$

- If we want posterior distribution just for variable θ_1 , $p(\theta_1|d, H, I)$,

then we marginalize:

$$p(\theta_1|d, H, I) = \int_{\theta_2^{\min}}^{\theta_2^{\max}} \dots \int_{\theta_N^{\min}}^{\theta_N^{\max}} p(\theta_1, \dots, \theta_N|d, H, I) d\theta_2 \dots d\theta_N$$

- Mean of a 1D posterior:

$$\begin{aligned}\mu &= E[\theta] \\ &= \int_{\theta^{\min}}^{\theta^{\max}} \theta p(\theta|d, H, I) d\theta\end{aligned}$$

- Variance of a 1D posterior:

$$\begin{aligned}\sigma^2 &= E[(\theta - \mu)^2] \\ &= \int_{\theta^{\min}}^{\theta^{\max}} (\theta - \mu)^2 p(\theta|d, H, I) d\theta\end{aligned}$$

- Means for N variables:

$$\begin{aligned}\mu_i &= E[\theta_i] \\ &= \int_{\theta_1^{\min}}^{\theta_1^{\max}} \dots \int_{\theta_N^{\min}}^{\theta_N^{\max}} \theta_i p(\theta_1, \dots, \theta_N|d, H, I) d\theta_1 \dots d\theta_N\end{aligned}$$

- Covariance matrix:

$$\begin{aligned}\Sigma_{ij} &\equiv E[(\theta_i - \mu_i)(\theta_j - \mu_j)] \\ &= \int_{\theta_1^{\min}}^{\theta_1^{\max}} \dots \int_{\theta_N^{\min}}^{\theta_N^{\max}} (\theta_i - \mu_i)(\theta_j - \mu_j) p(\theta_1, \dots, \theta_N|d, H, I) d\theta_1 \dots d\theta_N\end{aligned}$$

- *Confidence interval* is the smallest interval within whose limits a fraction γ of the posterior is contained:

$$\gamma = \int_{\theta^{\text{lo}}}^{\theta^{\text{hi}}} p(\theta|d, H, I) d\theta$$

where $\theta^{\text{hi}} - \theta^{\text{lo}}$ is minimal

- In most literature γ is taken to be 0.68 or 0.95, roughly corresponding to 1-sigma and 2-sigma intervals of Gaussian distribution

- Multi-dimensional confidence intervals:

$$\begin{aligned} \gamma_{\theta_1} &= \int_{\theta_1^{\text{lo}}}^{\theta_1^{\text{hi}}} p(\theta_1|d, H, I) d\theta_1 \\ &= \int_{\theta_1^{\text{lo}}}^{\theta_1^{\text{hi}}} \int_{\theta_2^{\text{min}}}^{\theta_2^{\text{max}}} \dots \int_{\theta_N^{\text{min}}}^{\theta_N^{\text{max}}} p(\theta_1, \dots, \theta_N|d, H, I) d\theta_1 \dots d\theta_N \end{aligned}$$

Hypothesis testing

- Estimating parameters is possible if generative model known
- If we want to compare possible generative models, e.g. X , Y : calculate posterior probabilities $p(X|d, I)$ and $p(Y|d, I)$

- Bayes' theorem:

$$p(X|d, I) = \frac{p(d|X, I)p(X|I)}{p(d|I)}$$

- Compute *odds ratio*

$$\begin{aligned} O_Y^X &\equiv \frac{p(X|d, I)}{p(Y|d, I)} \\ &= \frac{p(d|X, I) p(X|I)}{p(d|Y, I) p(Y|I)} \end{aligned}$$

where factors of $p(d|I)$ have canceled out

- $p(X|I)/p(Y|I)$ ratio of prior odds
- $p(d|X, I)/p(d|Y, I)$ ratio of evidences, or *Bayes factor* $B_Y^X = \frac{p(d|X, I)}{p(d|Y, I)}$

Hypothesis testing

- Hypotheses usually have parameters associated with them
- Bayes theorem relating posterior to likelihood:

$$p(\theta|d, H, I) = \frac{p(d|\theta, H, I)p(\theta|H, I)}{p(d|H, I)}$$

or

$$p(\theta|d, H, I)p(d|H, I) = p(d|\theta, H, I)p(\theta|H, I)$$

- Marginalize both sides over parameter(s):

$$\int p(\theta|d, H, I)p(d|H, I)d\theta = \int p(d|\theta, H, I)p(\theta|H, I)d\theta$$

Note that $p(d|H, I)$ independent of parameter(s), and posterior $p(\theta|d, H, I)$ normalized by definition, hence left hand side:

$$\int p(\theta|d, H, I)p(d|H, I)d\theta = p(d|H, I) \int p(\theta|d, H, I)d\theta = p(d|H, I)$$

Therefore evidence is given by

$$p(d|H, I) = \int p(d|\theta, H, I)p(\theta|H, I)d\theta$$

Hypothesis testing

- Odds ratio

$$\begin{aligned}O_Y^X &\equiv \frac{p(X|d, I)}{p(Y|d, I)} \\ &= \frac{p(d|X, I) p(X|I)}{p(d|Y, I) p(Y|I)}\end{aligned}$$

Bayes factor

$$B_Y^X = \frac{p(d|X, I)}{p(d|Y, I)}$$

Marginalized evidences e.g. $p(d|X, I) = \int p(d|\theta, X, I)p(\theta|X, I)d\theta$

- Hypotheses can have arbitrary number of free parameters
 - Does model that fits data the best give the highest evidence?
 - If so, model with more parameters would give highest evidence even if incorrect!

Occam's razor

- For simplicity, compare two generative hypotheses:
 - X has no free parameters
 - Y has one free parameter, λ

Will Y automatically be favored over X ?

- Odds ratio $O_Y^X = \frac{p(d|X, I) p(X|I)}{p(d|Y, I) p(Y|I)}$

- Evidence for X is straightforward, but for Y :

$$p(d|Y, I) = \int p(d|\lambda, Y, I) p(\lambda|Y, I) d\lambda$$

Assume flat prior for $\lambda \in [\lambda_{\min}, \lambda_{\max}]$:

$$p(\lambda|Y, I) = \frac{1}{\lambda_{\max} - \lambda_{\min}}, \quad \text{for } \lambda_{\min} \leq \lambda \leq \lambda_{\max}$$

Occam's razor

- Evidence for Y :

$$p(d|Y, I) = \int p(d|\lambda, Y, I)p(\lambda|Y, I)d\lambda$$

- Flat prior:

$$p(\lambda|Y, I) = \frac{1}{\lambda_{\max} - \lambda_{\min}}, \quad \text{for } \lambda_{\min} \leq \lambda \leq \lambda_{\max}$$

- For definiteness, assume likelihood of the form

$$p(d|\lambda, Y, I) = p(d|\lambda_0, Y, I) \exp \left[-\frac{(\lambda - \lambda_0)^2}{2\sigma_\lambda^2} \right]$$

- Evidence for Y :

$$\begin{aligned} p(d|Y, I) &= \int p(d|\lambda, Y, I)p(\lambda|Y, I)d\lambda \\ &= \int \frac{1}{\lambda_{\max} - \lambda_{\min}} p(d|\lambda_0, Y, I) \exp \left[-\frac{(\lambda - \lambda_0)^2}{2\sigma_\lambda^2} \right] d\lambda \\ &= \frac{p(d|\lambda_0, Y, I)}{\lambda_{\max} - \lambda_{\min}} \int \exp \left[-\frac{(\lambda - \lambda_0)^2}{2\sigma_\lambda^2} \right] d\lambda \\ &= p(d|\lambda_0, Y, I) \frac{\sigma_\lambda \sqrt{2\pi}}{\lambda_{\max} - \lambda_{\min}}. \end{aligned}$$

Occam's razor

- Evidence for Y :

$$p(d|Y, I) = p(d|\lambda_0, Y, I) \frac{\sigma_\lambda \sqrt{2\pi}}{\lambda_{\max} - \lambda_{\min}}$$

- Hence odds ratio becomes:

$$O_Y^X = \frac{p(X|I)}{p(Y|I)} \frac{p(d|X, I)}{p(d|\lambda_0, Y, I)} \frac{\lambda_{\max} - \lambda_{\min}}{\sigma_\lambda \sqrt{2\pi}}$$

where

- $p(X|I)/p(Y|I)$ ratio of prior odds; can be set to 1 in this example
- $p(d|X, I)/p(d|\lambda_0, Y, I)$ just compares best fits; will usually be < 1
- $(\lambda_{\max} - \lambda_{\min})/(\sigma_\lambda \sqrt{2\pi})$ penalizes Y if experimental uncertainty on λ much smaller than prior range
 - Will tend to be the case if λ not needed!

Occam's Razor:

“It is vain to do with more what can be done with fewer”

Likelihood principle

- Suppose experiment with generative hypothesis H , corresponding set of N parameters θ , observed data d , and background information I
- Then posterior of θ can be expressed using Bayes' theorem:

$$p(\theta|d, H, I) = \frac{p(d|\theta, H, I)p(\theta|H, I)}{p(d|H, I)}$$

- Only factor in RHS that depends on d and involves the θ is the likelihood $p(d|\theta, H, I)$
 - The likelihood function $p(d|\theta, H, I)$ contains all the information about the parameters θ that is present in the data
 - Only need to focus on the likelihood

Nested sampling

- Parameter estimation requires computing the posterior density distribution from likelihood and prior using Bayes' theorem:

$$p(\boldsymbol{\theta}|d, H, I) = \frac{p(d|\boldsymbol{\theta}, H, I)p(\boldsymbol{\theta}|H, I)}{p(d|H, I)}$$

- Often the parameter space has high dimensionality (e.g. 15 for quasi-circular binary inspiral), making it computationally challenging to map out the likelihood
- Similarly calculation of evidence integral over high-dimensional space:

$$\begin{aligned} p(d|H, I) &= \int d^N \boldsymbol{\theta} p(d|\boldsymbol{\theta}, H, I)p(\boldsymbol{\theta}|H, I) \\ &= \int d^N \boldsymbol{\theta} L(\boldsymbol{\theta})\pi(\boldsymbol{\theta}), \end{aligned}$$

- Efficient way of obtaining both: *nested sampling*

Nested sampling: basic idea

$$\begin{aligned} p(d|H, I) &= \int d^N \boldsymbol{\theta} p(d|\boldsymbol{\theta}, H, I) p(\boldsymbol{\theta}|H, I) \\ &= \int d^N \boldsymbol{\theta} L(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}), \end{aligned}$$

- Nested sampling computes the evidence by rewriting the above integration in terms of a single scalar called *prior mass* X
- “Fraction of volume with likelihood greater than λ ”

Mathematically:

$$X(\lambda) \equiv \int \int \cdots \int_{L(\boldsymbol{\theta}) > \lambda} \pi(\boldsymbol{\theta}) d^N \boldsymbol{\theta}$$

Element of prior mass: $dX = \pi(\boldsymbol{\theta}) d^N \boldsymbol{\theta}$

- Since prior is normalized, $X \in [0, 1]$

- Lower bound $X = 0$:

surface within which no higher likelihood; $\lambda = L_{\max}$

- Upper bound $X = 1$:

surface within which all points higher likelihood; $\lambda = L_{\min}$

Nested sampling: basic idea

$$\begin{aligned} p(d|H, I) &= \int d^N \boldsymbol{\theta} p(d|\boldsymbol{\theta}, H, I) p(\boldsymbol{\theta}|H, I) \\ &= \int d^N \boldsymbol{\theta} L(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}), \end{aligned}$$

- Rewrite as

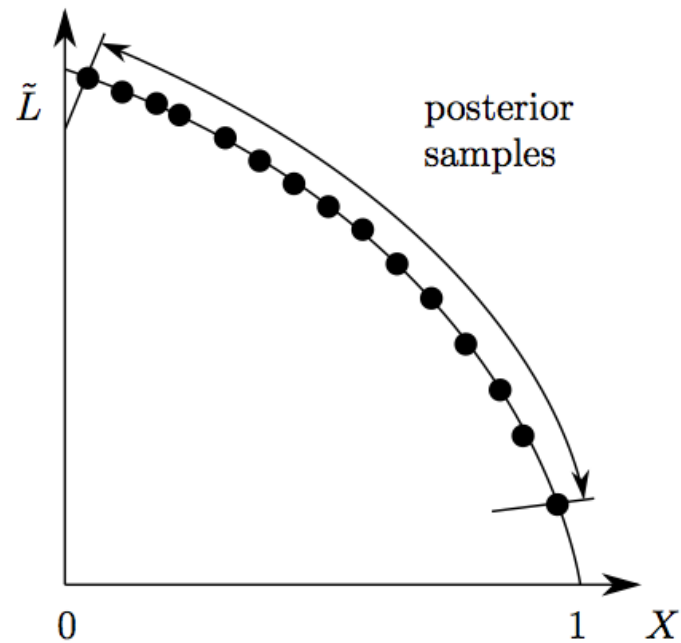
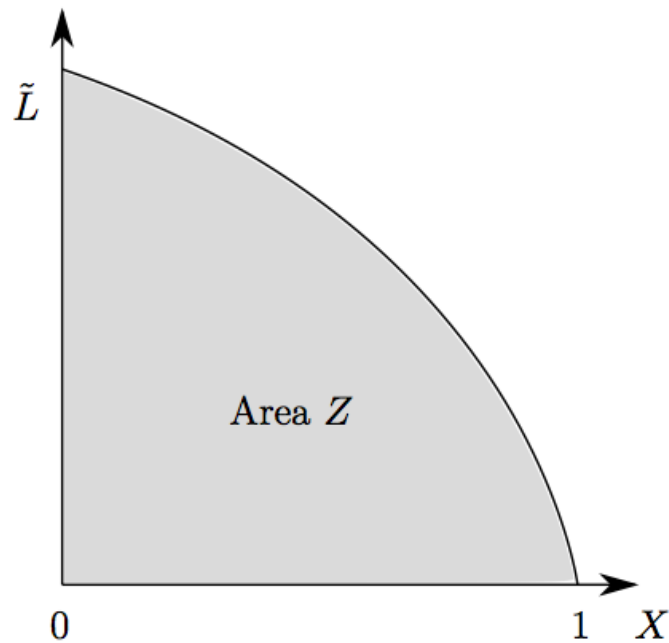
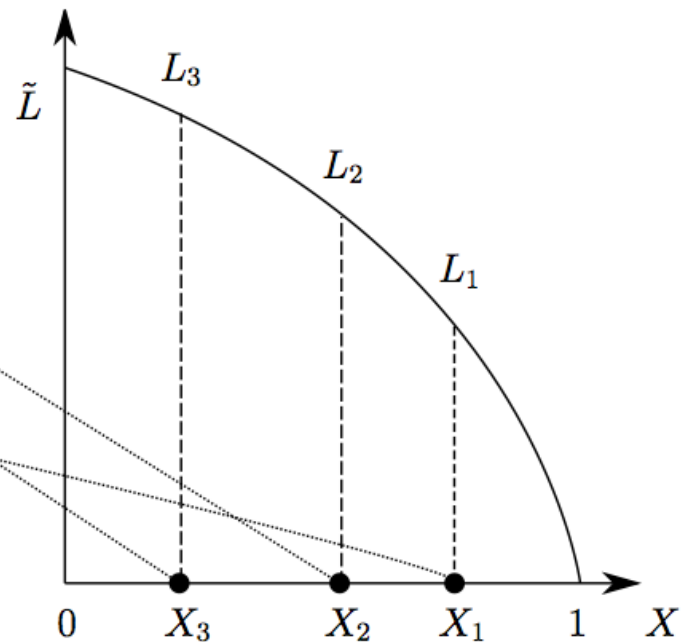
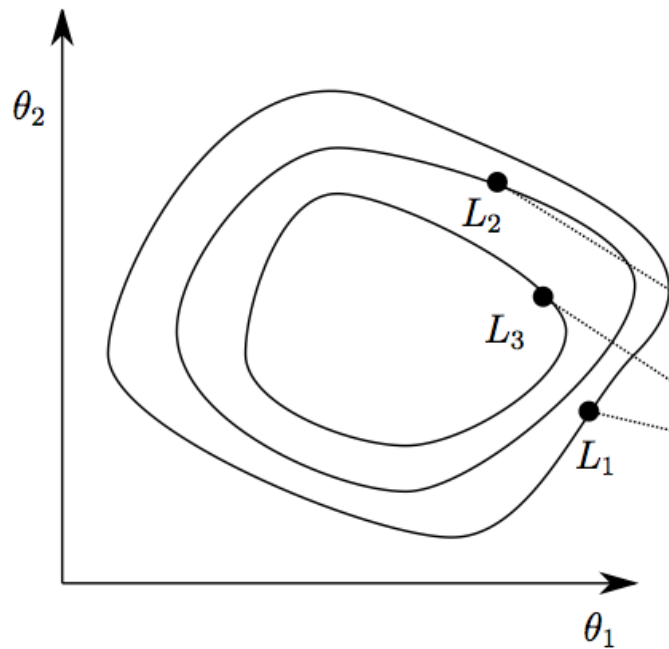
$$\begin{aligned} Z &= \int \int \cdots \int L(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d^N \boldsymbol{\theta} \\ &= \int \tilde{L}(X) dX. \end{aligned}$$

- Posterior obtained trivially from

$$\tilde{P}(X) = \frac{\tilde{L}(X)}{Z}$$

- Idea behind nested sampling: construct the function $\tilde{L}(X)$ by progressively finding locations in parameter space with higher likelihood and associated progressively smaller prior mass
 - Then use above formulae for evidence, posterior

Nested sampling: schematically



Nested sampling: the algorithm

- Drop M samples across parameter space, sampled from the prior
These are called “live points”
 - Each has likelihood associated with it
 - Associated with surface s.t. likelihood higher at boundary
 - Uniformly sampled in prior mass between 0 and 1
- Discard sample with lowest likelihood L_0 , i.e. highest prior mass X_0
 - Replace by new live point, sampled from the prior, which has with smaller likelihood
 - New point with lowest likelihood L_1 must have $X_1 < X_0$
 - *Statistically* assign value for X_1
- Repeat the step above

Nested sampling: the algorithm

- Having discarded the old highest-likelihood point with prior mass X_0 , how do we statistically assign a prior mass X_1 to the new highest-likelihood point?
- Probability that the surface with highest prior mass is at $X = \chi$ is joint probability that none of the samples have prior mass $> \chi$

$$P(X_i < \chi) = \prod_{i=1}^M \int_0^{\chi} dX_i = \prod_{i=1}^M \chi = \chi^M$$

- Probability *density* that highest of M samples has prior mass χ

$$P(\chi, M) = M\chi^{M-1}$$

- Define *shrinkage ratio* between new and old highest prior mass:

$$t = X_1/X_0$$

This has same probability density:

$$P(t, M) = Mt^{M-1}$$

- Hence we assign X_1 by drawing a shrinkage ratio from the above distribution

Nested sampling

- At first step: set $X = 1$
- At k^{th} iteration: live point with largest prior mass has

$$X_k = \prod_{j=1}^k t_j$$

- Recall distribution of shrinkage ratios:

$$P(t, M) = Mt^{M-1}$$

Mean and standard deviation of $\log(t)$:

$$\log t = (-1 \pm 1)/M$$

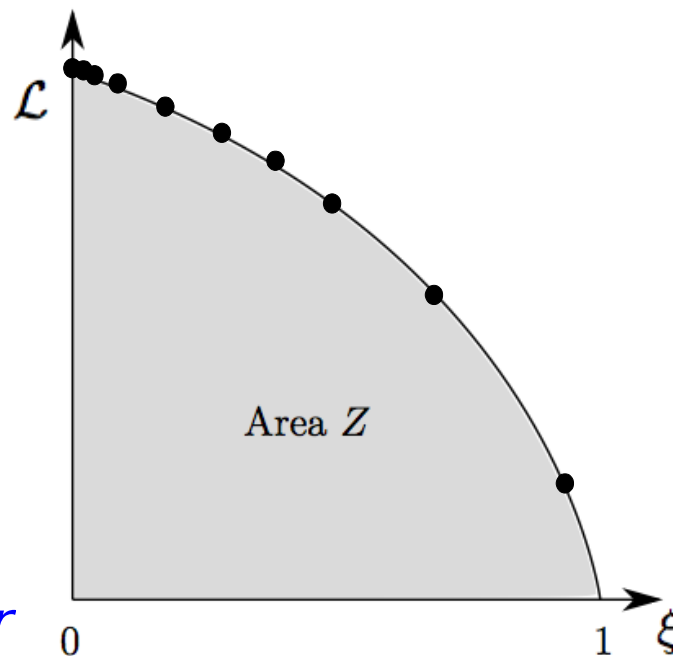
- Hence $\log(X_k)$ has mean and stdev

$$\log X_k = (k \pm \sqrt{k})/M$$

Hence mean values go like

$$X_k = \exp(-k/M)$$

- *Very quickly reaches prior mass where likelihood is largest*
- *Errors decrease exponentially*
- *Larger number of live points is better*



Nested sampling: termination condition

- No obvious choice for ending the sampling process
 - Use practical guidelines
- Estimate *information* as function of evidence and likelihood:

$$\mathcal{H} = \int P(X) \ln (P(X)) dX$$
$$\approx \sum_k \frac{L_k}{Z} \ln \frac{L_k}{Z} \Delta X_k,$$

Terminate when $X = e^{-\mathcal{H}}$

- Or, can estimate amount of evidence yet to be accumulated and compare with evidence already accumulated

Terminate when $L_{\max} X_{\text{cur}} < \alpha Z_{\text{cur}}$ where α is user-specified

Nested sampling: accuracy

- Take termination condition

$$X = e^{-\mathcal{H}}$$

- Means go like

$$X_k = \exp(-k/M)$$

“Terminate when count k exceeds $M\mathcal{H}$ ”

- Evidence:

$$Z = \int \tilde{L}(X) dX \approx \sum_k L_k \Delta X_k$$

Recall

$$\log X_k = (k \pm \sqrt{k})/M$$

Hence uncertainty on the evidence:

$$\Delta \log Z = \sqrt{\frac{\mathcal{H}}{M}}$$

- In gravitational-wave applications, with a few thousand live points this is typically $O(10^{-1})$ whereas for detectable signal $\log Z = O(10^2)$

Application to gravitational waves

- Compute evidence for hypothesis that there is a signal in the data, \mathcal{H}_S :

$$p(d|\mathcal{H}_S, I) = Z = \int \tilde{L}(X) dX \approx \sum_k L_k \Delta X_k$$

- Typical growth of $\tilde{L}(X)$: usually convenient to consider logarithm

