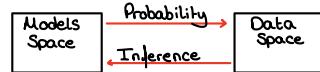


Advanced Statistical Methods for Data Science

Statistics: Collection, analysis and presentation of data aims to extract information from available data.

Probability: forward modelling, quantifying likelihood of events under uncertainty.

Inference: Drawing Conclusions from data, often probabilistic, to support/reject hypothesis



Bayesian Statistics

Probability interpreted as representing our state of knowledge / quantification of personal belief about something.

Bayes Theorem:

$$P(A, B) = P(A|B)P(B) \quad \left. \quad P(A|B) = \frac{P(B|A)P(A)}{P(B)} \right\}$$

$$P(B, A) = P(B|A)P(A)$$

↳ uses Symmetry of 'and' operator (,)

For Inference

$$P(M|D, I) = \frac{P(D|M, I)P(M|I)}{P(D|I)}$$

M - Model (values of all parameters)

D - Data

I - All of the assumptions made during analysis

rewritten:

$$P(M|D) = \frac{L(D|M)\pi(M)}{Z}$$

Posterior - Probability our model, M, is correct given observed data, D and prior information, I

$$P(M|D, I) \quad (\text{Bayesian only})$$

Likelihood - Probability of obtaining observed D data given our model M is true.

$$P(D|M, I) = L(D|M) \quad (\text{frequentist as well})$$

Prior - Probability our model, M, is correct before any data, based on prior information, I

$$P(M|I) = \pi(M)$$

Evidence - total probability of observed data across all models.

$$P(D|I) = Z(I) = \int d(\text{parameters}) \pi(\text{parameters}) L(D|\text{parameters})$$

Acts as a normalisation constant.

Likelihood: $L(D|M)$

'Probability of observed data conditioned on a particular choice of the model and any model parameters'

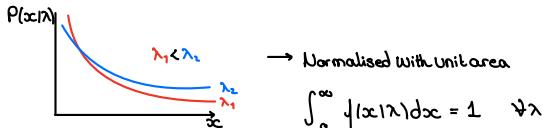
2 Versions of Likelihood

Probability Distribution $L(x|\theta) / P(x|\theta)$: $\theta \rightarrow x$

For a given $\theta \rightarrow$ function of x

Probability distribution for the data conditioned on a particular value of the model's distribution.

$$\text{eg: } P(x|\lambda) = \frac{\exp(-x/\lambda)}{\lambda} \rightarrow P(x|\lambda) = \frac{\exp(-\sum_i x_i/\lambda)}{\lambda^N}$$

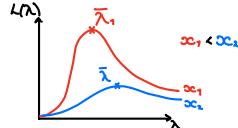


Function of model parameters $L(\theta) : x \rightarrow \theta$

For a given $x \rightarrow$ function of θ (Inference / Inverse Problem)

$$\hat{\theta} \text{ where } \frac{\partial}{\partial \theta} \log(L(\theta)) = 0$$

$$\text{eg: } L(\lambda) = \frac{\exp(-\sum_i x_i/\lambda)}{\lambda^N} \rightarrow \hat{\lambda} = \frac{\sum_i x_i}{N}$$



Both can be considered as slices through the surface $L(x|\theta)$

↳ For fixed θ : probability distribution of data conditioned on model parameters

For fixed x : likelihood of distribution parameters conditioned on measured result.

Fisher Information

quantifies the information that a observation/data point, x

gives about a unknown parameter, θ .

High Fisher Information → Data gives lots information on model parameter
Likelihood Dominated Regime

Score: Gradient of log-likelihood function wrt model parameter

$$S(\theta) = \frac{\partial \log L(x|\theta)}{\partial \theta}$$

↳ Measures sensitivity of log-likelihood to changes in parameters

→ likelihood has large discriminatory power

For True Parameters: θ_{true}

↳ Typically $S(\theta_{true}) \neq 0$ exactly due to statistical fluctuations in data

But: Averaged over all many samples x :

$$\begin{aligned} E_x[S(\theta_{true})] &= \int dx S(\theta_{true}) L(x|\theta_{true}) \quad \text{The true distribution of } x \text{ for datasets} \\ &= \int dx \frac{\partial \log L(x|\theta)}{\partial \theta} \Big|_{\theta_{true}} L(x|\theta_{true}) = \int \frac{\partial L(x|\theta)}{\partial \theta} \Big|_{\theta_{true}} \frac{1}{L(x|\theta_{true})} dx \\ &= \int \frac{\partial L(x|\theta_{true})}{\partial \theta} dx = \frac{\partial}{\partial \theta} \int L(x|\theta_{true}) dx = \frac{\partial}{\partial \theta} 1 = 0 \quad \text{Must be } \theta = \theta_{true} \end{aligned}$$

$$\therefore E_x[S(\theta_{true})] = 0$$

Fisher Information, $I(\theta)$: Variance of Score at θ_{true} wrt x , $\text{Var}_{x \sim} (S(\theta)) > 0$.

$$I(\theta) = E_x [S(\theta)^2 | \theta_{true}] - E_x [S(\theta) | \theta_{true}]^2$$

$$I(\theta) = \int dx L(x|\theta) S(\theta)^2 = \int dx L(x|\theta) \left(\frac{\partial \log L(x|\theta)}{\partial \theta} \right)^2$$

Peaked/Sensitive likelihood → Large variance of Score → Larger Fisher Information to model parameters.

Alternative Definition

$$I(\theta) = -E \left[\frac{\partial^2 \log L(x|\theta)}{\partial \theta^2} \Big|_{\theta=\theta_{true}} \right] = - \int dx \frac{\partial^2 \log L(x|\theta)}{\partial \theta^2} L(x|\theta)$$

Proof:

$$\begin{aligned} I(\theta) &= - \int dx \frac{\partial^2 \log L(x|\theta)}{\partial \theta^2} L(x|\theta) = - \int dx \left\{ \frac{\partial}{\partial \theta} \left[\frac{\partial}{\partial \theta} \log L(x|\theta) \right] \right\} L(x|\theta) \\ &= - \int dx \left\{ \frac{\partial}{\partial \theta} \left[\frac{1}{L(x|\theta)} \frac{\partial}{\partial \theta} L(x|\theta) \right] \right\} L(x|\theta) = - \int dx \left\{ \frac{1}{L(x|\theta)} \frac{\partial^2}{\partial \theta^2} L(x|\theta) - \frac{1}{L(x|\theta)^2} \left(\frac{\partial}{\partial \theta} L(x|\theta) \right)^2 \right\} L(x|\theta) \\ &= \int dx \left(\frac{1}{L(x|\theta)} \frac{\partial}{\partial \theta} L(x|\theta) \right)^2 L(x|\theta) - \underbrace{\int dx \frac{\partial^2}{\partial \theta^2} L(x|\theta)}_{\delta^2 \theta: \int dx L(x|\theta) = 0} = \int dx \left(\frac{\partial}{\partial \theta} \log L(x|\theta) \right)^2 L(x|\theta) \\ &= \int dx S(\theta)^2 L(x|\theta) \end{aligned}$$

QED

Transformation of Fisher Information

↳ $I(\theta)$ is not invariant under change in model parameters.

$$\phi = \phi(\theta) \rightarrow I(\theta) = \int dx L(x|\theta) \left(\frac{\partial}{\partial \theta} \log L(x|\theta) \right)^2$$

$$\rightarrow I(\phi(\theta)) = \int dx L(x|\phi(\theta)) \left(\frac{\partial \phi}{\partial \theta} \right)^2 \left(\frac{\partial}{\partial \theta} \log L(x|\phi(\theta)) \right)^2$$

$$I(\phi) = I(\theta) \left(\frac{\partial \phi}{\partial \theta} \right)^2$$

Using Fisher Information in experiment design:

Consider Taylor expansion of log likelihood:

$$\log L(x|\theta) = \log L(x|\theta_{\text{true}}) + (\theta - \theta_{\text{true}}) \frac{\partial \log L(x|\theta)}{\partial \theta} \Big|_{\theta=\theta_{\text{true}}} + \frac{1}{2} (\theta - \theta_{\text{true}})^2 \frac{\partial^2 \log L(x|\theta)}{\partial \theta^2} \Big|_{\theta=\theta_{\text{true}}}$$

$$\begin{aligned} E_x [\log L(x|\theta)] &= \underbrace{E_x [\log L(x|\theta_{\text{true}})]}_{\text{Constant}} + (\theta - \theta_{\text{true}}) \underbrace{E_x [S(\theta)]}_{E[S(\theta)] = 0} + \frac{1}{2} (\theta - \theta_{\text{true}})^2 E_x \left[\frac{\partial^2 \log L(x|\theta)}{\partial \theta^2} \right]_{\theta=\theta_{\text{true}}} \\ &= \text{Const} - \frac{1}{2} (\theta - \theta_{\text{true}})^2 I(\theta) \end{aligned}$$

This step cancels out log: $E_x [\log L(x|\theta)]$

$$\exp [E_x [\log L(x|\theta)]] \propto \exp \left(-\frac{1}{2} \frac{(\theta - \theta_{\text{true}})^2}{\sigma^2} \right) \quad \text{where } \sigma = \frac{1}{\sqrt{I(\theta)}} \rightarrow \text{High Information - Narrow width}$$

Not dependent on data can be evaluated without experiment.

Shows likelihood behaves as a (multivariate) Gaussian about θ_{true} .

The Prior, $\pi(\theta)$

Describes our state of knowledge/belief about model parameters before performing the experiment.

Typically parameters are continuous → prior is a probability distribution function.

Key points/Advantages: Encodes prior knowledge/assumptions

Can penalise extreme values (regularisation)

For low data, prior dominates, guiding predictions.

Priors should accurately represent state of knowledge/ignorance

↳ Informative Prior: Strong belief based on prior knowledge (large effect, dominates)

Non-Informative/Ignorance Priors: Assumes no/little prior knowledge (low effect, data dominates)

} Subject to knowledge of analyst.

Examples of Ignorance Priors

Uniform - Invariant under translation

Translation Invariance:

$$\pi(x) dx = \pi(x+\Delta x) d(x+\Delta x) \rightarrow \pi(x) \propto \text{Constant.}$$

Normalisation Condition

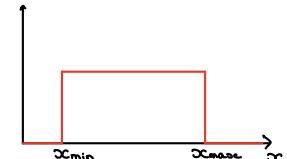
As prior is a p.d.f. → must be properly normalised

↳ For Uniform distribution: Requires x_{\min} and x_{\max} defined. [↓ Omitted - improper prior]

If ignorant - Select value with negligible support (likelihood)

$$\pi(x) = \frac{1_{(x_{\min}, x_{\max})}(x)}{x_{\max} - x_{\min}}$$

Indicator Function: $1_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{else} \end{cases}$



Log Uniform - Invariant under change units/scaling.

Note: Special case Jeffreys Prior

Scaling Invariance:

$$\text{TC}(\alpha) d\alpha = \text{TC}(\alpha \cdot c) d(\alpha \cdot c) \rightarrow \text{TC}(\alpha) \propto \frac{1}{\alpha}$$

Comparison:

→ Rather than being uniform in linear space it is uniform in log space

$\text{TC}(\log \alpha) \propto \text{constant}$ → useful for parameters that span multiple orders of magnitude.

Gives equal probability mass to each order of magnitude.

↳ ie 0.1-1, 1-10, 10-100 all have equal probability

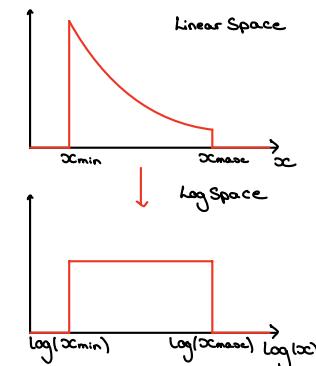
ensure both small and large values fairly represented.

Weaknesses of Uniform:

↳ overrepresents biases large values

For greater weighting for 100-1000 than 0.1-1 due to uniform in linear space

$$\text{TC}(\alpha) = \frac{1}{\alpha \log(\alpha_{\max}/\alpha_{\min})}$$



Jeffrey Priors

Issues with Uniform Priors when reparametrised:

For two equivalent but equally valid parameters α and β related by invertible transform: $\beta = \phi(\alpha)$

Analysis A: Parameter α

Likelihood, $f(\alpha | \alpha)$

Uniform Prior, $\text{TC}_\alpha(\alpha)$

$\xrightarrow{\phi}$ Analysis B: Parameter $\beta = \phi(\alpha)$

Likelihood $f(\beta | \beta) = g(\beta | \phi(\alpha))$

Does not convert to uniform

Induced density $| \frac{d\beta}{d\alpha} | \text{TC}_\alpha(\alpha)$

Inconsistency → Uniform prior not uniform/ignorance prior

Solution:

Jeffrey Prior on respective parameters

$$\left. \begin{array}{l} \text{Square root of } \text{TC}_\alpha(\alpha) \propto \sqrt{I(\alpha)} \\ \text{Fisher Information } I(\alpha) \end{array} \right\} \text{TC}_\alpha(\alpha) \propto \sqrt{I(\alpha)}$$

Advantage: Consistency under parameter transformation

'True non-informative prior'

Proof: Must conserve $\text{TC}(\alpha) d\alpha = \text{TC}(\beta) d\beta$

$$\begin{aligned} \text{TC}(\beta) &= \text{TC}(\alpha) \frac{d\alpha}{d\beta} \propto \sqrt{I(\alpha) \left(\frac{d\alpha}{d\beta} \right)^2} \\ &\propto \sqrt{I(\beta)} \end{aligned}$$

By transformation
property of fisher
information

Example-Jeffrey Prior Selection

For $\alpha \in \mathbb{R}$ with likelihood $L(\alpha | N, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\alpha-N)^2}{2\sigma^2}\right)$

Jeffrey Prior for N : $\text{TC}(N) \propto \sqrt{I(N)}$

$$I(N) = \int d\alpha L(\alpha | N, \sigma) S(N)^2 \quad \text{where } S(N) = \frac{\partial}{\partial N} \log(L(\alpha | N, \sigma)) = \frac{\alpha - N}{\sigma^2}$$

$$= \int d\alpha L(\alpha | N, \sigma) \left(\frac{\alpha - N}{\sigma^2} \right)^2 = \int d\alpha \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\alpha-N)^2}{2\sigma^2}\right) \left(\frac{\alpha - N}{\sigma^2} \right)^2$$

$$= \frac{1}{\sigma^2}$$

$$\therefore \text{TC}(N) \propto \sqrt{I(N)} \propto \text{Const}$$

- For this example
Jeffrey Prior, $\text{TC}(N)$ - Uniform

Similarly for prior on σ :

$$\text{TC}(\sigma) = \frac{1}{\sigma}$$

Jeffrey Prior $\text{TC}(\sigma)$ - Log uniform.

Conjugate Priors (update knowledge without change of distribution)

A form of prior in which posterior distribution belongs to the same family as the prior after updating with data (likelihood).

$$\frac{P(\theta|x)}{\text{Posterior}} \propto \frac{L(x|\theta)}{\text{Likelihood}} \frac{P(\theta)}{\text{Prior}}$$

$$\text{Prior} \sim \text{Distribution } A(\alpha) \xrightarrow{\text{Likelihood}} \text{Posterior} \sim \text{Distribution } (\alpha)$$

The likelihood: updates state of knowledge using data

For Conjugate Prior: Simply update parameters of distribution $\alpha \rightarrow \alpha'$

Unique to each likelihood:

i.e. Poisson Likelihood \rightarrow Gamma Prior

Proof in example

Advantages

- Simple form of Bayesian update (change parameters)
- Allow Calculations of Posterior to be Closed form
- Useful for iterative updates.

Issues:

- For realistic/complex likelihoods often unavailable.

Effects/Validity of Varying Priors

- The choice of prior effects posterior and thus results

Multiple results given subjective beliefs/knowledge

Posterior does not give single answer but 'updates state of belief'

'degree of belief' of outcome \rightarrow 'willingness to bet' on outcome

Role of Prior vs Data

Low/weak data limit: Prior Dominated Regime

Likelihood function is broad, does not constrain posterior

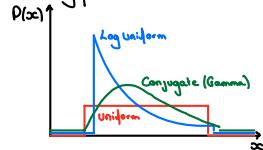
Prior has larger effect \rightarrow Posterior sensitive to subjective choice

High/Strong data limit: Likelihood Dominated Regime

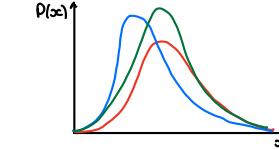
Likelihood function is sharply peaked, constraining/concentrating posterior

Prior has small effect \rightarrow Posterior robust to subjective choice

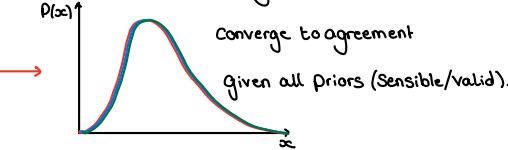
Different Priors:



Effect on Posterior (weak data)



Effect on Posterior (strong data)



Iterative Bayesian Inference.

Update Cycle

Can be performed

Iteratively

- Prior - represents our knowledge before data
- Likelihood - provides update information (applied with evidence to normalise)
- Posterior - represents updated knowledge after data/evidence

Iterative Process

update beliefs sequentially: Step n's posterior \rightarrow Step n+1's prior

- Allows knowledge to accumulate

1) Initial Prior, $P(\theta)$

2) 1st Update with data, D_1 , $P(\theta|D_1) = \frac{P(D_1|\theta)P(\theta)}{P(D_1)}$

3) 2nd Update with data, D_2 , $P(\theta|D_1, D_2) = \frac{P(D_2|\theta)P(\theta|D_1)}{P(D_2)}$

4) Repeating for D_3, D_4 , $P(\theta|D_1, D_2, D_3) = \frac{P(D_3|\theta)P(\theta|D_1, D_2)}{P(D_3)}$

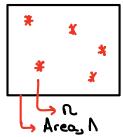
IS EQUIVALENT?

Advantages:

- 1) System only has to retain current posterior
ie for large datasets where storing all data is infeasible
 - 2) Ensures constant updating of model for data arriving sequentially (ie real world)
 - 3) Adaptive Bayesian Updating:
 - ↳ For distributions that change over time, constantly updating / tracking distribution dynamically
- IS there a higher weighting to new evidence → Should there be for large data sets ie dilution/forgetting.

Astronomy Example:

Wish to determine number density, S , stars per square degree:



Likelihood: $L(\text{Data} | \text{Model}) \approx L(n | S) = \frac{(AS)^n e^{-AS}}{n!}$

$(n \sim \text{Poisson})$

Prior Choice 1: Uniform $0 < S < 20 \text{ deg}^{-2}$

$$\text{TT}(S) = \frac{\mathbb{1}_{(0,20)}(S)}{20}$$

Associated Evidence $Z(n) = \int_{-\infty}^{\infty} P(n | S) \text{TT}(S) dS$

Posterior

$$P(S | n) = \frac{(AS)^n e^{-AS}}{20 n!} \frac{\mathbb{1}_{(0,20)}(S)}{Z}$$

$$= \int_{-\infty}^{\infty} \frac{(AS)^n e^{-AS}}{n!} \frac{\mathbb{1}_{(0,20)}(S)}{20} dS = \int_0^{20} \frac{(AS)^n e^{-AS}}{20 n!} dS$$

Prior Choice 2: Log Uniform

$$\text{TT}(S) = \frac{\mathbb{1}_{(S_{\min}, S_{\max})}(S)}{S \log(S_{\max}/S_{\min})}$$

Associated Evidence $Z(n) = \int_{-\infty}^{\infty} P(n | S) \text{TT}(S) dS$

Posterior:

$$P(S | n) = \frac{(AS)^n e^{-AS}}{n! Z} \cdot \frac{\mathbb{1}_{(S_{\min}, S_{\max})}(S)}{S \log(S_{\max}/S_{\min})}$$

$$= \int \frac{(AS)^n e^{-AS}}{n!} \frac{\mathbb{1}_{(S_{\min}, S_{\max})}(S)}{S \log(S_{\max}/S_{\min})} dS$$

Prior Choice 3: Conjugate Prior

For Poisson likelihood: Gamma Prior

$$\text{TT}(S) = \begin{cases} AS^{k-1} e^{-AS} / \Gamma(k) \theta^k & \text{if } S > 0 \\ 0 & \text{otherwise} \end{cases}$$

Mean: $E[\alpha] = k\theta$

Variance: $V[\alpha] = k\theta^2$

$\Gamma(k)$ - Gamma Function

Posterior

$$P(S | n) = \frac{(AS)^n e^{-AS}}{n!} \cdot \frac{AS^{k-1} e^{-AS/\theta}}{\Gamma(k) \theta^k} = \frac{AS^{n+k-1} e^{-AS(\theta+1)}}{n! \Gamma(k) \theta^k} \quad \left. \begin{array}{l} \text{Gamma Distribution:} \\ k' = n+k \quad \theta' = \frac{\theta}{A\theta+1} \end{array} \right\}$$

Iterative Bayesian Inference Example:

Using conjugate prior for simple iteration steps.

Initial Prior $S \sim \text{Gamma}(k, \theta)$ $\rightarrow k' = k + n_1, \theta' = \frac{\theta}{A_1 \theta + 1}$

Initial Posterior $S | n_1 \sim \text{Gamma}(k', \theta')$

Second Posterior $S | n_1, n_2 \sim \text{Gamma}(k'', \theta'')$

$$k'' = k' + n_2, \theta'' = \frac{\theta'}{A_2 \theta' + 1} \rightarrow k'' = k + n_1 + n_2, \theta'' = \frac{\theta}{(A_1 + A_2) \theta + 1}$$

Fisher Information for Experiment Selection:

Survey 1: Area, A_1 finding n_1 stars

Survey 2: Area, $A_2 (> A_1)$ finding n_2 stars (but only with probability detection - p)

Fisher Information Survey 1:

Likelihood $n_1|S \sim \text{Poisson}(A_1 S)$

$$P(n_1|S) = \frac{(A_1 S)^{n_1} \exp(-A_1 S)}{n_1!} \quad \text{Score} \rightarrow S_1(S) = \frac{\partial}{\partial S} \log P(n_1|S) = \frac{\partial}{\partial S} (n_1 \log(A_1 S) - A_1 S - \log n_1!) = \frac{n_1}{S} - A_1$$

Fisher Information

$$\begin{aligned} I_1(S) &= \sum_{n_1=0}^{\infty} L(n_1|S) S_1^2(S) = \sum_{n_1=0}^{\infty} L(n_1|S) \left(\frac{n_1}{S} - A_1 \right)^2 = \sum_{n_1=0}^{\infty} L(n_1|S) \left(\frac{n_1^2}{S^2} - \frac{2n_1 A_1}{S} + A_1^2 \right) \\ &\quad \text{Discrete-THussum} \\ &= \sum_{n_1=0}^{\infty} L(n_1|S) \left(\frac{n_1^2}{S^2} - \frac{2n_1 A_1}{S} + A_1^2 \right) = \frac{1}{S^2} \sum L(n_1|S) n_1^2 - \frac{2A_1}{S} \sum L(n_1|S) n_1 + A_1^2 \\ &= \frac{E[n_1^2]}{S^2} - \frac{2A_1}{S} E[n_1] + A_1^2 = \frac{\text{Var}(n_1) + \bar{n}_1^2}{S^2} - \frac{2A_1}{S} \bar{n}_1 + A_1^2 \end{aligned}$$

Poisson Mean: $\bar{n}_1 = A_1 S$ Variance: $A_1 S$

$$= \frac{A_1 S - A_1^2 S^2}{S^2} - \frac{2A_1}{S} A_1 S + A_1^2 = \frac{A_1}{S} + A_1^2 - 2A_1^2 + A_1^2 = \frac{A_1}{S} \quad \text{ie Information} \propto \text{Area.}$$

Fisher Information Survey 2:

True no. stars in A_2 , m : $m|S \sim \text{Poisson}(mS)$ $\rightarrow m$, a latent (unobserved variable)

Stars observed, n_2 : $n_2|m \sim \text{Binomial}(m, p)$

$$\begin{aligned} L(n_2|S) &= \sum_{m=n_2}^{\infty} P(n_2|m) P(m|S) = \sum_{m=n_2}^{\infty} m C_{n_2} p^{n_2} (1-p)^{m-n_2} \cdot \frac{(A_2 S)^m \exp(-A_2 S)}{m!} \\ &\quad \text{m} \geq n_2 \\ &= \frac{(A_2 p S)^{n_2} \exp(-A_2 S)}{n_2!} \sum_{N=0}^{\infty} \frac{((1-p) A_2 S)^N}{N!} \exp((1-p) A_2 S) \\ &= \frac{(A_2 p S)^{n_2} \exp(-A_2 p S)}{n_2!} \quad \text{Itself a Poisson} \sim \text{Poisson}(A_2 p S) \end{aligned}$$

Thus using previous results: Score: $S_2(S) = \frac{n_2}{S} - A_2$ Information: $I_2(S) = \frac{A_2 p}{S}$

If $A_2 p > A_1$, $I_2 > I_1$

Survey 2 preferable

Fisher Information from Multiple Experiments:

Information adds for independent experiments.

Given independence: $L(n_1, n_2|S) = L(n_1|S) L(n_2|S)$

Then: $I(S) = I_1(S) + I_2(S)$

Extracting Information from Posterior

For a high dimensional problem, Posterior is complete description about state of knowledge of model parameters, Θ_N

Marginal Distributions

Often wish to integrate out/marginalise 'nuisance parameter'

For $m+n$ parameters $\{\Theta_1, \Theta_2, \dots, \Theta_N, \phi_1, \phi_2, \dots, \phi_m\}$:

Marginalise over m parameters leaving N dim 'marginal posterior'

$$p(\Theta_1, \Theta_2, \dots, \Theta_N | \mathbf{x}) = \int d^m \phi P(\Theta_1, \Theta_2, \dots, \Theta_N, \phi_1, \phi_2, \dots, \phi_m)$$

Integrating over all parameters:

$$\int d^N \Theta_N P(\Theta_N | \mathbf{x}) = \text{Const}$$

If correctly normalised: $\text{Const} = 1$

If unnormalised posterior: $P(\Theta_N | \mathbf{x}) \propto L(\mathbf{x} | \Theta_N) \pi(\Theta_N)$

$$\text{Const} = \text{Bayesian Evidence, } Z = \int d^N \Theta_N P(\Theta_N | \mathbf{x}) = \int d^N \Theta_N L(\mathbf{x} | \Theta_N) \pi(\Theta_N)$$

Thus 'evidence' \approx 'marginal likelihood'

Point Estimates

Summarising posterior with single/few value Summary Statistics.

All 'loose information' from posterior which is full information.

Posterior Mean

For Θ_1 :

$$\langle \Theta_1 \rangle = \int d^N \Theta_N \Theta_1 P(\Theta_N | \mathbf{x}) \quad \text{from full } N\text{-dimensional posterior}$$
$$= \int d\Theta_1 \Theta_1 P(\Theta_1 | \mathbf{x}) \quad \text{from 1-D marginal posterior}$$

using defn of marginal

Posterior Mode - Θ_N^{MAP} Maximum a posteriori (MAP)

$$\frac{\partial P(\Theta_N | \mathbf{x})}{\partial \Theta_N} \Big|_{\Theta_N = \Theta_N^{\text{MAP}}} = 0$$

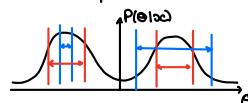
Marginal Posterior - integration over parameters can shift peak away from Θ_N^{MAP}

Median

$$\bar{\Theta} \text{ s.t. } \int_{-\infty}^{\bar{\Theta}} d\Theta P(\Theta | \mathbf{x}) = \frac{1}{2} \quad \text{Only valid from 1d marginal posterior}$$

Credible Intervals/Regions

A range/region in which model parameters lies with a particular probability level.



Issue: Not Unique

Common Choices in 1D

1) Narrowest (Collection of) Intervals that contain probability mass α .

Not invariant under reparametrization.

2) Highest Density Interval

Interval(s): $\text{Prob}(\Theta \in I(j_\alpha)) > \alpha$ s.t. $I(j_\alpha) = \{\Theta : P(\Theta | \mathbf{x}) > j_\alpha\}$ for largest j_α

3) Equal-tailed Interval:

$\Theta_{\min} < \Theta < \Theta_{\max} : \text{Prob}(\Theta < \Theta_{\min}) = \text{Prob}(\Theta > \Theta_{\max}) = \alpha/2$

Bayesian Computation

Generating Stochastic Samples from Target distribution, P (Posterior)

which can in turn help answer questions about distribution.

Monte Carlo:

Set of Stochastic Samples - n iid random variables - from Target Distribution, $P(x)$

$$x_i \stackrel{iid}{\sim} P \text{ for } i=0, \dots, n-1$$

↳ Typically sample of size 10^3 to 10^6

Used in Statistical Inference:

For P -Posterior: \rightarrow Posterior Samples \rightarrow Bayesian Inference

Insights (Deterministic Quantities):

1) Approximating full (d -dimensional) distribution

↳ d -dim histogram / Kernel density estimate

2) Approximating marginal distributions

↳ Generating lower dimensional histograms.

3) Computing Credible Intervals:

↳ Quantify uncertainties through intervals with given probabilistic content.

4) Approximate Integrals with Monte-Carlo Sum *

$$\langle \psi(x) \rangle = \int d\omega \psi(\omega) P(\omega) \quad \rightarrow \quad S_n = \frac{1}{n} \sum_{i=0}^{n-1} \psi(x_i)$$

$S_n \rightarrow \langle \psi(x) \rangle$ as $n \rightarrow \infty$ x_i : Equally weighted Stochastic Samples.

Error on estimator:

Variance $\propto 1/n$; Error $\propto 1/\sqrt{n}$ \rightarrow Relatively poor performance

↳ However holds for high dimensions \rightarrow becomes best case.

Stochastic Sampling

Constructing a random Sample over Sample Space, X

given probability distribution, P .

$$x \sim P, \quad x \in X$$

Possible (largely) Unfeasible:

Transform Sampling:

Maps from a known distribution Q \rightarrow Target Distribution, P

Given $x \sim Q$, invertible transform $y = f(x)$, $y \sim P$

$$P(y) = \left| \frac{dy}{dx} \right|^{-1} P(x) \quad \frac{dy}{dx} \text{ - Jacobian of Transform}$$

↳ 1D: Inverse CDF, \rightarrow ISSUE: Many dimensions: difficult to define

Inverse CDF Method:

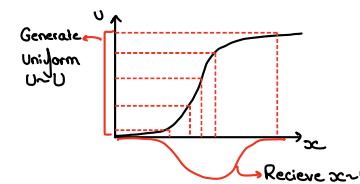
1) Uniform random variable $U \sim U$

2) Function / Transform

$$\text{CDF: } F(x) = \int_{-\infty}^x P(x') dx' = N \quad \text{monotonically increasing } \lim_{x \rightarrow -\infty} F(x) = 0, \lim_{x \rightarrow \infty} F(x) = 1 \rightarrow \text{has unique inverse}$$

PPF: $F^{-1}(u) = x$ Inverse CDF, Percentage Point Function

3) $x \sim P(x)$



Rejection Sampling (accept-reject)

Proposal Distribution, Q → which we can sample from efficiently

Q must have support everywhere P does → $P(x) \leq M Q(x) \forall x \in M$

1) Random Samples from both Proposal, Q and Uniform, U

↳ $(x, Mu) \sim Q, U \sim U$

2) Acceptance Condition

if $Mu < P(x) / Q(x)$: Accept

else : reject

Iterations ≈ M iterations per sample

As Proposal Q closer to target P , $M \rightarrow 1$

↳ Efficient production Samples

$$\text{Proof: } P(x = x \mid N < \frac{P(x)}{M Q(x)}) = \text{Prob}(N < \frac{P(x)}{M Q(x)} \mid x = x) \frac{\text{Prob}(x = x)}{\text{Prob}(N < \frac{P(x)}{M Q(x)})}$$

$$\text{Prob}(N < \frac{P(x)}{M Q(x)} \mid x = x) \rightarrow \frac{P(x)}{M Q(x)} \cdot 1$$

↳ As uniform

$$\text{Prob}(N < \frac{P(x)}{M Q(x)}) \rightarrow \int_{-\infty}^{\infty} dx \frac{P(x)}{M Q(x)} Q(x) = \frac{1}{M}$$

↳ As $x \sim Q$: expectation over x

$$\text{Overall: } P(x = x \mid N < \frac{P(x)}{M Q(x)}) = P(x).$$

Issue: Difficulties finding proposal which

↳ Easy to sample from

Multivariate: Very large M , inefficient

Efficiently encompasses target.

Importance Sampling

↳ Determines deterministic property directly from samples of proposed distribution

↳ Does not determine samples of target distribution

Using a proposed distribution, Q close to P (measured by ϵ) - with Q support everywhere.

For a set of proposal samples: $x_i \stackrel{iid}{\sim} Q \quad i=0, 1, \dots, n-1$

Assign weights: $W_i = \frac{P(x_i)}{Q(x_i)}$ with weighted samples

Use of weighted samples: $\{(x_i, W_i) \mid i=0, 1, \dots, n\}$

$$I = \int dx P(x) \approx \frac{1}{n} \sum_{i=0}^{n-1} W_i P(x_i)$$

↳ weight contribution of $x \sim Q$

↳ property of P → different from rejecting

Measuring quality of proposal:

↳ best when Q is close to P → greater W_i

$$\text{Number of effective samples, } N_{\text{eff}} = \left(\sum_{i=0}^{n-1} W_i \right)^2 / \sum_{i=0}^{n-1} W_i^2$$

$$\text{Efficiency: } \epsilon = \frac{N_{\text{eff}}}{n}$$

Markov Chain Monte Carlo

Overall aim:

1) Construct a markov chain which stationary distribution ≈ target distribution

↳ This will happen after chain reaches equilibrium after 'burn-in' phase.

2) Draw random samples from chain ≈ random samples from distribution

Markov Chains

↳ Overview:

For a given position/stage - transition probability, P defines probability of next move

Only dependent on current position (no memory of history)

After sufficiently evolving: Chain $\sim \pi$, stationary distribution

With correct transition probability P : $\pi \rightarrow \pi$

Markov Chain: Ordered sequence of random points x_0, x_1, \dots, x_n in sample space ($x \in \mathcal{X}$)

that satisfy Markov Property:

$$P(x_{i+1} | x_0, \dots, x_i) = P(x_{i+1} | x_i) = P(x_{i+1}, x_i)$$

Notation

Transition Probability

ie only dependent on previous point in chain

Time-Homogeneous Markov Chains: transition probability do not depend on position in

$$\text{chain/time: } P(x_{k+1} = x | x_k = y) = P(x_{j+1} = x | x_j = y) \quad \forall k, j$$

Time Homogeneous:

Chain reaches Stable Stationary distributions

used to Sample Constant target distributions.

Time Inhomogeneous

Chain distribution evolves dynamically

Non Stationary Underlying distribution (ie Climate Models)

Discrete Representation

For finite sample size, \mathcal{X} - transition probabilities can be described by matrix

$$\begin{array}{l} \text{inhomogeneous} \\ \text{P depends on } k \end{array} \quad P_{ab}^{(k)} = P(x_{k+1} = b | x_k = a) \quad a, b \in \mathcal{X} \quad P - N \times N \text{ matrix where } N = |\mathcal{X}| \\ \text{Each row - Probability distribution of next update} \end{math>$$

Sum of rows = 1

Dissecting Prob Dist - Time Homogeneous Markov Chain

For a given starting position:

$$\begin{aligned} P(x_1 | x_0) &= \int dx_1 P(x_1 | x_0, x_0) P(x_0 | x_0) \quad (\text{Law Total Prob}) \\ &= \int dx_1 P(x_1 | x_0) P(x_0 | x_0) \quad (\text{Markov Property}) \\ &= \int dx_1 P(x_1, x_0) \quad (\text{Transition Probabilities}) \end{aligned}$$

By induction:

$$P(x_i | x_0) = \int dx_{i-1} \int dx_{i-2} \dots \int dx_1 P(x_i, x_{i-1}) P(x_{i-1}, x_{i-2}) \dots P(x_1, x_0)$$

'Distribution i^{th} point given x_0 '

Deciding Transition Probability (Mappings):

Irreducibility: for any starting point $x_0 \in \mathcal{X}$ and for any region $A \subset \mathcal{X}$:

$$\exists n \geq 1, n \in \mathbb{Z}^+ \quad \text{s.t. } \int_A dx_0 P(x_n | x_0) > 0$$

↳ For any starting point, all points in space have a non zero probability
of being reached, albeit not necessarily in 1 step. (go anywhere tendency)

ergodicity

Time-Homogeneous and Irreducible

↳ Chain never stops moving / Converge to point

But chain distribution may converge → Limiting distribution (only Stationary dist).

Limiting Distribution: time homogenous Markov Chain approaches limiting distribution, λ

on Space \mathcal{X} if: $\lim_{n \rightarrow \infty} P(x_n | x_0) = \lambda(x_0)$ → if exists, unique

↳ If Markov Chain is ergodic: influence of x_0 disappears over time

↳ Stationary vs Limited.

If $x_0 \sim \lambda$: Chain starts in limiting distribution (equilibrium)

→ If not drawn from $\lambda(x)$

→ Chain requires 'burn-in period' before converging

x_0 from λ : Integrate (weighted) over all possible values x_0

$$P(x_0) = \int dx_0 P(x_0 | x_0) \quad \lambda(x_0) = \lambda(x_0) \int dx_0 P(x_0 | x_0) = \lambda(x_0)$$

↳ If the Chain Starts in Stationary distribution → remains there

Reality: Do not know $\lambda(x)$ → run for a bit

Equilibrium reached before collecting samples

Stationary Distribution:

For a Markov chain with transition probabilities $P(x';x)$, a distribution

on X is a stationary distribution if:

$$\pi(x') = \int dx \pi(x) P(x',x) \rightarrow \text{Stay in same distribution after step}$$

Stationary vs Limited

Limited \subset Stationary

Stationary: Non unique - if starts in $\pi(x)$, remains there

i.e. Stable distribution depends on starting position, x_0

Limited: Unique - Convergent distribution of chain

i.e. all points converge to this

Ergodicity / Irreducibility

↳ Markov chain has 1 stationary distribution
which is the limiting distribution

Non Ergodicity / Irreducibility

↳ Cannot have limiting distribution.

Detailed Balance: A markov chain with transition probabilities $P(x, x')$,

is said to follow detailed balance wrt π if:

$$\pi(x) P(x, x') = \pi(x') P(x', x)$$

prob flow $x \rightarrow x'$ prob flow $x' \rightarrow x$

↳ if equal - flow of probability is symmetric

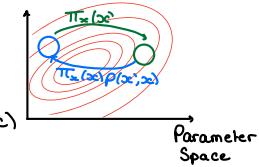
→ no net probability flow

→ stable/stationary

Desire $\text{Prob}(x \in A \text{ and } x' \in B) = \text{Prob}(x \in B \text{ and } x' \in A)$

$$P(x \in A) P(x' \in B | x) = P(x \in B) P(x' \in A | x)$$

$$\int_{x \in A} \pi(x) \int_{x' \in B} P(x', x) = \int_{x \in B} \pi(x) \int_{x' \in A} P(x', x)$$



Lemma #1: If Markov chain satisfies detailed balance wrt π $\rightarrow \pi$ is a stationary distribution of Markov chain.

→ Stricter condition
than necessary.

Proof: DB condition

$$\pi(x) P(x, x') = \pi(x') P(x', x)$$

$$\text{Integrate wrt } x' \text{ over } X: \int_X dx' \pi(x) P(x', x) = \int_X dx' \pi(x') P(x, x')$$

$$\pi(x) = \underbrace{\int P(x' | x) dx'}_1 = \int_X dx' \pi(x') P(x, x')$$

$$\pi(x) = \int_X dx' \pi(x') P(x, x') \rightarrow \text{definition of stationary, QED}$$

Lemma #2: If Markov chain is homogeneous and irreducible and satisfies detailed balance wrt π $\rightarrow \pi$ is unique stationary/Limiting distribution

Result:

Plan to find transition probability $P(x, x')$ for a time-homogeneous, irreducible Markov chain

St. Satisfies detailed balance where stationary distribution, $\pi = P$, Target Distribution

result: Chain will approach P (after evolving for sufficient iterations - burn-in period)

Gibbs Sampling:

Samples are generated by iteratively updating one parameter (dimension) at a time, drawing from its 1D conditional distribution given the current values of all other parameters.

Full-High Dimensional Target Distribution: (d dimensional)

$$P(x^0, x^1, \dots, x^{d-1})$$

1D Conditional for Parameter x^k

$$P(x^k | x^0, \dots, x^{k-1}, x^{k+1}, \dots, x^{d-1}) = \frac{\text{All values, except variable } x^k}{\int dx^k P(x^0, x^1, \dots, x^{d-1})}$$

Notation: $P(x^k | x^{k-1})$
ie $x^{k-1} = x^0, \dots, x^{k-1}, x^{k+1}, \dots, x^{d-1}$

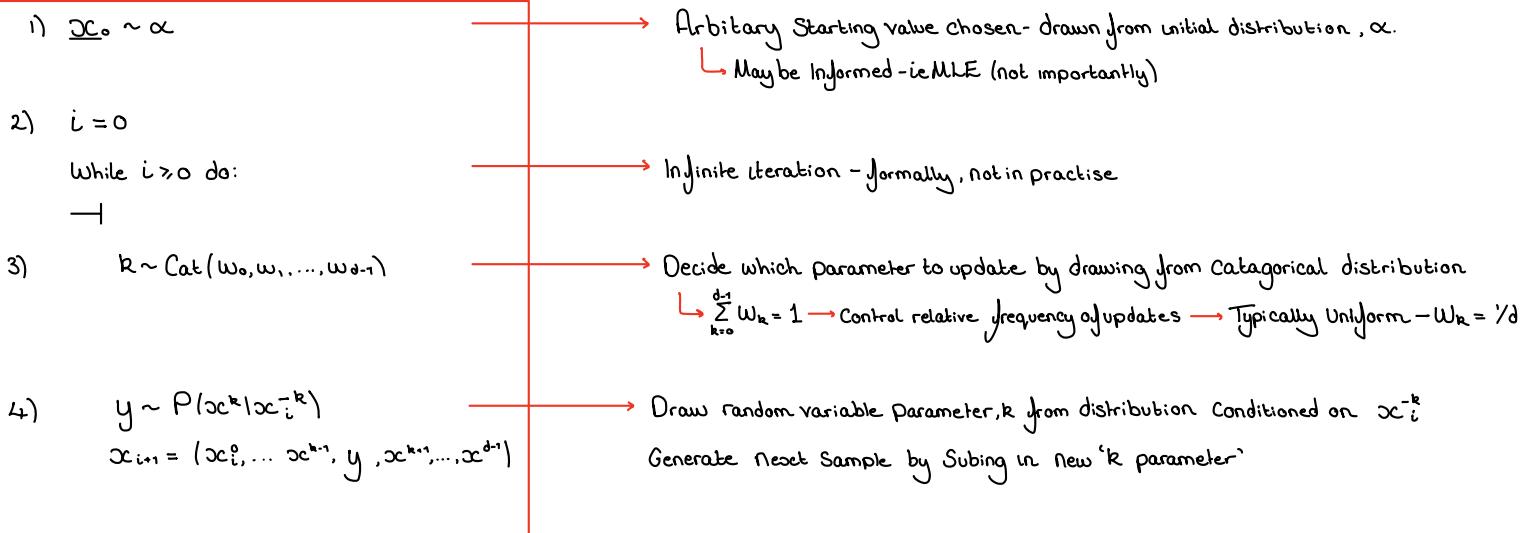
Assumes easily Sample from 1D conditional distribution

↳ Rejection Sampling

Transform Sampling (PPF)

Embedded/Secondary Markov Chain

Gibbs Algorithm Steps



Theorem: Gibbs Algorithm produces Markov Chain x_0, x_1, \dots that

Satisfies detailed balance condition $\pi \propto P$

When would you use non-Uniform weights?

- 1) Assign higher weighting to slower Converging params that mix slowly / don't vary much between successive Iterations
- 2) Some parameters posterior much more peaked than others - do not explore all values frequently
↳ high Samples to correctly estimate Variance.

Proof: Transition probability defined by Gibbs

$$P(y|x) = \sum_k w_k S^{(d-1)} (y^k - x^k) P(y^k | x_{-k})$$

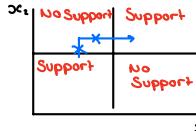
1-D conditional distributions
Sum overall possible
Fixed/Constant-dirac delta dist
ie $\int x^k \neq y^k, P=0$

Using detailed balance eqn:

$$\begin{aligned}
 P(x) P(y|x) &= P(y) P(x|y) \\
 \text{LHS } P(x) P(y|x) &= \sum_k P(x^k | x_{-k}) P(x_{-k}) w_k S^{(d-1)} (y^k - x^k) P(y^k | x_{-k}) \\
 &\quad \xrightarrow{\text{by Dirac Delta: } P(y^k) = P(x^k)} \\
 \text{RHS } P(y) P(y|x) &= \sum_k P(y^k | y_{-k}) P(y_{-k}) w_k S^{(d-1)} (y^k - x^k) P(x^k | y_{-k}) \\
 &= \sum_k P(y^k | x_{-k}) P(x_{-k}) w_k S^{(d-1)} (y^k - x^k) P(x^k | x_{-k}) \\
 &= \sum_k P(x^k | x_{-k}) P(x_{-k}) w_k S^{(d-1)} (y^k - x^k) P(y^k | x_{-k}) \quad \text{QED } \text{LHS} = \text{RHS}.
 \end{aligned}$$

↳ Detailed balance ensures Stationary distribution but doesn't guarantee irreducibility / limiting distribution. ie

↳ Typically produces irreducible markov chain - not guaranteed.



ie No gibbs chain will allow transition (horizontal/vertical) between Supported regions

Inefficient Samples:

↳ by only updating 1 parameter at a time:

Common to take every $(d\text{-dim})^n$ point in Gibbs Chain

→ does not guarantee no duplicates due to $k \sim \text{Cat}(\dots)$ but reduces them (Similar to Sweep but not guaranteed new value cat)

Gibbs Sampling Example

2D target distribution: $P(x, y) = y \exp(-[xy + y])$ with $0 < x, y < \infty$

$$\underline{x \text{ Conditional: }} P(x|y) = \frac{\exp(-y-x)}{\int_0^\infty \exp(-y-x) dx} = y \exp(-y-x) \quad (\text{exponential shape parameter of 1})$$

$$\underline{y \text{ Conditional: }} P(y|x) \propto y \exp(-[xy + y]) \rightarrow y|x \sim \text{Gamma}\left(k=2, \theta=\frac{1}{x+1}\right) \quad \text{Gamma}(k, \theta) = \frac{y^{k-1} \exp(-y/\theta)}{\theta^k \Gamma(k)}$$

0th Iteration: $(x_0, y_0) = (1, 1)$ - arbitrary

1st Iteration: Randomly choose $k \in \{0, 1\}$ → Pick $k=0$, x -coordinate:

$$\text{Draw } x \sim \frac{\exp(-y_0 x)}{y_0} \rightarrow (x_1, y_1) = (x_0, y_1 = y_0)$$

2nd Iteration: Randomly choose $k \in \{0, 1\}$ → Pick $k=1$, y -coordinate:

$$\text{Draw } y \sim \text{Gamma}\left(2, \frac{1}{x_1+1}\right) \rightarrow (x_2, y_2) = (x_2 = x_1, y_2)$$

etc:

Gibbs 'Sweep' Sampling

↳ Every Component gets updated (in a fixed order) at every iteration (in a Sweep)

↳ Components updated by drawing from 1-d Conditional

Conditioned on parameter values updated so far

Does Not Satisfy detailed balance → Does Converge to target distribution.

Lemma: Gibbs Sweep algorithm produces a Markov chain x_0, x_1, \dots that has the target distribution P as unique stationary distribution.

Algorithm:

$$x_0 \sim a$$

$$i = 0$$

while $i > 0$ do

for $k=0$ to $d-1$ do:

$$y^k \sim P(x^k | y^0, \dots, y^{k-1}, x^{k+1}, \dots, x^{d-1})$$

end for

→ Already updated in this Sweep

$$x_{i+1} = (y^0, \dots, y^{d-1})$$

$$i = i + 1$$

Proof:

For 2d Case: $P(x) = P(x^0, x^1)$ be the target distribution

Transition: $(x^0, x^1) \rightarrow (y^0, x^1) \rightarrow (y^0, y^1)$

$$P(y^0, y^1, [x^0, x^1]) = P(y^0|x^1) P(y^1|y^0) = \frac{P(y^0, x^1)}{\int da P(a, x^1)} \cdot \frac{P(y^0, y^1)}{\int db P(y^0, b)}$$

Defn Stationary dist:

$$\pi(x^1) = \int dx^0 \pi(x^0) P(x^0, x^1)$$

$$\begin{aligned}
 \int dx P(x) P(y, x) &= \int dx^0 \int dx^1 P(x^0, x^1) \rho([y^0, y^1], [x^0, x^1]) \\
 &= \frac{P(y^0, y^1)}{\int dx P(y, x)} \int dx^1 \frac{P(y^0, x^1)}{\int dx P(x, x^1)} \int dx^0 P(x^0, x^1) \\
 &= P(y^0, y^1) = P(y) \quad \therefore P(y) = \int dx P(x) P(y, x) \quad \text{Stationary Conditioned Satisfied.} \\
 &\quad \text{Can be extended to } d \text{ dimension.}
 \end{aligned}$$

Blocked Gibbs Sampling

- ↳ Groups multiple parameters into blocks of parameters, $\Theta^N = (x^1, \dots, x^d)$
- allows multiple parameters (of block) to be updated at once
- ↳ Requires: blocks, Θ^N to be chosen s.t. possible to sample from higher dimensional Conditional probability, $P(\Theta^N | \Theta^{-N})$

Gibbs Sampling with Ordered overrelaxation

- ↳ Randomly walking - Slow exploration of parameter space
- OOR → introduces negatively correlated transitions
 - ↳ encourages chain to make larger moves

Order Statistics

- ↳ Draw multiple possible values from given sample
- Sort them and select new value based on relative position (enforcing negative correlation)

Metropolis - Hastings Sampling

Advantage: does not require 1D conditional distribution.

The Proposal Distribution, $Q(y|x)$

A d -dimensional distribution which generates candidates for next chain position.

- ↳ e.g. commonly: $y \sim N(x, \Sigma)$ - Gaussian (Symmetric proposal)

Following Probabilistic rule:

Proposal is either accepted ($x_{i+1} = y$) or rejected ($x_{i+1} = x_i$)

Acceptance Probability:

$$a = \min \left(1, \frac{P(y)}{P(x)} \frac{Q(x|y)}{Q(y|x)} \right)$$

Overall transition probabilities:

- ↳ mixture of proposal distribution, $Q(y|x)$ with fraction a (acceptance), and delta function with fraction $(1-a)$ (rejection).
- ↳ $P(y|x) = a(y|x) Q(y|x) + \delta^d(y-x) \int dy' [1 - a(y',x)] Q(y'|x)$
- ↳ $P(\text{accepting} | \text{proposed } y) P(\text{proposing } y)$
- ↳ Sum over all rejected proposals

Algorithm:

$$x_0 \sim x$$

$$i = 0$$

while $i > 0$:

$$y \sim Q(y|x_i)$$

$$a = P(y) Q(x_i|y) / P(x_i) Q(y|x_i)$$

$$U \sim U$$

$$\text{if } U < a:$$

$$x_{i+1} = y$$

else:

$$x_{i+1} = x_i$$

Note: MH only depends on ratio $P(y)/P(x)$

↳ Thus does not require normalised Target distribution

Do not require evidence (only L, π_C, P)

Stuck points:

Repeat entries occur when proposal rejected ($x_{i+1} = x_i$)

↳ Acceptance fraction: $\frac{\text{No. accepted proposals}}{\text{Total no. iterations}}$

Intuition:

'if' uphill in distribution: $P(y) > P(x) \rightarrow a(y|x) = 1$

'if' 'little downhill' in distribution: $P(y) < P(x) \rightarrow a(y|x) < 1$

↳ Accepted given portion of time.

'favours uphill tendency but can make downhill steps'

Metropolis Sampling:

A special case of MH where proposal is symmetric, i.e.:

$$Q(y|x) = Q(x|y)$$

$$a(y|x) = \min \left(1, \frac{P(y)}{P(x)} \right)$$

The case for gaussian with identity covariance

Theorem: MH algorithm produces Markov Chain x_0, x_1, \dots that

Satisfies the detailed balance Condition $\pi = P$

Proof:

$$\pi(x) P(y, x) = \pi(y) P(x, y)$$

$$\hookrightarrow P(x) a(y, x) Q(y|x) + \delta^d(y-x) P(x) \int dy [1 - a(y, x)] Q(y|x) = P(y) a(x, y) Q(x|y) + \delta^d(x-y) P(y) \int dy [1 - a(y, x)] Q(y|y)$$

by $\delta^d(y-x)$; $y \leftrightarrow x$ thus equivalent and cancel

Must Show: $P(x) a(y, x) Q(y|x) = P(y) a(x, y) Q(x|y)$

By definition: $a = \min(1, \frac{P(x) Q(y|x)}{P(y) Q(x|y)})$ let $a(x, y) = \frac{P(x) Q(y|x)}{P(y) Q(x|y)}$ $\rightarrow a(y, x) = 1$
due to 'min' term.

Overall Satisfied: $\left[\text{Similarly if } a(y, x) = \frac{P(y) Q(x|y)}{P(x) Q(y|x)} \text{ and } a(x, y) = 1 \right]$

Gibbs \rightarrow MH Algorithms

Gibbs Transition Probability $Q(y|x) = \sum w_k \delta^{(d-1)}(y^k - x^k) P(y^k|x^k)$

MH Acceptance Junction: $a(y, x) = \min\left(1, \frac{P(y) Q(x|y)}{P(x) Q(y|x)}\right)$

(for Gibbs)

$$\begin{aligned} &= \min\left[1, \frac{\sum_k w_k \frac{P(x^k|y^k) \delta^{(d-1)}(y^k - x^k)}{P(x)}}{\sum_k w_k \frac{P(y^k|x^k) \delta^{(d-1)}(x^k - y^k)}{P(y)}} \right] \\ &= \min\left[1, \frac{\sum_k w_k \frac{\delta^{(d-1)}(y^k - x^k)}{P(x^k)}}{\sum_k w_k \frac{\delta^{(d-1)}(x^k - y^k)}{P(y^k)}} \right] = 1 \end{aligned}$$

P(x) = $P(x^k|x^k) P(x^k)$
P(y) = $P(y^k|y^k) P(y^k)$

Always Unity / always accepted.

Gibbs — Special case of MH where one uses specific proposals built from the 1D Conditionals of the target

— with an acceptance of unity.

Limitations of Gibbs and MH

\hookrightarrow Exhibit random walk behaviour \rightarrow follows diffusion type behaviour

May take many iterations to 'diffuse' slowly around sample space

Diffusion Equator

'Root Mean Squared' Distance moved $\sim \sqrt{N_{\text{steps}}}$

'Stuck' in Extreme regions

Gibbs — Large No. Small Steps due to 1D conditions

MH — Large Proposal dist — Large No. rejections

Small Proposal dist — Small Steps

\hookrightarrow especially bad for high dimensions

Hamiltonian/Hybrid Monte Carlo

\hookrightarrow Exploits gradient information about the target distribution $(\nabla \log P)$ \rightarrow Treats distribution as potential well.

Proposes a transition to a new position y by drawing a random momentum $p \sim Q$ and integrating Hamiltonian equations forward.

✓: Only requires a function proportional to the target PDF. (not normalised)

X: require gradient of Target distribution wrt parameters.

\hookrightarrow $-\log P(x)$ wrt x : Can only be applied to smooth target distributions.

Hamiltonian Dynamics

Positions - Points $x \in X$ in Sample Space

Potential/ Log Position Distribution (Explain- Logarithm of target PDF)

$$\log P(x) = -E(x) + \text{const}$$

$E(x)$ - potential

Force (acting on a particle)

$$\text{Force} = -\nabla E(x)$$

Momenta -

Variables ($p \in T_x X$) in tangent Space. - Same dimensions as x

Kinetic/ Log Momenta Distribution, Q on Tangent Space, $T_x X$

Can be chosen freely but easy to sample from - ie multivariate gaussian

$$\log Q(p) = -K(p) + \text{const}$$

Kinetic Energy, $K(p)$

$$K(p) = \frac{1}{2} p^T M^{-1} p$$

M - mass matrix, Positive Semidefinite

Phase Space, (x, p)

Position, momentum pairs - dimensionality: x_2 original x

Canonical Distribution, R

distribution on phase space (-ve of Hamiltonian)

$$\log R(x, p) = -H(x, p) + \text{const}$$

Hamiltonian: $H(x, p) = E(x) + K(p)$

$$H(x, p) = -(\log Q(p) + \log P(x)) + \text{const}$$

$$= -(\log \{Q(p)P(x)\}) + \text{const'}$$

Sampling

\hookrightarrow Sample $(x, p) \sim R$ in 2d dimensions (rather than $x \sim P$)

Discard momentum, p variable

$$\text{Equivalence as: } P(x) = \int dp R(x, p)$$

Time Parameters: (introduced parameters)

Position and Momentum premuted to function of $t \rightarrow x(t), p(t)$

$x(t), p(t)$ evolve by Hamiltonian Dynamics:

$$\frac{dx^k}{dt} = \frac{\partial H}{\partial p^k} = (M^{-1} \cdot p)^k$$

$$\frac{dp^k}{dt} = -\frac{\partial H}{\partial x^k} = -\nabla E(x)$$

Deterministic Map, T_s

$$T_s: (x(t), p(t)) \longrightarrow (x(t+s), p(t+s))$$

Properties of Hamiltonian Dynamics

\hookrightarrow Reversibility

\rightarrow Hamiltonian dynamics is reversible: T_s has an inverse T_{-s}

$$\left. \begin{array}{l} T_s(x(t), p(t)) \longrightarrow x(t+s), p(t+s) \\ T_{-s}(x(t+s), -p(t+s)) \longrightarrow x(t), -p(t) \end{array} \right\} \begin{array}{l} \text{reverse momenta, get back to the same} \\ \text{initial point with same momenta (in reverse)} \end{array}$$

Volume Preservation/Liouville's Theorem

\rightarrow Consider all points (x, p) in region of Phase Space, A (volume V_A)

B is the image of A under T_s ; the its volume $V_B = V_A$

Hamiltonian Conservation

Hamiltonian dynamics Conserved Hamiltonian ($\frac{dH}{dt} = 0$)

Numerical Integration/Discrete Steps

↳ exact Solutions often not feasible

Must use Symplectic Integrator → Exactly time reversible
Volume Preserving.

Commonly: Leapfrog Method (Symplectic)

1) Half Step in momentum: $\rho(t + \frac{1}{2}\Delta t) = \rho(t) + \frac{1}{2}\Delta t \frac{d\rho}{dt}$
 $= \rho(t) - \frac{1}{2}\Delta t \nabla E \Big|_{x=x(t)}$ → Energy term evaluated at initial position, $x(t)$

2) Full Step in Position: $x(t + \Delta t) = x(t) + \Delta t \frac{dx}{dt}$
 $= x(t) + \Delta t (M^{-1} \cdot \rho(t + \frac{1}{2}\Delta t))$ → Momentum term evaluated at half step

3) Half Step in momentum: $\rho(t + \Delta t) = \rho(t + \frac{1}{2}\Delta t) + \frac{1}{2}\Delta t \frac{d\rho}{dt}$
 $= \rho(t + \frac{1}{2}\Delta t) - \frac{1}{2}\Delta t \nabla E \Big|_{x=x(t + \Delta t)}$ → Energy term evaluated at $\frac{1}{2}$ step position, $x(t + \frac{1}{2}\Delta t)$

Properties of Leapfrog Dynamics (match that of Hamiltonian).

↳ Reversibility:

Leapfrog dynamics are exactly reversible — for all Δt

Forwards:

$$x(t), \rho(t) \rightarrow x(t + \Delta t), \rho(t + \Delta t)$$

$$1) \rho(t + \frac{\Delta t}{2}) = \rho(t) - \frac{\Delta t}{2} \nabla E(x(t))$$

$$2) x(t + \Delta t) = x(t) + \Delta t M^{-1} \rho(t + \frac{\Delta t}{2})$$

$$3) \rho(t + \Delta t) = \rho(t + \frac{\Delta t}{2}) - \frac{\Delta t}{2} \nabla E(x(t + \Delta t))$$

Reverse - Reverse Momentum $\rho(t + \Delta t) \rightarrow -\rho(t + \Delta t)$

$$1) \rho'(t + \frac{\Delta t}{2}) = -\rho(t + \Delta t) - \frac{\Delta t}{2} \nabla E(x(t + \Delta t))$$

 $= -\rho(t + \frac{\Delta t}{2}) - \frac{\Delta t}{2} \nabla E(x(t + \Delta t)) - \frac{\Delta t}{2} \nabla E(x(t + \Delta t))$
 $= -\rho(t + \frac{\Delta t}{2})$

$$2) x'(t) = x(t + \Delta t) + \Delta t M^{-1} \rho'(t + \frac{\Delta t}{2})$$

 $= x(t) + \Delta t M^{-1} \rho(t + \frac{\Delta t}{2}) + \Delta t M^{-1} (-\rho(t + \frac{\Delta t}{2})) = x(t)$

$$3) \rho'(t) = \rho(t + \frac{\Delta t}{2}) - \frac{\Delta t}{2} \nabla E(x(t))$$

 $= -\rho(t + \frac{\Delta t}{2}) - \frac{\Delta t}{2} \nabla E(x(t)) = -\rho(t) + \frac{\Delta t}{2} \nabla E(x(t)) - \frac{\Delta t}{2} \nabla E(x(t))$
 $= -\rho(t)$

$$x(t + \Delta t), -\rho(t + \Delta t) \rightarrow x(t), -\rho(t) \text{ QED.}$$

Volume Preserving (exactly)

Note all three steps are 'Shear transformations':

↳ each variables (x or ρ) changes by an amount that depends linearly on the other

e.g. Step 2: $(x) \rightarrow (x') = (x + \Delta t f(\rho))$ with $\text{Jacobi} \ \mathcal{J} = \begin{vmatrix} \frac{\partial x'}{\partial x} & \frac{\partial x'}{\partial \rho} \\ \frac{\partial p'}{\partial x} & \frac{\partial p'}{\partial \rho} \end{vmatrix} = \begin{vmatrix} 1 & \Delta t \frac{\partial f}{\partial \rho} \\ 0 & 1 \end{vmatrix} = 1$

↳ All steps - Unit determinant. Thus overall volume preserving

Hamiltonian Conservation (approximately for Small Δt)

Not perfect conservation due to errors from discrete steps.

$$H(x(t + \Delta t), \rho(t + \Delta t)) = H(x(t), \rho(t)) + \mathcal{O}(\Delta t^2)$$

Algorithm

$$\mathbf{x}_0 \sim \mathbf{x}$$

$$\mathbf{L} = 0$$

While $i \geq 0$:

$$\mathbf{p} \sim Q \quad \xrightarrow{\text{Draw momentum from distribution, } Q}$$

$$\mathbf{x} = \mathbf{x}_i$$

$$H_{\text{initial}} = H(\mathbf{x}, \mathbf{p}) \quad \xrightarrow{\text{Initial energy}}$$

for $t=1$ to L do:

$$\mathbf{x}, \mathbf{p} \rightarrow \text{Leapfrog}(\mathbf{x}, \mathbf{p}, \Delta t, M) \rightarrow \text{Propagate Forwards with Leapfrog}$$

$$H_{\text{final}} = H(\mathbf{x}, \mathbf{p}) \quad \xrightarrow{\text{Final energy}}$$

$$\alpha = \exp(H_{\text{initial}} - H_{\text{final}}) \quad \xrightarrow{\text{Measure of Change in Hamiltonian (1 if equal)}}$$

Measure of 'failure' of discrete integration to conserve H .

$$U \sim N \quad \xrightarrow{\text{Draw random variable from uniform } [0, 1]}$$

$U \leq \alpha$: MH Accept/reject process:

$$\text{Accept with probability } \min(1, \alpha)$$

$$\alpha > 1 : H_{\text{initial}} > H_{\text{final}} \rightarrow 100\% \text{ acceptance}$$

$$\alpha < 1 : H_{\text{initial}} < H_{\text{final}} \rightarrow < 100\% \text{ (but close) Acceptance}$$

$$H(t) \approx H(t + \Delta t) \text{ to } \mathcal{O}(\epsilon^2)$$

Increase in H is unfavourable

HMC \leftrightarrow Blocked Gibbs + Metropolis Hastings

$$\text{With augmented distribution: } R(\mathbf{x}, \mathbf{p}) = P(\mathbf{x})Q(\mathbf{p})$$

\hookrightarrow R is Separable: $P(\mathbf{x}), Q(\mathbf{p})$ - Marginal and Conditional distribution

$$P(\mathbf{x}) = R(\mathbf{x}|\mathbf{p})$$

$$Q(\mathbf{p}) = R(\mathbf{p}|\mathbf{x})$$

Advantages:

\hookrightarrow For large L , small Δt :

Large Steps with high acceptance probability

When Sampling $(\mathbf{x}_i, \mathbf{p}_i) \stackrel{\text{iid}}{\sim} R$

\hookrightarrow Blocked Gibbs algorithm: Block 1: position variables

Block 2: Momentum variables

$$\mathbf{p}_{i+1} \sim R(\mathbf{p}|\mathbf{x}_i) \equiv Q(\mathbf{p})$$

$$\mathbf{x}_{i+1} \sim R(\mathbf{x}|\mathbf{p}_{i+1}) \equiv P(\mathbf{x})$$

\hookrightarrow MH Step: Proposal of position: Leapfrog with L steps (integrate)

\hookrightarrow Acceptance: $\min(1, \alpha)$

Theorem: HMC algorithm produces Markov Chain $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots$ that Satisfies

detailed balance with $\text{TC} = P$

Proj: Shown by proving HMC Special case of Metropolis.

i) Suppose initial position \mathbf{x}

ii) Momentum- Draw random \mathbf{p} and evolve (numerically Integrate) forwards- Leapfrog

\hookrightarrow Deterministic Integration (ie fixed output for given input)

\therefore HMC provides a proposal distribution.

$$y, q, = \mathcal{I}(\mathbf{x}, \mathbf{p}) \quad \text{dependant on Probability } (p) \text{ from dist } Q$$

$$\text{Proposal } (y|\mathbf{x}) = \frac{Q(\mathbf{p})}{|\frac{\partial y}{\partial p}|} \quad \text{Where is this from}$$

\therefore By exact reversibility of leapfrog Hamiltonian Dynamics.

$$\mathbf{x}, \mathbf{p} = \mathcal{I}(y, -q)$$

$$\text{Proposal } (\mathbf{x}|y) = \frac{Q(-q)}{|\frac{\partial \mathbf{x}}{\partial q}|}$$

By leapfrog hamiltonian dynamics: Volume preserving

$$\frac{\partial y}{\partial p} = \frac{\partial \mathbf{x}}{\partial q}$$

Given $Q(p) = Q(-p)$

$$\frac{\text{Proposal}(x|y)}{\text{Proposal}(y|x)} = \frac{Q(p)}{Q(q)}$$

Thus acceptance probability

$$\alpha(y|x) = \min\left(1, \frac{P(y) \text{Proposal}(x|y)}{P(x) \text{Proposal}(y|x)}\right) = \min\left(1, \frac{P(y) Q(q)}{P(x) Q(p)}\right) = \min\left(1, \frac{\exp(-H(y, q))}{\exp(-H(x, p))}\right)$$

$$\alpha(y|x) = \min\left(1, \exp[H(x, p) - H(y, q)]\right) = \min\left(1, \exp[H_{\text{initial}} - H_{\text{final}}]\right)$$

even for large proposal steps \rightarrow High acceptance probability

$\alpha = \min\left(1, \exp(H_{\text{initial}} - H_{\text{final}})\right)$ Still close to 1 as discretised Hamiltonian dynamics approximately conserves Hamiltonian

Hyperparameters of HMC:

time step Δt \rightarrow if too large, Leapfrog integration inaccurate (H not conserved - α very low)

number of integrations, L \rightarrow if too large, unnecessary integration - particle oscillates about potential well

\rightarrow if $L \Delta t$ too small - does not explore space quickly.

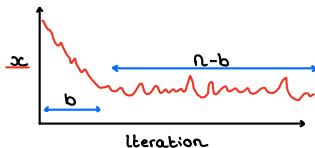
mass matrix, M (often identity matrix)

No U-Turn Sampler:

- Solves problem of tuning L by adaptively stopping Hamiltonian trajectory when it double backs on its self. $(q_{\text{new}} - q_{\text{old}}) \cdot q_{\text{new}} < 0$
- Once halted, point selected by weighted probabilities randomly.
 \rightarrow avoids unnecessary computation while maintaining exploration efficiency

Convergence Diagnostics

Burn In Period:



b - burn in: 'Initial transient part of the Markov Chains evolution - remembers starting position. Discarded as not representative of chain / target distribution'

Gelman-Rubin Statistic

Checks for consistency between intra and inter chain estimates of the variance of an arbitrary function, $f(x)$.

For N_c chains of length n (post burn-in)

Intra-Chain Variance: Variance within each individual MCMC chain.

Inter-Chain Variance: Variance between multiple independent chains.

δx_{ij} - i denotes position $i = 0, 1, \dots, n-1$

- j denotes chain $j = 0, 1, \dots, N_c-1$

GR Statistic: $R_j = \sqrt{\frac{B}{nW}}$ \rightarrow \sqrt{n} allows comparability - as B - Variance of chains of length $n \times N_c$
 W - Variance of chains of length N_c

B - Between-Chain (Inter chain) Variances (Variance of Chain means $\times n$ - Variance of all samples)

$$B = \frac{n}{N_c-1} \sum_{j=0}^{N_c-1} (\mathbb{E}[f_j] - \bar{f})^2 \rightarrow \text{Why } \times n: \text{estimating Variance underlying distribution rather than means:}$$

Corrects for $\text{Var}(\bar{x}) = \sigma^2/n$

Mean of j^{th} chain:

$$\mathbb{E}[f_j] = \frac{1}{n} \sum_{i=0}^{n-1} f(x_{ij})$$

Overall Sample Mean:

$$\bar{f} = \frac{1}{N_c} \sum_{j=0}^{N_c-1} \mathbb{E}[f_j]$$

W - Within-Chain (intra chain) variances (Mean of Chain variances)

$$W = \frac{1}{N_c} \sum_{j=0}^{N_c-1} \text{Var}[f]_j$$

Variances of j^{th} chain: $\text{Var}[f]_j = \frac{1}{n-1} \sum_{i=0}^{n-1} (f(x_{i,j}) - E[f]_j)^2$

Evaluating R_j :

R_j - depends on function f

↳ Typically evaluated on individual parameters: $f(x) = x^k$ d -dimensional

R_k - Convergence param for k^{th} param. $k = 0, 1, \dots, d-1$ Parameter Space

↳ Parameters converge/explore at different rates.

R_j ideally 1: unreasonable

$R_j \leq 1.2 \quad \forall k$ (dimensions) → expected deviation due to stochasticity.

AutoCorrelation Length:

Once a MCMC Chain have reached equilibrium - wish to draw independent Samples from it.

Adjacent Samples in MCMC Chain - not independent. (each correlated with predecessor)

Thinning - Take every m Samples

↳ Measure of Correlation:

AutoCorrelation Length - Integrated AutoCorrelation Time (IAT), T_j

↳ Length of time for Chain Samples to be approximately independent.

Expectation $f(x)$: $\bar{f} = \langle f \rangle = \int f(x) P(x) dx = \frac{1}{n} \sum_{i=0}^n f(x_i)$

If Chain were Independent: $x_i \sim \text{iid}$

$$\text{Var}[S_n] = \frac{1}{n} \sigma_f^2$$

$$\sigma_f^2 = \int dx (f(x) - \bar{f})^2 P(x)$$

↳ however Variance larger by T_j

In fact:

autocorrelation function at lag I :

$$P_I(I) = \frac{C(I)}{C(0)}$$

where $C(I) = \frac{1}{n-I} \sum_{i=0}^{n-I-1} (f(x_i) - \bar{f})(f(x_{i+I}) - \bar{f})$

$$T_j = 1 + 2 \sum_{I=0}^{n-1} P_I(I)$$

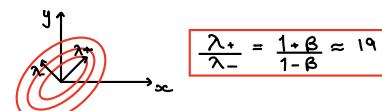
typically calculated for each parameter, T_k

$$\text{IAT: } T = \max_k (T_k)$$

For Safety often $m = 10T$

Case Study: Strongly Correlated ($\beta=0.9$) two-Dimensional Gaussian

$$P = N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \beta \\ \beta & 1 \end{pmatrix}\right) \propto \exp\left(-\frac{1}{2}(x-N)\Sigma^{-1}(x-N)^T\right) = \exp\left(-\frac{x_1^2 + y^2 - 2\beta x_1 y_1}{2[1-\beta^2]}\right)$$



Gibbs Sampling

1-D Conditional distributions:

$$P(x|y) \propto \exp\left(-\frac{(x-\beta y)^2}{2[1-\beta^2]}\right) ; \quad x|y \sim N(\beta y, 1-\beta^2)$$

$$P(y|x) \propto \exp\left(-\frac{(y-\beta x)^2}{2[1-\beta^2]}\right) ; \quad y|x \sim N(\beta x, 1-\beta^2)$$

Improvement: Define new rotated variables

$$(u, v) = R_{\pi/4}(x, y)$$

Metropolis - Hastings Sampling (Metropolis)

Proposal Distribution

$$Q = N(\Omega, \Sigma)$$

$$\text{1) } \Sigma_{\text{small}} = \begin{pmatrix} 0.05 & 0 \\ 0 & 0.05 \end{pmatrix} \quad \Sigma_{\text{large}} = \begin{pmatrix} 1.5 & 0 \\ 0 & 1.5 \end{pmatrix} \quad \Sigma_{\text{corr}} = \begin{pmatrix} 1 & \beta \\ \beta & 1 \end{pmatrix}$$

↳ Effectively using Samples from P
to produce Samples from P.

MH Sampling: Weighing up:

Too Small a proposal: High acceptance rate
Slow random walk behaviour

Too Large a proposal: Low acceptance rate
bigger Steps.

} Balance by minimising IAT, T

Hamiltonian

Requires partial derivative of the energy function:

$$\log P(x) = -E(x) + \text{const}$$

$$E(x, y) = \frac{\alpha x^2 + y^2 - 2\beta xy}{2[1-\beta^2]} \rightarrow \frac{\partial E}{\partial x} = \frac{\alpha - \beta y}{1-\beta^2} \quad , \quad \frac{\partial E}{\partial y} = \frac{y - \beta x}{1-\beta^2}$$

Free parameters: α , β , γ

↳ Optimised to allow large steps in target distribution with high acceptance rate.

Performance Metrics Table

Algorithm	IAT, T	Acceptance Freq, g	No. ud Samples, N_{eff} (% of τ)	time/iteration	time / sample
:	:	:	:	:	:
:	:	:	:	:	:
:	:	:	:	:	:

Bayesian Model Selection

↳ Comparison of models

Propose 2 models:

	Prior	Likelihood.
Model A: λ_A	$\pi(\lambda_A A)$	$L(D \lambda_A, A)$
Model B: λ_B	$\pi(\lambda_B B)$	$L(D \lambda_B, B)$

Biases in MLE Ratio:

$$\text{MLE Ratio} = \frac{\max_{\lambda_A} L(D | \lambda_A, A)}{\max_{\lambda_B} L(D | \lambda_B, B)} = \frac{L(D | \hat{\lambda}_A, A)}{L(D | \hat{\lambda}_B, B)} \rightarrow \text{MLE Ratio} \begin{cases} > 1: \text{Model A 'better'} \\ < 1: \text{Model B 'better'} \end{cases}$$

↳ Issues:

↳ Fails to account for complexity of model

A model with many parameters → generally able to fit data better → Higher value of MLE

Susceptible to prioritising overfitting through more complex model.

[Parsimonious - model that fits data well with small no. parameters]

The Odds Ratio

↳ Alternatively Bayesian approach: Posterior Odds Ratio

$$\Theta_{A,B} = \frac{1}{\Theta_{B,A}}$$

$$\text{Posterior Odds Ratio, } \Theta_{A,B} = \frac{P(A|D)}{P(B|D)} = \frac{P(D|A) \times P(A)}{P(D|B) \times P(B)}$$

Marginalise over
Parameters

Prior Odds Ratio - Encompasses prior belief of models (A: prior)
1, Unity for fair/unbiased Comparison

$$\Theta_{A,B} = \frac{\int d\lambda_A P(D|\lambda_A, A) \times P(A)}{\int d\lambda_B P(D|\lambda_B, B) \times P(B)}$$

Product rule -

Achieves Likelihood \times Prior

$$\Theta_{A,B} = \frac{\int d\lambda_A P(D|\lambda_A, A) P(\lambda_A|A) \times P(A)}{\int d\lambda_B P(D|\lambda_B, B) P(\lambda_B|B) \times P(B)} = \frac{\int d\lambda_A \mathcal{L}(D|\lambda_A, A) \pi(\lambda_A|A) \times P(A)}{\int d\lambda_B \mathcal{L}(D|\lambda_B, B) \pi(\lambda_B|B) \times P(B)}$$

Note: $\int d\lambda_A \mathcal{L}(D|\lambda_A, A) \pi(\lambda_A|A) = Z_A$
 $\int d\lambda_B \mathcal{L}(D|\lambda_B, B) \pi(\lambda_B|B) = Z_B$

$$\Theta_{A,B} = \frac{Z_A}{Z_B} \times \frac{P(A)}{P(B)}$$

Ratio of model evidences updates our
relative state of beliefs of two models

Influence of Prior Odds Ratio.

↳ Prior odds ratio generally insignificant in comparison to evidence ratio

↳ i.e. If model fits data significantly differently - evidence ratio has greater influence

If both models perform similarly (Evidence Ratio ≈ 1) or increase does not justify complexity:

Prior belief dominates.

Example:

$$\text{Model A: } P(D|A) = \prod_{i=1}^6 \frac{\exp(-\frac{1}{2}y_i^2)}{\sqrt{2\pi}} = 1.41 \times 10^{-4} \rightarrow \text{Gaussian Mean: 0, Variance: 1}$$

(no free parameters - no marginalisation integral)

$$\text{Max MLE/Mean MLE} = 1.41 \times 10^{-4} \quad \text{Evidence} = \text{Likelihood}$$

$$\text{Model B: } P(D|\lambda, B) = \prod_{i=1}^6 \frac{\exp(-\frac{1}{2}(y_i - \lambda)^2)}{\sqrt{2\pi}} \rightarrow \text{Gaussian Mean: } \lambda$$

(λ -only free parameter)

$$\text{Max MLE} (\lambda = \hat{\lambda} = 1) = 2.84 \times 10^{-3}$$

↳ Calculate Evidence:

$$\text{Prior: } P(\lambda|B) = \frac{1}{2\Lambda} \mathbb{1}_{(\lambda, \Lambda)}(\lambda) \text{ where } \Lambda \gg 1.$$

$$P(B|D) = Z_B = \int_{-\Lambda}^{\Lambda} d\lambda \frac{1}{2\Lambda} \prod_{i=1}^6 \left[\frac{\exp(-\frac{1}{2}(y_i - \lambda)^2)}{\sqrt{2\pi}} \right]$$

Limit of a wide, uninformative prior, $\Lambda \rightarrow \infty$:

$$P(B|D) \approx \frac{1}{2\Lambda} \int_{-\infty}^{\infty} d\lambda \frac{1}{\sqrt{2\pi}} \left[\frac{\exp(-\frac{1}{2}(y_i - \lambda)^2)}{\sqrt{2\pi}} \right] = \frac{1.45 \times 10^{-3}}{\Lambda}$$

Posterior Odds Ratio:

$$\Theta_{A,B} = \frac{P(A|D)}{P(B|D)} \times \frac{P(A)}{P(B)} = \frac{1.41 \times 10^{-4}}{1.45 \times 10^{-3}} \times 1 = 0.0973 \Lambda$$

↳ Effected by both prior odds ratio and model priors.

↳ Penalises both large no. parameters and weakly constrained a-priori

Computing the Bayesian evidence

Bayes Theorem: $P(\alpha|d) = \frac{P(d|\alpha, M) P(\alpha|M)}{P(d|M)}$

Bayesian Evidence: $Z = P(d|u) = \int_{\alpha} d\alpha L(d|\alpha) \pi(\alpha)$

→ Integration over full space of model parameters, $\alpha \in \alpha$

Methods for Computation:

In general - numerically evaluating a complicated/high dimensional integral over the full parameter space

→ Analytic Computation

↳ Often only feasible for well defined - low dimensional posterior. (after a conjugate prior)

Laplace Approximation.

↳ Unnormalised Posterior Distribution - $\hat{P}(\hat{\alpha}) = \hat{P}(D|\hat{\alpha}) = L(\hat{\alpha}) \pi(\hat{\alpha})$

Normalising Evidence - $Z = \hat{P}(D) = \int d\hat{\alpha} \hat{P}(\hat{\alpha})$

Suppose $\hat{P}(\hat{\alpha})$ has maximum at: $\hat{\alpha}$

Expand $\log \hat{P}(\hat{\alpha})$ about $\hat{\alpha}$ (general case: multivariate D-dimensional problem)

$$\log \hat{P}(\hat{\alpha}) \approx \log \hat{P}(\hat{\alpha}) - \frac{1}{2} C_{ij} (\alpha_i - \hat{\alpha}_i)(\alpha_j - \hat{\alpha}_j) + O([\hat{\alpha} - \hat{\alpha}]^3)$$

↳ No linear term as expanded about maximum

↳ Ignoring higher order terms.

C - Hessian Matrix

↳ Negative matrix of Second derivatives of log-posterior:

$$C_{ij} = -\frac{\partial^2}{\partial \alpha_i \partial \alpha_j} \bigg|_{\hat{\alpha} = \hat{\alpha}}$$

Exponential:

$$\hat{P}(\alpha) \approx \hat{P}(\hat{\alpha}) \exp\left(-\frac{C}{2}(\alpha - \hat{\alpha})^2\right)$$

Unnormalised Posterior \approx Gaussian
(well defined numerical integral)

$$Z \approx \hat{P}(\hat{\alpha}) \int \frac{1}{2\pi} \frac{1}{\det C} d\alpha$$

→ Semi-Analytic Method

↳ Max, $\hat{\alpha}$ and C_{ij} 's derivatives must be found.

Note: Not invariant under change of parameters

$$\alpha \rightarrow \alpha' = y(\alpha) \quad \text{then} \quad C \rightarrow C' = \frac{C}{(\frac{dy}{d\alpha})^2}$$

$$\text{Proj: (1-D case)} \quad C = -\frac{\partial^2}{\partial \alpha^2} \bigg|_{\alpha=\hat{\alpha}} \log P(\alpha) \quad C' = -\frac{\partial^2}{\partial y^2} \bigg|_{y=\hat{y}} \log P(y)$$

$$\frac{\partial^2}{\partial \alpha^2} \log P(\alpha) = \frac{\partial}{\partial \alpha} \left(\frac{\partial}{\partial y} \frac{\partial y}{\partial \alpha} \right) = \frac{\partial^2 y}{\partial \alpha^2} \cdot \frac{\partial}{\partial y} + \frac{\partial y}{\partial \alpha} \frac{\partial}{\partial \alpha} \left(\frac{\partial}{\partial y} \right) = \frac{\partial^2 y}{\partial \alpha^2} \cdot \frac{\partial}{\partial y} + \frac{\partial y}{\partial \alpha} \cdot \left(\frac{\partial}{\partial y} \frac{\partial y}{\partial \alpha} \right) \frac{\partial}{\partial y} = \frac{\partial^2 y}{\partial \alpha^2} \cdot \frac{\partial}{\partial y} + \left(\frac{\partial y}{\partial \alpha} \right)^2 \frac{\partial^2}{\partial y^2}$$

$$\frac{\partial^2}{\partial \alpha^2} \log P(\alpha) = \frac{\partial^2 y}{\partial \alpha^2} \frac{\partial}{\partial y} \log P(y(\alpha)) + \left(\frac{\partial y}{\partial \alpha} \right)^2 \frac{\partial^2}{\partial y^2} \log P(y(\alpha)) = \left(\frac{\partial y}{\partial \alpha} \right)^2 \frac{\partial^2}{\partial y^2} \log P(y(\alpha))$$

$$C = \left(\frac{\partial y}{\partial \alpha} \right)^2 C' \rightarrow C' = \frac{C}{\left(\frac{\partial y}{\partial \alpha} \right)^2}$$

Savage-Dickey Density Ratio

↳ Allows calculation of evidence ratio (Bayes factor) (not Z explicitly) for nested models:

$$\text{i.e. } B_{1,2} = \frac{P(d|M_1)}{P(d|M_2)} = \frac{Z_1}{Z_2}$$

Requires: i) 1 model to be nested inside other

ii) Consistent priors

Useful for: Bayes factor allows for model comparison - quantifies

a model's ability to explain observed data over another

$B_{1,2} > 1$: Data favours M_1

$B_{1,2} < 1$: Data favours M_2

Nested Models:

$M_1: \emptyset$ (lower dimensionality)

$M_2: \varepsilon, \emptyset$ (higher dimensionality)

$$L(d|\emptyset, M_1) = L(d|\varepsilon=0, \emptyset, M_2)$$

A Simpler model, M_1 can be recovered from a more complex model, M_2 by fixing one or more parameters.

Consistent Priors:

Priors on shared parameters of model must be consistent.

$$\text{TC}(\emptyset|M_1) = \text{TC}(\emptyset|\varepsilon=0, M_2)$$

Typically uses a Separable Prior:

$$\text{TC}(\emptyset|M_1) = g(\emptyset) \text{ and } \text{TC}(\varepsilon, \emptyset|M_2) = f(\varepsilon)g(\emptyset)$$

Evidence for Simpler model:

$$Z_{M_1} = P(d|M_1) \xrightarrow{\text{data given model over all parameter space.}}$$

$$= \int d\emptyset L(d|\emptyset, M_1) \text{TC}(\emptyset|M_1) \xrightarrow{\text{Nested models + Consistent priors}}$$

$$= \int d\emptyset L(d|\emptyset, \varepsilon=0, M_2) \text{TC}(\emptyset|\varepsilon=0, M_2)$$

$$= P(d|\varepsilon=0, M_2)$$

) Bayes Theorem

$$= \frac{P(\varepsilon=0|d, M_2)}{P(\varepsilon=0|M_2)} \frac{P(d|M_2)}{P(d|M_2)}$$

Evidence of Model 2, Z_{M_2}

$$Z_{M_1} = \frac{P(\varepsilon=0|d, M_2)}{P(\varepsilon=0|M_2)} Z_{M_2}$$

$$B_{1,2} = \frac{Z_{M_1}}{Z_{M_2}} = \frac{P(\varepsilon=0|d, M_2)}{P(\varepsilon=0|M_2)}$$

PDF of Prior

↳ Traditionally Simple analytical function - easy to evaluate.

PDF of Posterior

↳ Typically harder to evaluate

Option: MCMC to Sample posterior of model, M_2 : $(\varepsilon_i, \emptyset_i) \sim P(\varepsilon, \emptyset|d, M_2)$ for $i=1, 2, \dots, N$

Discard \emptyset_i and use remaining to estimate $P(\varepsilon|d, M_2)$ - i.e. Kernel Density estimation (KDE)

- Evaluate $P(\varepsilon=0|d, M_2)$

↳ Bayes Factor: $\frac{\text{Posterior PDF}(\varepsilon=0) \text{ of } M_2}{\text{Prior PDF}(\varepsilon=0) \text{ of } M_2}$

Thermodynamic Integration:

Uses Statistical mechanics Concepts

→ Annealing

→ Controlling the temperature of the system.

$\beta^{-1} = k_B T \rightarrow$ Annealing / Inverse temperature Parameter, β

$$0 < \beta < 1$$

Annealing Likelihood

$$L(\infty|d, \beta) = L(d|\infty)^\beta$$

$\beta=1$ (Low temp limit): Recover true likelihood

$$\rightarrow Z(\beta=1) = Z$$

$\beta=0$ (High temp limit): Flat Likelihood (Posterior \approx Prior) $\rightarrow Z(\beta=0) = 1$ (Prior is normalised)

Modified Posterior:

$$P(\infty|d, \beta) = \frac{L(d|\infty)^\beta \text{TC}(\infty)}{Z(\beta)}$$

$$Z(\beta) = \int_{\infty} d\infty L(d|\infty)^\beta \text{TC}(\infty)$$

β - a smooth switch between prior and posterior by turning likelihood on and off.

$$\frac{d}{d\beta} \log(z(\beta)) = \frac{1}{z(\beta)} \frac{dz}{d\beta}$$

$$L(d\omega) = \exp(\beta \log(d\omega))$$

$$= \frac{1}{z(\beta)} \int_{\mathcal{X}} d\omega \pi(\omega) \frac{d}{d\beta} L(d\omega)^{\beta} = \frac{1}{z(\beta)} \int_{\mathcal{X}} d\omega \pi(\omega) L(d\omega)^{\beta} \log L(d\omega)$$

$$\frac{d}{d\beta} \log(z(\beta)) = \int_{\mathcal{X}} d\omega P(\omega|d, \beta) \log L(d\omega) = E_{\omega} [\log L(d\omega) | \beta]$$

'Expectation of $\log(L(d\omega))$ over all realisations of model parameters, ω drawn from modified Bayesian posterior at a fixed value inverse temperature, β '

$$\log(z) = \int_0^1 d\beta E_{\omega} [\log L(d\omega) | \beta] = [\log z(\beta)]_0^1 = \log \left(\frac{z(\beta=1)}{z(\beta=0)} \right) = \log \left(\frac{z}{1} \right)$$

Requires MCMC to approximate expectations.

MCMC Implementation:

Temperature Ladder: $(\beta_0, \beta_1, \dots, \beta_M)$

Series of increasing inverse temp: $0 \leq \beta_N < 1$ for $N = 0, 1, \dots, M$: $\beta_0 = 0$, $\beta_M = 1$

MCMC per Ladder Step:

For each β_N ,

Run MCMC Simulation to Sample from $P(\omega|d, \beta_N)$ (annealing posterior)

N posterior Samples: $\omega_i^{(\beta_N)} \stackrel{iid}{\sim} P(\omega|d, \beta_N)$ for $i=0, 1, \dots, N-1$

Estimate Expectation: $E_{\omega} [\log L(d\omega) | \beta_N] \approx \frac{1}{N} \sum_{i=0}^{N-1} \log L(d\omega_i^{(\beta_N)})$

Evaluate Integral:
 $\log(z) = \frac{1}{2} \sum_{N=1}^M (E_{\omega} [\log L(d\omega) | \beta_N] + E_{\omega} [\log L(d\omega) | \beta_{N-1}]) \Delta \beta_N$ where $\Delta \beta_N = \beta_N - \beta_{N-1}$
 (trapezium rule)

Computationally expensive: Many temperatures - Bladder Steps (large M)
 Many Chain steps for large Sample Size.

Improvements:

Use info from previous temp to inform Starting position of next

Reduce burn in period.

Parallel-Tempered MCMC - run in parallel and exchange info during evolution.

Nested Sampling (MC (Stochastic Sampling) algorithm)

Primarily an Integration algorithm, also produces posterior Samples.

→ Requires use of other Stochastic Sampling Algorithm internally (nested).

$\Xi(L)$ - 'what fraction of the prior volume has a Likelihood $> L$:

i.e. probability Γ_V from prior gives good fit ($> L$) to data

Aim for: $Z = \int_{\mathcal{X}} d\omega L(d\omega) \pi(\omega)$ (integral over entire parameter Space, \mathcal{X}).

Maximum Likelihood, $L_{\max} = \max_{\omega \in \mathcal{X}} L(d\omega)$

Prior probability/Mass:

$$\Xi(L) = \int_{\{\omega | L(d\omega) > L\}} d\omega \pi(\omega)$$

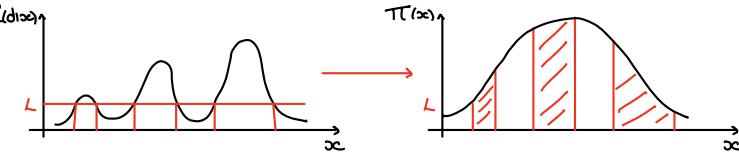
Integral of Prior over region dependent on Likelihood.

Monotonically decreasing function with:

$$\Xi(0) = 1, \quad \Xi(L_{\max}) = 0$$

'The Prior probability associated with points with a Likelihood greater than a given value L '

Region of integration does not need to be simply connected i.e multimodal.



Note: both L and L referred to as likelihood.

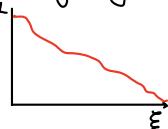
Assuming ξ is continuous (a likelihood is not completely flat):

$\xi(L)$ - Strictly decreasing

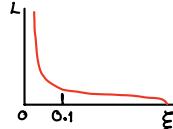
Thus invertible: $L(\xi)$ - Single values, well defined

$L(\xi)$ - likelihood for a given prior mass

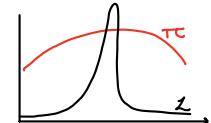
In low dimensions:



In high dimensions:



Most prior volume (parameter space) gives terrible fit to data - Likelihood - Implies very well peaked likelihood relative to Prior



Example:

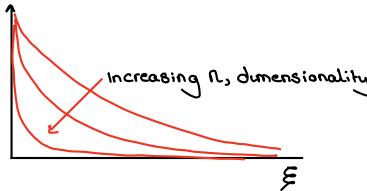
$$\xi \in \mathbb{R}^n \quad \text{Prior: } \text{ΤΤ}(\xi) = \frac{1}{V_n R^n} \begin{cases} 1 & \text{if } \|\xi\| < R \\ 0 & \text{otherwise} \end{cases} \quad \text{Likelihood: } L(d\xi) = \exp\left(-\frac{\|\xi\|^2}{2}\right)$$

$$\text{For given } L: \quad L = \exp\left(-\frac{\|\xi\|^2}{2}\right) \quad r = \sqrt{-2 \ln(L)}$$

$$\text{Area of prior: Volume} \times \text{Density} = V_n r^n \times \text{ΤΤ}(r)$$

$$\text{Fraction } \xi(L) = \left(\frac{r}{R}\right)^n = \left(\frac{\sqrt{-2 \ln(L)}}{R}\right)^n$$

$$L(\xi) = \exp\left(-\frac{R^2}{2} \xi^{2n}\right)$$



Considering Small range of likelihood: $(L, L+dL)$

$d\xi$ - The prior mass associated with likelihood values in this range

$$\begin{aligned} d\xi &= \xi(L) - \xi(L+dL) \\ &= \int_{\{\xi | L < L(d\xi) < L+dL\}} d\xi \text{ ΤΤ}(\xi) \end{aligned} \quad \text{Prior Mass for Likelihood Values } L, L+dL$$

Thus:

$$Z = \int L(\xi) \text{ ΤΤ}(\xi) d\xi \quad \text{Integral of likelihood weighted by prior. (High dimensional)}$$

$$dz = L \, d\xi \quad \text{Discrete formulation - Likelihood, } L \times \text{Prior Mass, } L \, d\xi$$

$$Z = \int_0^1 dz \, L(\xi) \quad \text{Total Evidence - Full likelihood range. (1 Dimensional)}$$

Determining Integral Numerically:

Decreasing Sequence of M points in Prior Volume: $0 < \xi_M < \xi_{M-1} < \dots < \xi_0 < 1$

With associated (threshold) likelihoods: $L_n = L(\xi_n)$

Numerical Integral
(Trapezium Rule)

$$Z \approx \frac{1}{2} \sum_{n=1}^M (L_{n-1} + L_n) \Delta \xi_n \quad \text{where } \Delta \xi_n = \xi_{n-1} - \xi_n$$

Define $\xi_0 = 1$ $L_0 = 0$.

Equivalent formulation:

$$Z \approx \sum_{n=1}^N w_n L_n$$

$$\text{Weights } w_n = (\xi_{n-1} - \xi_n)/2$$

$$\xi_{N+1} = \xi_N$$

Nested Sampling Algorithm

For $j=0$ to $N_{\text{live}}-1$:

$$\boldsymbol{x}_j \sim \pi(\boldsymbol{x})$$

$$L_j = \mathcal{L}(\boldsymbol{d}|\boldsymbol{x}_j)$$

$$Z=0, \text{Samples} \rightarrow [], i=0$$

While Condition:

$$\text{idx} = \text{argmin}(L_0, \dots, L_{N-1})$$

$$L_i = L_{\text{idx}}$$

$$\xi = \text{exp}(-i/N_{\text{live}})$$

$$w_i = \frac{1}{2} [\text{exp}(-(i-1)/N_{\text{live}}) - \text{exp}(-(i+1)/N_{\text{live}})]$$

$$Z = Z + w_i L_i$$

$$\boldsymbol{x} \sim \pi(\boldsymbol{x}) \mathcal{L}(\boldsymbol{d}|\boldsymbol{x}) > L_i$$

$$\boldsymbol{x}_{\text{disc}} = \boldsymbol{x}_i$$

$$L_{\text{disc}} = \mathcal{L}(\boldsymbol{d}|\boldsymbol{x})$$

Continue = Stopping Condition

$$i = i+1$$

Initialising Parameter Space location of N_{live} live points:

by drawing independent samples from prior.

Calculate likelihood of data for each parameter values. $\mathcal{L}(\boldsymbol{d}|\boldsymbol{x}_j)$

N_{live} - Sufficiently large (10^3 - 10^4)

For each iteration, identify live point with lowest Likelihood.

Replace it with new sample with higher likelihood

i.e. draw new/replacement live point from Constrained prior

$$\pi(\boldsymbol{x} | \mathcal{L}(\boldsymbol{d}|\boldsymbol{x}) > L_i)$$

Store discarded point - forming a list of discarded points

Relation (Probabilistic) between \mathcal{L} and ξ

Update numerical integration of evidence using trapezium weighted sum.

Relation of \mathcal{L} and ξ - $\mathcal{L}(\xi)$ or $\xi(\mathcal{L})$ - Probabilistic

For each iteration prior mass shrinks by $\xi_i = t \xi_{i-1}$

Where t random variable ($0 < t < 1$) following

$$\text{distribution: } P(t) = \begin{cases} N_{\text{live}} t^{N_{\text{live}}-1} & \text{for } 0 < t < 1 \\ 0 & \text{Otherwise} \end{cases}$$

Which has:

$$E[\log t] = -1/N_{\text{live}} \quad \text{Var}[\log t] = 1/N_{\text{live}}^2$$

This can be propagated forward to estimate error in Evidence from Nested Sampling

$$\xi_i \approx \text{exp}(-i/N_{\text{live}})$$

Justification

t - the mass fraction remains after removing lowest-likelihood sample.

For each iteration - N_{live} points sampled uniformly in remaining prior mass. $(0, 1)$

Remaining live points - drawn from remaining mass

New threshold ξ_i - is largest among N_{live} points (largest ξ value/prior mass function).

Mathematical Derivation

$$N_{\text{live}}: x_i \sim U(0,1) \rightarrow T = \max(x_i)$$

$$P(t) = \frac{d}{dt} \int_{-\infty}^t dt P(t)$$

$$\text{CDF: } P(T \leq t) = P(\max(x_i) \leq t) \text{ iid}$$

$$= P(x_1 \leq t) \cdot P(x_2 \leq t) \dots P(x_{N_{\text{live}}} \leq t)$$

$$= (P(x_i \leq t))^{N_{\text{live}}} = t^{N_{\text{live}}}$$

PDF: Derivative of CDF

$$P(t) = \frac{d}{dt} P(T \leq t) = \frac{d}{dt} t^{N_{\text{live}}}$$

$$P(t) = N_{\text{live}} t^{N_{\text{live}}-1}$$

Stopping Criterion

ALS algorithm iterates - Nine points cluster to $L_{\text{peak}} = L_{\text{max}}$

Each Step: L_i increases and ξ_i decreases to 0

More accurate Z updates

Convergence $\Delta Z \rightarrow 0$:

Remaining Evidence Contribution

$\Delta Z = L_* \xi_i - L_{\text{max}} \log \text{live points}$

When $|\Delta Z| = L_* \xi_i < \text{tolerance}$

Stopping Condition (L_*, ξ_i, tol)

If $L_* \xi_i < \text{tol}$:

Continue = False

else:

Continue = True

Sampling from Constrained Prior:

1) Rejection Sampling from the Prior:

↳ Try to Sample from unconstrained prior.

Acceptance Criterion: $L > L'$ at sampled point.

Each Step - Prior probability of Condition decreases. (acceptance fraction).

2) MCMC (ie Metropolis Hastings)

3) Rejection Sampling from bounding distribution

4) Galilean Monte Carlo

Posterior Predictive Checks (Sense Check)

Prior odds ratio gives preferred model from finite list of possible models.

But: Preferred Model from Subset $\not\rightarrow$ Correct Model

Given a model M with parameters Θ and observed data d_{obs} ,

the Bayesian Posterior can be defined: $\text{Posterior}(\Theta | d_{\text{obs}}, M)$

Can define a generative model, and draw 'predicted' / unobserved data, d_{pred} from it. Such as:

Posterior Predictive Distribution

$\Theta | d_{\text{obs}}, M \sim \text{Posterior}(\Theta | d_{\text{obs}}, M)$

$d_{\text{pred}} | \Theta, M \sim \text{Likelihood}(d_{\text{pred}} | \Theta, M)$

Overall:

$d_{\text{pred}} \sim \text{PPD}(d_{\text{pred}} | d_{\text{obs}}, \Theta, M)$

Note: Prior Predictive Distribution

$\Theta | M \sim \text{Prior}(\Theta | M)$

$d_{\text{pred}} | \Theta, M \sim \text{Likelihood}(d_{\text{pred}} | \Theta, M)$

Once generated data, d_{pred} :

Compare y_{pred} with y_{obs} to check Consistency

Commonly: Simple visual inspection (plots).

Flaws in PPC:

Uses y_{obs} twice: y_{pred} is generated using y_{obs} .

y_{pred} is compared with y_{obs} .

Limitations - Shouldn't use for:

→ Do not use to perform model comparison

Or over-engineer Statistics on them

Hierarchical Bayesian Models:

Gaussian Processes

An infinite-dimensional generalisation of finite-dimensional Gaussian distribution.

Recap: Univariate Gaussian Distribution

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

μ : mean
 σ^2 : Variance (>0)

$$x \sim N(\mu, \sigma^2)$$

Multivariate Gaussian Distribution

$$p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

μ - Mean Vector ($\in \mathbb{R}^d$)

Σ - $d \times d$ Covariance matrix

(Symmetric and +ve definite)

$$x \sim N(\mu, \Sigma)$$

Note: Positive Definite

A real, Symmetric matrix Σ is positive definite if: $Z^T \Sigma \cdot Z > 0$

(Positive Semi Definite: $Z^T \Sigma \cdot Z \geq 0$) i.e. eigenvalues of Σ $\lambda_i > 0 \forall i$

Exercise: Show that for a real, Symmetric matrix is positive definite if and only if all eigenvalues λ are strictly greater than zero, $\lambda > 0$.

For a Positive, Symmetric Matrix:

Spectral Theorem: Real, Symmetric matrix A diagonalised by Orthogonal matrix.

$$A = Q \Lambda Q^T$$

Q - Orthogonal matrix

Λ - diagonal matrix.

$$x^T A x = x^T Q \Lambda Q^T x$$

$$= (Q^T x)^T \Lambda (Q^T x)$$

$$= y^T \Lambda y = \sum_{i=1}^n \lambda_i y_i^2 \quad \text{For } x^T A x > 0 : \sum_{i=1}^n \lambda_i y_i^2 > 0$$

Gaussian Processes

defines distributions over random functions rather than finite no. random variables.
(distribution of functions).

Consider Set S , and function, j :

$j: S \rightarrow \mathbb{R}$ - Maps points in Set S to real numbers, \mathbb{R}

Sampling over functions is random (i.e. multiple functions randomly selected from).

↳ Gaussian Process: distribution over functions.

Mathematically:

For any Set S , a Gaussian process assigns a random value $j(x)$ to each point $x \in S$ such that for a finite $n \in \mathbb{N}$ and any collection of points $\{x_i\}_{i=1, \dots, n} \in S$ the vector $(j(x_1), \dots, j(x_n)) \in \mathbb{R}^n$ is distributed as a multivariate (ndim) Gaussian vector.

Advantages

- Easy interpolation in any no. dimensions/Spaces with non-trivial topology.
- Computational cost
- Memory Scales $\Theta(n^3)$

For any finite set of inputs: $\{x_1, \dots, x_n\}$

Corresponding outputs j follow mv Gauss Dist

Distribution of Functions, j , from Gaussian Process, GP:

$$j \sim GP(\mu, k)$$

μ - mean functions

k - Covariance functions

Set, S

↳ Either finite or infinite

If S finite - GP multivariate Gaussian

Typically S continuous - $S = \mathbb{R}$

Use of Gaussian Processes:

↳ GP allows us to consider infinite dimensional smooth functions while only working with finite sets of points

Infinite Dimensional Smooth Functions \longleftrightarrow Finite Dimensional Random Vectors

Defining Gaussian Processes:

For a given Set of points (x_1, x_2, \dots, x_n)

$$y \sim N(N(x), K(x, x)) \sim GP(N(x), K(x, x))$$

Note: Kernel Function \approx Covariance Function

Can be arbitrary.

N - mean vector

$$N = (N(x_1), N(x_2), \dots, N(x_n))$$

K - Covariance Matrix

(Kernel/Covariance function)

$$K = \begin{pmatrix} K(x_1, x_1) & \dots & K(x_1, x_n) \\ \vdots & \ddots & \vdots \\ K(x_n, x_1) & \dots & K(x_n, x_n) \end{pmatrix} \quad K_{ij} = k(x_i, x_j)$$

must be Symmetric positive (semi) definite function

if $n \times n$ matrix with Components $K_{ij} = k(x_i, x_j)$ is

positive definite for all sets $n \in \mathbb{N}$ points x_1, \dots, x_n in S .

i.e all matrices K from points in S : must be positive definite

Examples of GP:

Gaussian Process Regression:

Allow easy interpolation of data in any dimension and non-trivial topologies.

Allows inclusion of prior information (as Bayesian in design).

However: Computationally expensive.

Examples of Kernel Functions:

Squared Exponential Kernel

$$K_{SE}(x, x') = \sigma_f^2 \exp\left(-\frac{(x-x')^2}{2L^2}\right)$$

$\sigma_f^2 > 0$: overall amplitude parameter

$L > 0$: length scale parameter

Smooth - desirable property

Note: 1D definition ($S = \mathbb{R}$)

Generalised by replacing $(x-x')$ with Suitable distance metric

Linear Kernel

$$K_{\text{Linear}}(x, x') = \sigma_f^2 x x'$$

$\sigma_f^2 > 0$: overall amplitude parameter

Proof of Symmetric Positive Semi-definite.

Symmetric: $K_{\text{Linear}}(x, x') = K_{\text{Linear}}(x', x)$ as $x x' = x' x$.

Semi-definite: Consider Covariance matrix $K = \sigma_f^2 x x^T$

Eigenvalue eqn: $K e = \lambda e$

$$\sigma_f^2 x x^T e = \lambda e$$

↳ Singular eigen vector: x

Singular eigenvalue: $\lambda = \sigma_f^2 x x^T$ or $\sigma_f^2 \|x\|^2$

↳ All the rows are linearly dependent as

they are simply $x \sim x$, i.e. vector by a scalar.

This makes it a rank=1 matrix,

Otherwise for any eigenvector orthogonal to x , v ($x^T v = 0$)

$$K \cdot v = \sigma_f^2 x x^T v = \sigma_f^2 x \cdot 0 = 0. \quad \text{Thus eigenvalue 0.}$$

Brownian Kernel:

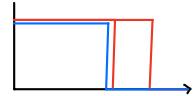
$$K_{\text{Brownian}}(x, x') = \sigma_f^2 \min(x, x')$$

Proof of Positive definite Kernel Junction:

K Brownian rewritten:

$$K_{\text{Brownian}}(\mathbf{x}, \mathbf{x}') = \sigma_j^2 \int_0^{\infty} dt H(t; \mathbf{x}) H(t; \mathbf{x}')$$

Where $H(t; \mathbf{x}) = \begin{cases} 1 & \text{if } t \leq \mathbf{x} \\ 0 & \text{if } t > \mathbf{x} \end{cases}$



Aim to Show non-negative for any set of points $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$

for any vector $\mathbf{z} \in \mathbb{R}^n$

$$\begin{aligned} \mathbf{z}^T \cdot \mathbf{K} \cdot \mathbf{z} &= \sum_{i,j} z_i z_j K_{ij} \\ &= \sum_{i,j} \sigma_j^2 z_i z_j \int_0^{\infty} dt H(t; \mathbf{x}_i) H(t; \mathbf{x}_j) \\ &= \sigma_j^2 \int_0^{\infty} dt \underbrace{\sum_{i,j} z_i z_j H(t; \mathbf{x}_i) H(t; \mathbf{x}_j)}_{\left(\sum_i z_i H(t; \mathbf{x}_i) \sum_j z_j H(t; \mathbf{x}_j) \right)} = \left(\sum_i z_i H(t; \mathbf{x}_i) \right)^2 \\ &= \sigma_j^2 \left(\sum_i z_i H(t; \mathbf{x}_i) \right)^2 \geq 0. \end{aligned}$$

Covariance Functions:

Given an example 'Seed Covariance function', it can be adjusted

accordingly:

Summation:

For any two Covariance functions: $K_1(\mathbf{x}, \mathbf{x}')$ and $K_2(\mathbf{x}, \mathbf{x}')$

Their Sum: $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}')$ is also a Covariance function

Product:

For any two Covariance functions $K_1(\mathbf{x}, \mathbf{x}')$ and $K_2(\mathbf{x}, \mathbf{x}')$

Their Product: $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') K_2(\mathbf{x}, \mathbf{x}')$ is also a Covariance function.

Warping:

For any Covariance function $K(u, u')$ on the Space S' and a warping function

$u: S \rightarrow S'$: the function $K(\mathbf{x}, \mathbf{x}') = K(u(\mathbf{x}), u(\mathbf{x}'))$ is a new Covariance function on the Space S . (u pulls back Kernel S' to S). Where $\mathbf{x}, \mathbf{x}' \in S$, $u, u' \in S'$

Renormalisation:

For any Covariance function $K(\mathbf{x}, \mathbf{x}')$ and function $\alpha: S \rightarrow \mathbb{R}$ the

function $\alpha(\mathbf{x}) \alpha(\mathbf{x}') K(\mathbf{x}, \mathbf{x}')$ is a new Covariance function

Proof of warping junction:

Show $K(\mathbf{x}, \mathbf{x}')$ is positive definite: $\mathbf{z}^T \cdot \mathbf{K} \cdot \mathbf{z} = \sum_{i,j} z_i z_j K_{ij}$

$$\begin{aligned} &= \sum_{i,j} z_i z_j K(u(\mathbf{x}_i), u(\mathbf{x}_j)) \quad \text{for any choice } \mathbf{x}_i \\ &= \sum_{i,j} z_i z_j K(u_i, u_j) \\ &= \sum_{i,j} z_i z_j K(u_i, u_j) \quad \text{where } u_i = u(\mathbf{x}_i) \\ &\geq 0 \text{ as } K \text{ is positive semi definite.} \end{aligned}$$

Example with Periodic Kernel:

Let $S' = \mathbb{R}^2$ using 2D Square Exp Kernel $K((\mathbf{x}), (\mathbf{x}')) = \exp\left(\frac{-(\mathbf{x}-\mathbf{x}')^2 - (y-y')^2}{2L^2}\right)$

Trying to induce/pull back Gaussian Process onto Unit Circle/New Space:

$S = S_1$ (circle)

$$U(\emptyset) = \begin{pmatrix} \cos(\emptyset) \\ \sin(\emptyset) \end{pmatrix}$$

Define a new Kernel K on S , by

$$K(\emptyset, \emptyset') = K(U(\emptyset), U(\emptyset'))$$

$$= \exp\left(-\frac{1}{2l^2} \left[(\cos(\emptyset) - \cos(\emptyset'))^2 + (\sin(\emptyset) - \sin(\emptyset'))^2 \right]\right)$$

$$= \exp\left(-\frac{2}{l^2} \sin^2\left(\frac{\emptyset - \emptyset'}{2}\right)\right)$$

Periodic Kernel

proved to be Semi definite assuming

Squared exponential

Focusing on the Real Line: $S = \mathbb{R}^n$

Stationary Kernel

Only depends on the difference of the inputs:

$$K(x, x') = K(\tau) \quad \text{where } \tau = x - x'$$

A Stationary Kernel is invariant under translations

$$K(x, x') = k(x + \Delta, x' + \Delta) \quad \text{for all } \Delta \in \mathbb{R}^n.$$

Squared Exponential Kernel:

↳ Both Stationary and isotropic.

Isotropic Kernel

A Stationary distribution that only depends on Euclidean Distance

$$\text{i.e. } k(\underline{r}) = k(r) \quad r = \|\underline{r}\|$$

Invariant under both translations and rotations

Properties of Stationary Kernels:

i) By Symmetry Property: $K(x, x') = K(x', x)$

$$\rightarrow K(-\underline{r}) = K(\underline{r})$$

ii) Positivity: $K(\tau = 0) = K(x, x) \geq 0$

iii) boundedness: $K(\tau) \leq K(\tau = 0)$ for all $\tau \in \mathbb{R}^d$

$$K(\tau) = \begin{pmatrix} K(0) & K(\tau) \\ K(\tau) & K(0) \end{pmatrix} \rightarrow \text{Eigenvalues to be positive:}$$

must be: $K(\tau) \leq K(0)$

Proving Stationary Kernels are positive (semi) definite:

Lemma: A function $K: \mathbb{R}^n \rightarrow \mathbb{R}$, denoted $K(\tau)$ can be used to define a Semi-definite Stationary Kernel: $k: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ as $k(x, x') = K(x - x')$ if the Fourier transform $\tilde{K}(y) > 0$ for all $y \in \mathbb{R}^n$

Fourier transform:

$$\tilde{K}(y) = \int K(\tau) e^{(2\pi i y \cdot \tau)}$$

Proof: Suppose the Fourier transform always exists:

$$K(\tau) = \int_{\mathbb{R}^n} dy \tilde{K}(y) \exp(-2\pi i y \cdot \tau)$$

Now Considering:

$$\begin{aligned} \underline{z}^T \cdot K \cdot \underline{z} &= \sum z_i z_j K(x_i, x_j) \\ &= \sum z_i z_j K(x_i - x_j) \quad \text{Stationary Kernel} \\ &= \sum_{i,j} z_i z_j \int dy \tilde{K}(y) \exp(-2\pi i y \cdot (x_i - x_j)) \end{aligned}$$

$$= \int dy \sum_{i,j} z_i z_j \tilde{K}(y) \exp(-2\pi i y \cdot x_i) \exp(2\pi i y \cdot x_j)$$

Complex Conjugate

→ Generalised Result of Bochner's Theorem

A continuous function $K(\tau)$ is positive semi-definite if and only if it is the Fourier transform of a finite non-negative Borel measure N on \mathbb{R}^d .

$$\text{i.e. } K(\tau) = \int_{\mathbb{R}^n} dy N(y) \exp(-2\pi i y \cdot \tau)$$

$$= \int df \tilde{K}(f) \left| \sum z_i \exp(-2\pi i f \cdot x_i) \right|^2 \geq 0 \quad \boxed{\tilde{K}(f) \geq 0}$$

Thus $\tilde{K}(f) \geq 0$ for all f : The Stationary Kernel is positive semi definite.

The Squared Exponential Kernel:

We first show it is a Stationary Kernel, $K_{SE} = \sigma_f^2 \exp\left(-\frac{(x-x')^2}{2l^2}\right) = \sigma_f^2 \exp\left(-\frac{\tau^2}{2l^2}\right)$
 i.e. $K(x, x') = K(x + \Delta, x' + \Delta)$

$$\begin{aligned}
 \text{Taking its Fourier transform: } \tilde{K}(f) &= \int_{-\infty}^{\infty} \sigma_f^2 \exp\left(-\frac{\tau^2}{2l^2}\right) \exp(2\pi i f \tau) d\tau \\
 &= \sigma_f^2 \int_{-\infty}^{\infty} \exp\left(-\frac{\tau^2}{2l^2} + 2\pi i f \tau\right) d\tau \xrightarrow[\text{Completing the Square}]{-\frac{\tau^2}{2l^2} + 2\pi i f \tau = -\frac{1}{2l^2} [\tau^2 + 4\pi i f l^2 \tau]} \\
 &\quad -\frac{1}{2l^2} \left[\tau^2 + 4\pi i f l^2 \tau \right] \\
 &= -\frac{1}{2l^2} \left[\tau + 2\pi i f l^2 \right]^2 - 2\pi f l^2 \\
 &= \sigma_f^2 \exp\left(2\pi^2 f^2 l^2\right) \int_{-\infty}^{\infty} \exp\left\{-\frac{1}{2l^2} \left[\tau + 2\pi i f l^2 \right]^2\right\} d\tau \\
 &= \sigma_f^2 \exp\left(2\pi^2 f^2 l^2\right) \int_{-\infty}^{\infty} \exp\left\{-\frac{u^2}{2l^2}\right\} du = \sigma_f^2 \exp\left(2\pi^2 f^2 l^2\right) \times \sqrt{2\pi} \\
 &\quad \xrightarrow[\text{L} \sqrt{2\pi} : \text{Gaussian Integral}]{\int_{-\infty}^{\infty} \exp\left\{-\frac{u^2}{2\sigma^2}\right\} du = \sigma \sqrt{2\pi}} \\
 \tilde{K}(f) &= \sigma_f^2 \sqrt{2\pi} \exp\left(-2\pi^2 f^2 l^2\right) \\
 &\geq 0 \text{ for all } f
 \end{aligned}$$

Gaussian Process Regression:

(Interpolation + Extrapolation)

Lemma: Let $x \sim N(\mu, \Sigma)$. Split the vector $x \in \mathbb{R}^n$ into $x^T = (x_1^T, x_2^T)$

where $x_1 \in \mathbb{R}^{n_1}$ and $x_2 \in \mathbb{R}^{n_2}$ and $n_1 + n_2 = n$.

Similarly Split $N^T = (N_1^T, N_2^T)$ and $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$ where $\Sigma_{11} = n_1 \times n_1$, $\Sigma_{12} = n_1 \times n_2$, $\Sigma_{22} = n_2 \times n_2$.

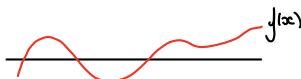
The distribution of x_1 conditioned on a fixed value of x_2 is:

$$x_1 | x_2 \sim N(N_1 + \Sigma_{12} \cdot \Sigma_{22}^{-1} \cdot (x_2 - N_2), \Sigma_1 - \Sigma_{12} \cdot \Sigma_{22}^{-1} \cdot \Sigma_{12}^T)$$

Distribution of Subset₁ with respect to fixed Subset₂

In Words - If you take a slice through an n -dimensional Multivariate Gaussian along any m -dimensional hyperplane ($m < n$) then the pdf along this slice has the shape of another Gaussian: $x_1 | x_2$

Tackling Regression:



Suppose we have measured the value of the function at n points, $\underline{x} = (x_1, x_2, \dots, x_n)$ and the measured values are $\underline{y} = (y(x_1), \dots, y(x_n))$ and we wish to predict the value $y(x_*)$ at some new location x_* .

Prior on Space of Functions: Use Zero-mean Gaussian Process (assume you have scaled to zero-mean)

$y \sim GP(0, K)$ Any General Kernel.

Def of Gaussian Process:

apply to (x_*, x_1, \dots, x_n)

Prior means:

$$(y(x_*), y_1, \dots, y_n) \sim N(0, K)$$

$$\text{Where } K' = \begin{pmatrix} K_{**} & K_*^T \\ K_* & K \end{pmatrix} \quad \text{Where } \begin{aligned} K_{**} &= k(x_*, x_*) && \text{Single value, } \mathbb{R} \\ k_* &= k(x_*, x_i) && i=1, \dots, n \end{aligned}$$

$$K = K(x_i, x_j) \quad i, j = 1, \dots, n \quad \text{Matrix, } \mathbb{R}^{n \times n}$$

$$= (k(x_*, x_1), k(x_*, x_2), \dots, k(x_*, x_n)) \\ = \begin{pmatrix} k(x_*, x_1) & \dots & k(x_*, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix}$$

$$y(x_*), y \sim N(0, \begin{pmatrix} K_{**} & K_*^T \\ K_* & K \end{pmatrix})$$

Determining the Conditional distribution:

$$y(x_*) | y = N(E[y(x_*)], \text{Var}[y(x_*)])$$

$$E[y(x_*)] = K_*^T \cdot K^{-1} \cdot y$$

$$\text{Var}[y(x_*)] = K_{**} - K_*^T \cdot K^{-1} \cdot K_*$$

View Point: 'Gaussian Processes are a type of Conjugate Prior'

Prior: $y \sim GP(0, K)$

Posterior: $y(x_*) | y \sim N(E[y(x_*)], \text{Var}[y(x_*)])$

This Conjugate Prior as Posterior distributed in the Same Way: Gaussian Process.

Example Conjugate prior with update of parameters

$$N \rightarrow N': E[y(x_*)] = K_*^T \cdot K^{-1} \cdot y$$

$$K \rightarrow K': \text{Var}[y(x_*)] = K_{**} - K_*^T \cdot K^{-1} \cdot K_*$$

Considering Measurement Values, $y(x_i)$, with associated Gaussian errors, σ_i :

Adding a diagonal term to Covariance matrix.

$$\sigma^2 = (\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)^T$$

$$E[y(x_*)] = K_*^T \cdot (K + \text{diag}(\sigma^2))^{-1} \cdot y$$

$$\text{Var}[y(x_*)] = K_{**} - K_*^T \cdot (K + \text{diag}(\sigma^2))^{-1} \cdot K_*$$

Additionally done to help Numerical Stability.

Evidence of Gaussian Process.

$$\log(L) = \log P(y) = -\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{N}{2} \log(2\pi)$$

Can use this to help inform Covariance function.

Gaussian Process Review: ('mean square': Continuity and Differentiability).

Properties of Gaussian Process are determined by k :

'terms like continuity etc become hard to define traditionally as we are discussing finite sampled points'

Restrict $\mathbf{x} \in \mathbb{R}^n$, Let x_1, x_2, \dots approach/Converge x_* if $|x_i - x_*| \rightarrow 0$ as $i \rightarrow \infty$

Mean Square Continuous

f is mean square continuous if $E[|f(x_n) - f(x_*)|^2] \rightarrow 0$ as $n \rightarrow \infty$

If this holds for all $x_* \in \mathbb{R}^n$ then we say f is MS continuous everywhere

Relating this condition to Kernel function:

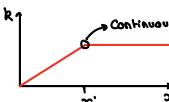
f is Mean Square Continuous at x_* iff Kernel function $k(x, x_*)$ is 'continuous' in $x = x' = x_*$ (and $\mu(x)$ is continuous at $x = x_*$).

For a Continuous Kernel: $k(\tau) = k(x, x')$ where $\tau = x - x'$

Simply requires ensuring $K(\tau)$ is continuous (translation invariance). Can show it once.

$$R_{SE} = \sigma_f^2 \exp\left(-\frac{\tau^2}{2L^2}\right) \quad \text{Continuous } \tau=0$$

$$k_{\text{Brownian}}(x, x') = \sigma_f^2 \min(x, x') \quad \text{also continuous.}$$



Mean Square Differentiability

$$\frac{\partial f}{\partial x_i} = \text{limit in MS} \left(\frac{f(x + \varepsilon e_i) - f(x)}{\varepsilon} \right)$$

$$\text{If } f \sim \text{GP} \text{ then if it exists } \frac{\partial f}{\partial x_i} \sim \text{GP} \quad \text{with } K(x, x') = \frac{\partial^2 K(x, x')}{\partial x_i \partial x_i'}$$

The level of mean square differentiable of f is determined by level of differentiability of kernel k (with factor of k)

R_{brownian} - MS continuous but not differentiable.

Often k_{SE} - Over Smooths:

Alternative: Matern Kernel Function

$$k(x, x') = k_L(x) = \frac{\sigma_f^2 2^{1-N}}{\Gamma(N)} \left(\sqrt{2N} \frac{x}{L} \right)^N K_N \left(\sqrt{2N} \frac{x}{L} \right) \quad \text{modified Bessel Function}$$

Considered Generalisation k_{SE}

$L \rightarrow \infty$: Matern \rightarrow SE

L controls 'level of smoothness'.

Derivatives at $\tau=0$: This kernel is α times MS differentiable if $N > \alpha$. as kernel is int \downarrow (α) time differentiable

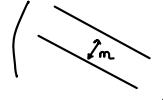
Computational Cost of GPR:

K^{-1} has dimensions $n \times n$: $\Theta(n^2)$ memory
 $\Theta(n^3)$ computation

Compact Support Kernel:

Kernels such as RSE require dense matrix ie all positions have values (but many very small).

Wendland Polynomials:

Leaves Band Diagonal K  memory $\Theta(n \cdot m)$

Hierarchical Bayesian Models (and Probabilistic Graphical Models)

Parameters Split across Several 'levels', reflected in a particular factorisation of the prior.

Example Problem 1:

Estimating the age of a Single Star: τ Using observations: t ,

Regular Bayesian Inference:

→ Gaussian Prior: $\frac{1}{\sqrt{2\pi\Delta^2}} \exp\left(-\frac{1}{2} \frac{(N-\tau)^2}{\Delta^2}\right)$

Gaussian Likelihood: $\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(t-\tau)^2}{\sigma^2}\right)$

Alternatively $\tau | N, \Delta \sim N(N, \Delta^2)$

$t | \tau, \sigma \sim N(\tau, \sigma^2)$

Prob Graphical Model

→ A way of illustrating joint distribution on every random variables. (Directed Acyclic Graph - DAG)

Fixed values: N, Δ, σ

↳ No circles

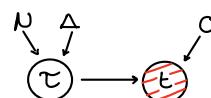
Random Variables: t, τ

↳ Shaded: observed variables

Empty: parameters/latent variables

Arrow: Denote order of conditional probabilities.

If generating fake data (what order?)



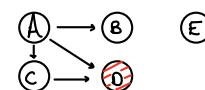
Bayesian Posterior:

$$P(\tau | t) = \frac{1}{P(t)} P(t | \tau) P(\tau)$$

$$= \frac{1}{Z} \pi(\tau) P(t | \tau)$$

$$\text{Where Evidence } Z(\sigma, N, \Delta) = \int d\tau \pi(\tau) P(t | \tau, \sigma, N, \Delta)$$

Example Problem 2:



$$P(A, B, C, D, E) = P(A) P(B | A) P(C | A) P(D | A, C) P(E)$$

Can also write: (for observed D)

$$P(A, B, C, D, E | D) = P(A, B, C, D, E | D) P(D)$$

$$P(A, B, C, D, E | D) = \underbrace{\frac{1}{P(D)} P(D | A, C)}_{\text{Posterior.}} \underbrace{P(A) P(B | A) P(C | A)}_{\text{Evidence Likelihood}} \underbrace{P(E)}_{\text{Prior}}$$

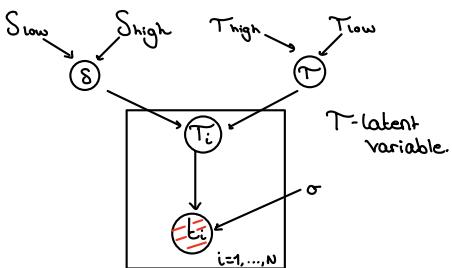
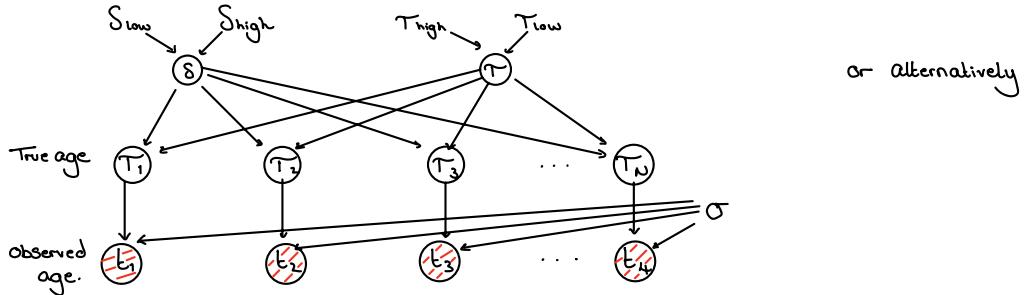
Example Hierarchical Model:

Assume all stars in a cluster formed around the same time.

Hyper Prior: $\begin{cases} \tau | \tau_{\text{low}}, \tau_{\text{high}} \sim \text{log uniform}(\tau_{\text{low}}, \tau_{\text{high}}) \\ \sigma | \sigma_{\text{low}}, \sigma_{\text{high}} \sim \text{log uniform}(\sigma_{\text{low}}, \sigma_{\text{high}}) \end{cases}$

Prior $\tau_i | \tau, \Delta \sim N(\tau, \sigma^2)$ S-Small scatter in ages within cluster

Likelihood: $t_i | \tau_i, \sigma \sim N(\tau_i, \sigma^2)$ independently for $i=1, \dots, N$



Bayes Theorem:

$$\begin{aligned}
 \text{Eq. 1: } P(\tau, \delta, \{\tau_i\}, \{t_i\}) &= P(\tau) P(\delta) P(\{\tau_i\} | \delta, \tau) P(\{t_i\} | \{\tau_i\}) \\
 &= \underbrace{\pi(\tau) \pi(\delta)}_{\text{Hyper Priors}} \underbrace{\frac{N_s}{\tau} \pi(\tau_i | \delta, \tau)}_{\text{Prior}} \underbrace{\frac{N_s}{\tau} \mathcal{L}(t_i; \tau_i)}_{\text{Likelihood.}}
 \end{aligned}$$

Eq. 2:

$$P(\tau, \delta, \{\tau_i\}, \{t_i\}) = P(\{t_i\}) P(\tau, \delta, \{\tau_i\} | \{t_i\})$$

Equating with posterior:

$$\underbrace{P(\tau, \delta, \{\tau_i\} | \{t_i\})}_{\text{Posterior}} = \frac{1}{\underbrace{P(\{t_i\})}_{\text{Evidence}}} \underbrace{\pi(\tau) \pi(\delta)}_{\text{Hyper Priors}} \underbrace{\frac{N_s}{\tau} \pi(\tau_i | \delta, \tau)}_{\text{Prior}} \underbrace{\frac{N_s}{\tau} \mathcal{L}(t_i; \tau_i)}_{\text{Likelihood.}}$$

Truely we are interested in \$\tau, \delta\$ (not individual \$\tau_i\$)

Marginalise over latent variables - \$\tau_i\$

$$P(\tau, \delta | \{t_i\}) = \frac{1}{P(\{t_i\})} \underbrace{\pi(\tau) \pi(\delta)}_{\text{Hyper Priors}} \underbrace{\frac{N_s}{\tau} \int \delta \tau' \pi(\tau' | \delta, \tau) \cdot \mathcal{L}(t_i | \tau')}_{\text{Product of evidence of Single models.}}$$