

Bayesian Theory

Study Note

Ian S.W. Ma

Winter 2020

Contents

1	Statistics Basics	1
2	Bayes Throrem	2
3	Conjugate Priors	3
4	Normal Distribution Priors	3
5	Jeffrey's Prior	4
6	Reference Prior	5
7	Point Estimates	5
8	Interval Estimates	6
9	Classic Hypothesis Testing	6
10	Bayesian Hypothesis Testing	6
11	Deterministic Numerical Integration	8
12	Monte Carlo Integration	9
13	Markov Chain Monte Carlo (MCMC)	11
14	Metropolis-Hasting Algorithm	13
15	MCMC Diagnostics	14

1 Statistics Basics

General Structure: Probability Space/Triple

Probability Space/Triple $(\Omega, \mathcal{F}, \mathbb{P})$

- Ω : **Sample Space**, set of all possible outcomes.
- \mathcal{F} : **Set of events** σ , each event is a subset of Ω
- \mathbb{P} Probability of event $\sigma \in \mathcal{F}$

Conditional Probability

Conditional probability of A given B : $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \Leftrightarrow \mathbb{P}(A|B)\mathbb{P}(B) = \mathbb{P}(A \cap B)$

General properties

- $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$
Events are mutually exclusive $\Rightarrow \mathbb{P}(A \cap B) = 0$
- $\mathbb{P}(A \cap B) = \mathbb{P}(A|B)\mathbb{P}(B)$
Events are independent $\Rightarrow \mathbb{P}(A|B) = \mathbb{P}(A)$
- Law of Total Probability

If events A_1, A_2, A_N are mutually exclusive and exhaustive ($\bigcup_{i=1}^N A_i = \Omega$) then:

$$\mathbb{P}(B) = \sum_{i=1}^N \mathbb{P}(A_i \cap B) = \sum_{i=1}^N \mathbb{P}(B|A_i)\mathbb{P}(A_i)$$

Random Variables

$X \sim$ Distribution	$\mathbb{E}[X]$	$\text{Var}[X]$
Normal(μ, σ^2)	μ	σ^2
Bernoulli(θ)	θ	$\theta(1 - \theta)$
Binomial(θ)	$n\theta$	$n\theta(1 - \theta)$
Poisson(μ)	μ	μ
Uniform(α, β)	$\frac{\alpha+\beta}{2}$	$\frac{(\alpha+\beta)^2}{12}$
Beta(α, β)	$\frac{\alpha}{\alpha+\beta}$	$\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$
Gamma(α, β)	$\frac{\alpha}{\beta}$	$\frac{\alpha}{\beta^2}$
Inv. Gamma(α, β)	$\frac{\beta}{\alpha-1}$	$\frac{\beta^2}{(\alpha-1)^2(\alpha-2)}$
Exponential(λ)	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$

Random Vectors

$\mathbf{X} = X_1, \dots, X_k \sim$ Distribution	$\mathbb{E}[X_i]$	$\text{Var}[X_i]$	$\text{Cov}(X_i, X_j)$
Multivariate Normal($\boldsymbol{\mu}, \Sigma$)	μ_i	$\Sigma_{i,i}$	$\Sigma_{i,j}$
Multinomial($n, \theta_1, \dots, \theta_k$)	$n\theta_i$	$n(1 - \theta_i)$	$-n\theta_i\theta_j$

2 Bayes Throrem

Discrete Case

$$\begin{aligned}
 \mathbb{P}(X = x_i | Y = y_j) &= \frac{\mathbb{P}(X = x_i, Y = y_j)}{\mathbb{P}(Y = y_j)} \\
 &= \frac{\mathbb{P}(Y = y_j | X = x_i) \mathbb{P}(X = x_i)}{\mathbb{P}(Y = y_j)} \\
 &= \frac{\text{Likelihood} \times \text{Prior}}{\text{Marginal}} \\
 &= \frac{\mathbb{P}(Y = y_j | X = x_i) \mathbb{P}(X = x_i)}{\sum_{k=1}^n \mathbb{P}(X = x_k, Y = y_j)} \\
 &= \frac{\mathbb{P}(Y = y_j | X = x_i) \mathbb{P}(X = x_i)}{\sum_{k=1}^n \mathbb{P}(Y = y_j | X = x_k) \mathbb{P}(X = x_k)} \\
 \Rightarrow \mathbb{P}((X|Y)) &= \frac{\mathbb{P}(Y|X) \mathbb{P}(X)}{\sum_x \mathbb{P}(X = x, Y)} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Marginal}}
 \end{aligned}$$

Continous Case

$$f(X|Y) = \frac{f(X, Y)}{g(Y)} = \frac{g(Y|X)f(X)}{\int g(Y|X)f(X)dX}$$

Bayesian Statistical Inference

- $L(\theta|y) = f(y|\theta)$: Likelihood
- $\pi(\theta)$: Prior
- $m(y) = p(y)$: Marginal distribution/Normalizing constant

$$\text{Posterior Distribution: } p(\theta|y) = \frac{f(y|\theta)\pi(\theta)}{\int p(\theta, y)d\theta} = \frac{f(y|\theta)\pi(\theta)}{\int f(y|\theta)\pi(\theta)d\theta} \propto f(y|\theta)\pi(\theta)$$

Predictive Distrubutions

$$\text{Prior Predictive Distribution: } p(y^{new}) = \int_{\theta} f(y^{new}|\theta)\pi(\theta)d\theta$$

$$\text{Posterior Predictive Distribution: } p(y^{new}|y^{old}) = \int_{\theta} f(y^{new}|\theta, y^{old})p(\theta|y^{old})d\theta$$

Bayes Theorem for Multiple Parameters $\Theta = \{\theta_1, \dots, \theta_q\}$

$$\mathbb{P}(\Theta|\mathbf{y}) = \frac{\mathbb{P}(\mathbf{y}|\Theta)\mathbb{P}(\Theta)}{\mathbb{P}(\mathbf{y})} \propto \mathbb{P}(\mathbf{y}|\Theta)\mathbb{P}(\Theta)$$

3 Conjugate Priors

Sample Distribution	Parameter	Prior
Bernoulli(θ)	θ	Beta(α, β)
Binomial(n, θ)	θ	Beta(α, β)
Poisson(θ)	θ	Gamma(α, β)
Exponential(θ)	θ	Gamma(α, β)
Normal(μ, σ_0^2)	μ	Normal(μ_h, σ_h^2)
Normal(μ_0, σ^2)	σ^2	Inverse Gamma(α, β)
Normal($\mu_0, 1/\tau$)	$\tau = 1/\sigma^2$	Gamma(α, β)
Multinomial($n, \theta_1, \dots, \theta_k$)	$\theta_1, \dots, \theta_k$	Dirichlet($\alpha_1, \dots, \alpha_k$)
Uniform($0, \theta$)	θ	Pareto(θ_m, α)

4 Normal Distribution Priors

Unknown μ , known σ^2

- Sample $\mathbf{y} \in \mathbb{R}^n$
- Prior $\mu \sim \text{Normal}(\mu_0, \sigma_0^2) = \text{Normal}\left(\mu_0, \frac{\sigma^2}{m}\right) \Rightarrow m = \sigma^2/\sigma_0^2$
- Posterior:

$$\mu|\mathbf{y}, \sigma^2 \sim \text{Normal}\left(\frac{m\mu_0 + n\bar{y}}{m+n}, \frac{\sigma^2}{m+n}\right) = \text{Normal}\left((1-w)\mu_0 + w\bar{y}, \frac{\sigma^2}{m+n}\right)$$

$$w = \frac{m}{m+n}$$

Unknown μ , known τ

- Sample $\mathbf{y} \in \mathbb{R}^n$
- Prior $\mu \sim \text{Normal}\left(\mu_0, \frac{\sigma^2}{m}\right) = \text{Normal}\left(\mu_0, \frac{1}{\tau m}\right)$
- Posterior:

$$\mu|\mathbf{y}, \sigma^2 \sim \text{Normal}\left(\frac{m\mu_0 + n\bar{y}}{m+n}, \frac{1}{\tau(m+n)}\right) = \text{Normal}\left((1-w)\mu_0 + w\bar{y}, \frac{1}{\tau(m+n)}\right)$$

$$w = \frac{m}{m+n}$$

Unknown τ , known μ

- Sample $\mathbf{y} \in \mathbb{R}^n$
- Prior $\tau \sim \text{Gamma}(\alpha, \beta)$
- Posterior:

$$\tau|\mathbf{y}, \mu \sim \text{Gamma}\left(\alpha + \frac{n}{2}, \beta + \frac{z^2}{2}\right)$$

$$z^2 = \sum_i (y_i - \mu)^2$$

Unknown σ^2 , known μ

- Sample $\mathbf{y} \in \mathbb{R}^n$
- Prior $\mu, \sigma^2 \sim \text{Normal}\left(\mu_0, \frac{\sigma^2}{\kappa}\right) \times \text{Inverse Gamma}(\alpha, \beta)$
- Posterior:

$$\mu, \sigma^2 | \mathbf{y} \sim \text{Normal}\left(\frac{\kappa\mu_0 + n\bar{y}}{\kappa + n}, \frac{\sigma^2}{\kappa + n}\right) \times \text{Inverse Gamma}\left(\alpha + \frac{n}{2}, \beta + \frac{(n-1)s^2}{2} + \frac{\kappa n(\bar{y} - \mu_0)^2}{2(\kappa + n)}\right)$$

$$s^2 = \sum_i (y_i - \bar{y})^2 / (n-1)$$

$\kappa = \sigma^2(\alpha - 1)/\beta$ or given (seriously it's not written anywhere Ken, I'm just guessing here man)

5 Jeffrey's Prior

Definition

Jeffery Prior: $\pi_{JP}(\theta) \propto \sqrt{I(\theta|y)}$

Fisher Information:

$$I(\theta|y) = \mathbb{E} \left[\left(\frac{d \log f(y|\theta)}{d\theta} \right)^2 \right] = -\mathbb{E} \left[\frac{d^2 \log f(y|\theta)}{d\theta^2} \right]$$

1:1 transformation of Jeffrey's Prior

$$\pi(\phi) = \pi_{JP}(\theta) \left| \frac{d\theta}{d\phi} \right|$$

Jeffrey's Prior for Multivariate Parameter Vector

Gradient of log likelihood:

$$S(\boldsymbol{\theta}) = \begin{bmatrix} \partial_{\theta_1} \log f(x) \\ \vdots \\ \partial_{\theta_n} \log f(x) \end{bmatrix}$$

Hessian Matrix of log likelihood:

$$H(\boldsymbol{\theta}) = \text{Jacobian}[S(\boldsymbol{\theta})] = \begin{bmatrix} \partial_{\theta_1} \partial_{\theta_1} \ln f(x) & \cdots & \partial_{\theta_1} \partial_{\theta_n} \ln f(x) \\ \vdots & \ddots & \vdots \\ \partial_{\theta_n} \partial_{\theta_1} \ln f(x) & \cdots & \partial_{\theta_n} \partial_{\theta_n} \ln f(x) \end{bmatrix}$$

Fisher Information: $I(\theta|x) = -\mathbb{E}[H(\theta)]$

Jeffrey's Prior: $\pi_{JP} \propto \sqrt{\det(I(\theta|x))}$

6 Reference Prior

Kullback-Leibler (KL) divergence

A measure of the difference between two pmfs/pdfs. When $KL(f, g) = 0$, the two distributions are identical.

Properties

- $KL(f, g) \neq KL(g, f)$
- $KL(fmg) \geq 0$

Continuous parameter x

$$KL(f, g) = \int \ln \left[\frac{f(x)}{g(x)} \right] dx = \mathbb{E}_f[\log f(X)] - \mathbb{E}_g[\log f(X)]$$

Discrete parameter x

$$KL(f, g) = \sum_{x \in X} \ln \left[\frac{f(x)}{g(x)} \right] = \mathbb{E}_f[\log f(X)] - \mathbb{E}_g[\log f(X)]$$

Reference Prior

Read week 4 notes (2.2)

7 Point Estimates

Components of Decision Theory with emphasis on parameter estimation

- State Space Θ , unknown true value $\theta \in \Theta$
- Action Space $\mathcal{A} \ni a$, sometimes $\mathcal{A} = \Theta$
- Sampling Distribution $f(\mathbf{y}|\theta)$
- Loss Function $\mathcal{L}(a|\theta)$
- Risk $R_\theta(a|\mathbf{y}) = \mathbb{E}_\theta[L(a|\theta, \mathbf{y})] = \int_{\theta \in \Theta} \mathcal{L}(a|\theta) p(\theta|\mathbf{y}) d\theta$
- Bayes Estimator of a parameter $\hat{\theta}_{BE} = \arg \min_{\hat{\theta} \in \Theta} R_\theta(\hat{\theta}|\mathbf{y})$

Common Loss Functions and Corresponding Estimators

- **Squared Error Loss:** $\mathcal{L}(\theta|\hat{\theta}) = (\theta - \hat{\theta})^2$, Bayes Estimator $\hat{\theta} = \mathbb{E}[\theta|\mathbf{y}]$
- **Absolute Error Loss:** $\mathcal{L}(\theta|\hat{\theta}) = |\theta - \hat{\theta}|$, Bayes Estimator $\hat{\theta} = \theta_{0.5}$, $(\mathbb{P}(\theta \leq \theta_{0.5}) = 0.5)$
- **0-1 Loss:** $\mathcal{L}(\theta|\hat{\theta}) = I(\theta \neq \hat{\theta})$, Bayes Estimator $\hat{\theta}$ is the posterior mode

8 Interval Estimates

$$P \times 100\% \text{ Bayesian Credible interval} = [LB, UB] \text{ where } \int_{LB}^{UB} p(\theta, \mathbf{y}) d\theta = P$$

If looking for one side-confidence bounds then $LB = -\infty$ or $UB = \infty$

Symmetric Credible Interval

$$[LB, UB] \text{ where } \mathbb{P}(\theta \leq LB) = \mathbb{P}(\theta \geq UB) = \alpha/2 = (1 - P)/2$$

Highest Posterior Density Interval (HPDI)

Credible Interval where all values outside the interval has a density smaller than any of the values in the interval.

Credible Regions

$$\iint p(\theta_1, \theta_2 | \mathbf{y}) d\theta_1 d\theta_2 = 1 - \alpha = P$$

9 Classic Hypothesis Testing

1. Assume hypothesis H_0 is true
2. Calculate test statistic $T(\mathbf{y}_{obs})$ based on observed sample data regarding to H_0 and H_1
3. Conditional on H_0 being true, $p\text{-value} = \mathbb{P}(T(\mathbf{y}) \text{ more extreme than } T(\mathbf{y}_{obs}) | \theta, H_0)$
4. Reject H_0 and accept H_1 for sufficiently small $p\text{-values}$, do not reject H_0 otherwise

10 Bayesian Hypothesis Testing

Suppose there are two hypotheses about parameter θ :

$$H_0 : \theta \in \Theta_0 \quad H_1 : \theta \in \Theta_1$$

where $\Theta_0 \cup \Theta_1 = \Theta$ and $\Theta_0 \cap \Theta_1 = \emptyset$.

The Bayesian approach specifies prior probabilities on each hypotheses:

$$\begin{aligned} p_0 &= \mathbb{P}(H_0 \text{ is true}) = \mathbb{P}(\theta \in \Theta_0) \\ p_1 &= \mathbb{P}(H_1 \text{ is true}) = \mathbb{P}(\theta \in \Theta_1) \end{aligned}$$

Where the posteriors are as follows:

$$\begin{aligned} \mathbb{P}(H_0 | \mathbf{y}) &= \mathbb{P}(\theta \in \Theta_0 | \mathbf{y}) \\ \mathbb{P}(H_1 | \mathbf{y}) &= \mathbb{P}(\theta \in \Theta_1 | \mathbf{y}) \end{aligned}$$

Where $\mathbb{P}(H_0 | \mathbf{y}) + \mathbb{P}(H_1 | \mathbf{y}) = 1$

Simple

- Single parameter values $H_0 : \theta = \theta_0$ $H_1 : \theta = \theta_1$
- Posteriors:

$$\begin{aligned}\mathbb{P}(H_0|\mathbf{y}) &= \mathbb{P}(\theta = \theta_0|\mathbf{y}) = \frac{f(\mathbf{y}|\theta_0)p_0}{m(\mathbf{y})} = \frac{f(\mathbf{y}|\theta_0)p_0}{f(\mathbf{y}|\theta_0)p_0 + f(\mathbf{y}|\theta_1)p_1} \\ \mathbb{P}(H_1|\mathbf{y}) &= \mathbb{P}(\theta = \theta_1|\mathbf{y}) = \frac{f(\mathbf{y}|\theta_1)p_1}{m(\mathbf{y})} = \frac{f(\mathbf{y}|\theta_1)p_1}{f(\mathbf{y}|\theta_0)p_0 + f(\mathbf{y}|\theta_1)p_1} \\ \mathbb{P}(H_0|\mathbf{y}) &= 1 - \mathbb{P}(H_1|\mathbf{y})\end{aligned}$$

- posterior odds of H_0 against H_1 :

$$\frac{\mathbb{P}(H_0|\mathbf{y})}{\mathbb{P}(H_1|\mathbf{y})} = \frac{f(\mathbf{y}|\theta_0)p_0}{f(\mathbf{y}|\theta_1)p_1}$$

Composite

- Single parameter values $H_0 : \theta \in \Theta_0$ $H_1 : \theta \in \Theta_1$
- Prior: $p_i = \int_{\theta \in \Theta_i} \pi(\theta) d\theta$
- Posteriors:

$$\begin{aligned}\mathbb{P}(H_i|\mathbf{y}) &= \frac{p(\mathbf{y}, H_i)}{m(\mathbf{y})} = \frac{p_i p(\mathbf{y}|H_i)}{m(\mathbf{y})} = \frac{p_i \int p(\mathbf{y}, \theta|H_i) d\theta}{m(\mathbf{y})} = \frac{p_i \int f(\mathbf{y}|\theta) \pi(\theta|H_i) d\theta}{m(\mathbf{y})} \\ &= \frac{p_i \int_{\theta \in \Theta_i} f(\mathbf{y}|\theta) \frac{\pi(\theta)}{p_i} d\theta}{m(\mathbf{y})} = \frac{\int_{\theta \in \Theta_i} f(\mathbf{y}|\theta) \pi(\theta) d\theta}{m(\mathbf{y})} = \int_{\theta \in \Theta_i} p(\theta|\mathbf{y}) d\theta = \mathbb{P}(\theta \in \Theta_i|\mathbf{y})\end{aligned}$$

- posterior odds of H_0 against H_1 :

$$\frac{\mathbb{P}(H_0|\mathbf{y})}{\mathbb{P}(H_1|\mathbf{y})} = \frac{\int_{\theta \in \Theta_0} f(\mathbf{y}|\theta) \pi(\theta) d\theta}{\int_{\theta \in \Theta_1} f(\mathbf{y}|\theta) \pi(\theta) d\theta} = \frac{\mathbb{P}(\theta \in \Theta_0|\mathbf{y})}{\mathbb{P}(\theta \in \Theta_1|\mathbf{y})} = \frac{\mathbb{P}(\theta \in \Theta_0|\mathbf{y})}{1 - \mathbb{P}(\theta \in \Theta_0|\mathbf{y})}$$

Multiple models

- Set of models M_1, \dots, M_K
- H_i : The correct model is M_i

$$\mathbb{P}(H_i|\mathbf{y}) = \frac{\mathbb{P}(H_i, \mathbf{y})}{\mathbb{P}(\mathbf{y})} = \frac{\mathbb{P}(H_i, \mathbf{y})}{\sum_{j=1}^K \mathbb{P}(H_j, \mathbf{y})}$$

Bayes Factor

- Prior odds for H_0 against $H_1 = p_0/p_1$
- Posterior odds for H_0 against $H_1 = \mathbb{P}(H_0|\mathbf{y})/\mathbb{P}(H_1|\mathbf{y})$
- Bayes Factor for H_0 against H_1 :

$$BF_{01} = \frac{\text{Posterior odds for } H_0 \text{ against } H_1}{\text{Prior odds for } H_0 \text{ against } H_1} = \frac{\mathbb{P}(H_0|\mathbf{y})/\mathbb{P}(H_1|\mathbf{y})}{p_0/p_1}$$

$$BF_{10} = \frac{1}{BF_{01}}$$

BF_{ij}	Interpretation
< 3	No edvidence for H_i over H_j
> 3	Positive edvidence for H_i
> 20	Strong edvidence for H_i
> 150	Very strong edvidence for H_i

Simple H_0 vs Simple H_1

$$BF_{01} = \frac{\mathbb{P}(H_0|\mathbf{y})/\mathbb{P}(H_1|\mathbf{y})}{p_0/p_1} = \frac{(f(\mathbf{y}|\theta_0)p_0)/(f(\mathbf{y}|\theta_1)p_1)}{p_0/p_1} = \frac{f(\mathbf{y}|\theta_0)}{f(\mathbf{y}|\theta_1)}$$

Composite H_0 vs Composite H_1

$$BF_{01} = \frac{\mathbb{P}(H_0|\mathbf{y})/\mathbb{P}(H_1|\mathbf{y})}{p_0/p_1} = \frac{\left[\int_{\theta \in \Theta_0} f(\mathbf{y}|\theta)\pi(\theta)d\theta \right] / \left[\int_{\theta \in \Theta_1} f(\mathbf{y}|\theta)\pi(\theta)d\theta \right]}{p_0/p_1} = \frac{\mathbb{P}(\theta \in \Theta_0|\mathbf{y})/\mathbb{P}(\theta \in \Theta_1|\mathbf{y})}{p_0/p_1}$$

Simple H_0 vs Composite H_1

$$BF_{01} = \frac{\mathbb{P}(H_0|\mathbf{y})/\mathbb{P}(H_1|\mathbf{y})}{p_0/p_1} = \frac{[f(\mathbf{y}|\theta_0)p_0] / \left[p_1 \int_{-\infty}^{\infty} f(\mathbf{y}|\theta)\pi(\theta)d\theta \right]}{p_0/p_1} = \frac{f(\mathbf{y}|\theta_0)}{\int_{-\infty}^{\infty} f(\mathbf{y}|\theta)\pi(\theta)d\theta} = \frac{f(\mathbf{y}|\theta_0)}{m(\mathbf{y})}$$

Multipes Hypotheses

Read 5B.7 cba to type

11 Deterministic Numerical Integration

$$\int_{x_i}^{x_{i+m}} f(x)dx \approx \sum_{j=1}^m w_{ij}f(x_{i+j})$$

12 Monte Carlo Integration

$$\text{Expectation } \mathbb{E}[\theta] = \int \theta p(\theta) d\theta$$

Samples of $\theta : \theta_1, \dots, \theta_N$

$$\Rightarrow \text{Monte Carlo Estimate } \hat{\mathbb{E}}[\theta] = \frac{1}{N} \sum_{i=1}^N \theta_i \approx \mathbb{E}[\theta]$$

By the Law of Large Numbers, $\mathbb{P}(\lim_{N \rightarrow \infty} \hat{\mathbb{E}}[\theta] = \mathbb{E}[\theta]) = 1$.

Many Bayesian integrals can be viewed as expectations, such as probabilities and normalising constants.

Direct Sampling

- Target Distribution, often the posterior $p(\theta|\mathbf{y})$
- Integrals of interest: $\mathbb{E}[h(\theta|\mathbf{y})] = \int h(\theta)p(\theta|\mathbf{y})d\mathbf{y}$
 - Common choices for $h(\theta)$:
 - $h(\theta) = \theta$, the posterior mean
 - $h(\theta) = (\theta - \mathbb{E}[\theta])^2$, the posterior variance
 - $h(\theta) = I(a \leq \theta \leq b)$, $\mathbb{P}(\theta \in [a, b])$
- Samples $\theta_1, \dots, \theta_N$
- Monte Carlo estimation of Integrals of interest $\mathbb{E}[h(\theta|\mathbf{y})] \approx \hat{\mathbb{E}}[h(\theta|\mathbf{y})] = \frac{1}{N} \sum h(\theta_i)$
- Monte Carlo Error = $\hat{\mathbb{E}}[h(\theta|\mathbf{y})] - \mathbb{E}[h(\theta|\mathbf{y})]$

Inverse Probability Integral Transform Method (Inverse PIT)

Continuous Random Variable

- Let Y be continuous random variables with cdf $F_Y(y) \in [0, 1]$
- Define new random variable $U = F_Y(Y) \sim \text{Uniform}(0, 1)$

The Inverse PIT method

- Transform a uniform variable U using the inverse CDF of Y ($g(U) = F_Y^{-1}(U)$). $g \sim Y$
- Therefore if one can evaluate F^{-1} for Y , then sample y can be generated by $y = F^{-1}(u)$, where $u \sim \text{Uniform}(0, 1)$

Discrete Random Variable

The Inverse PIT method

- Same for Continuous method where:

$$F^{-1}(u) = \min_y \{y : F_Y(y) \geq u\}$$

Rejection Sampling

Target distribution $p(\theta|\mathbf{y})$

- Generates independent samples from envelope distribution $g(\theta)$
- Only some generated values are kept/used, while the rest are discarded
- Often used for generating univariate random variables than a multivariate random vector

General Rejection Algorithm

For target distribution $p(\theta)$ and envelope $g(\theta)$, the (upper) bound M of p/g is defined:

$$M \geq \frac{p(\theta)}{g(\theta)} \quad \forall \theta \Leftrightarrow M g(\theta) \geq p(\theta)$$

The Algorithm to generate a single value from $p(\theta)$ is as follows:

1. Generate θ^* from $g(\theta)$
2. Generate u from $\text{Uniform}(0, M g(\theta^*))$
3. If $u \leq p(\theta^*)$ keep θ^* else reject and return step 1

The acceptance rate is $\frac{1}{M}$ which M is defined as above. For optimal acceptance rate, we choose $M_{opt} = \sup \left\{ \frac{p(\theta)}{g(\theta)} \right\}$

Importance Sampling

Similar to Rejection sampling. However, we adjust the values to correspond to the right distribution instead of rejecting them. Therefore all generated θ are kept, there's no rejection.

Suppose there is another distribution has the same support as $p(\theta|\mathbf{y})$, the integral then can be written as $\mathbb{E}_p[h(\theta|\mathbf{y})] = \mathbb{E}_g[h(\theta)w(\theta)]$ where $w(\theta) = p(\theta|\mathbf{y})/g(\theta)$ is the **importance ratio**.

The Monte Carlo estimates can be calculated $\hat{\mathbb{E}}_p[h(\theta|\mathbf{y})] \equiv \hat{\mathbb{E}}_g[h(\theta)w(\theta)] = \frac{1}{N} \sum h(\theta^i)w(\theta^i)$

Sampling Importance Re-Sampling

Algorithm for generating a sample size of n

1. Choosing $N > n$. generate N samples of $\theta^i, i = 1, \dots, N$ from the envelope distribution (importance sampler) $g(\theta)$
2. Calculate the importance ratio $w_i = p(\theta^i|\mathbf{y})/g(\theta)$, then scale by the sum of weights:

$$w^{*i} = \frac{w^i}{\sum_j^N w^j}$$

3. randomly sample n θ^i 's with replacement, with probabilities w^{*1}, \dots, w^{*N}

13 Markov Chain Monte Carlo (MCMC)

Overview

- **Dependent Samples:** $\theta|y$
- **Iterative Conditional Generation:** $g(\theta^i|\theta^{i-1})$
- **Burn-in:** The initial values $\theta^1, \dots, \theta^B$ are not distributed according to the target distribution $p(\theta)$, therefore would be discarded
- **Sample for Inference:** The additional $N - B$ samples $\theta^{B+1}, \dots, \theta^N$ are used for inference

Markov Chain

- **State Space S :** is a set of all possible states j
- **Process $\theta^t = j$:** process at time t is in state j

Definition

$$\begin{aligned}\mathbb{P}(\theta^0 = x_0, \dots, \theta^T = x_T) &= \mathbb{P}(\theta^T = x_T | \theta^0 = x_0, \dots, \theta^{T-1} = x_{T-1}) \\ &\quad \times \mathbb{P}(\theta^{T-1} = x_{T-1} | \theta^0 = x_0, \dots, \theta^{T-2} = x_{T-2}) \\ &\quad \times \mathbb{P}(\theta^0 = x_0)\end{aligned}$$

Markov Property: if $\theta^t | \theta^{t-1}$ is independent to other values, then $\mathbb{P}(\theta^T = x_T | \theta^0 = x_0, \dots, \theta^{T-1} = x_{T-1}) = \mathbb{P}(\theta^T = x_T | \theta^{T-1} = x_{T-1})$. Therefore:

$$\mathbb{P}(\theta^0 = x_0, \dots, \theta^T = x_T) = \prod_{t=1}^T \mathbb{P}(\theta^t = x_t | \theta^{t-1} = x_{t-1}) \times \mathbb{P}(\theta^0 = x_0)$$

The sequence of random variables $\theta^1, \dots, \theta^K$ with Markov Property is a Markov Chain

Terminology and Notation

- **One-Step Transition Probability** $p_{ij}^t = \mathbb{P}(\theta^t = j | \theta^{t-1} = i)$
- if $p_{ij}^t = p_{ij} \forall t$ then the chain is **time homogeneous** otherwise **time inhomogeneous**
- **Transition probability matrix** $(P_t)_{ij} = p_{ij}$, where each row sums to one
for a time homogeneous chain $P_t = P \forall t$

Limiting Behaviour of Markov Chain

- **Recurrent State:** state where a Markov Chain returns with probability 1
- **Nonnull State:** state where it will eventually recurrence in a finite expected time
- **Irreducible Markov Chain:** $\exists m \in \mathbb{N} \text{ s.t. } \mathbb{P}(\theta^{m+n} = j | \theta^n = i) > 0 \forall i$
- **Periodic Markov Chain:** $\exists m \in \mathbb{N} \text{ s.t. } \exists d \in \mathbb{N} \text{ s.t. } \mathbb{P}(\theta^{m+n} = j | \theta^n = i) > 0 \cap m/d \in \mathbb{Z}^+ \forall i$
otherwise **Aperiodic Markov Chain**
- **Ergodic Markov Chain:** Irreducible, aperiodic, all states are **nonnull** and **recurrent**

Marginal and Stationary Distribution

- **Marginal Probability Distribution** for state at time t is denoted π_t , which is a vector of probabilities that sum to one
- $(\pi_t)_i = \pi_t(i) = \mathbb{P}(\theta^t = i)$
- $\pi_{t+1}^\top = \pi_t^\top P \Leftrightarrow \pi_{t+1} = P^\top \pi_t$
- If $\pi^\top P = \pi^\top$ then π is a **stationary distribution**

Reversible Markov Chain

If a time homogeneous Markov Chain with tpm P and Stationary distribution π fulfill the **Detailed Balance equation** then it is a reversible Markov Chain.

Detailed Balance equation is written $\pi(i)p_{ij} = \pi(j)p_{ji}$

Key Theoretical Results

If a Markov Chain with tpm P and is irreducible, aperiodic and has a stationary distribution π , the the followings are true:

- π is unique, $\exists! \pi \text{ s.t. } \pi^\top P = \pi^\top$
- $\lim_{n \rightarrow \infty} \mathbb{P}(\theta^{t+n} = j | \theta^t = i) = \pi(j)$
- π is the follows solution to the following program:

$$\sum_{i \in S} \pi(i) = 1$$

$$\pi(j) = \sum_{i \in S} \pi(i)p_{ij}$$

Subject to : $\pi(j) \geq 0 \forall j \in S$

14 Metropolis-Hasting Algorithm

Steps

Iteration t given sample Θ^{t-1} :

1. Generate candidate value Θ^c from proposal distribution $q(\Theta^c|\Theta^{t-1})$
2. Calculate Metropolis Hasting Ratio (MHR) aka. acceptance rate if < 1 :

$$\text{MHR}(\Theta^{t-1}, \Theta^c) = \frac{p(\Theta^c)q(\Theta^{t-1}|\Theta^c)}{p(\Theta^{t-1})q(\Theta^c|\Theta^{t-1})}$$

3. Generate $u \sim \text{Uniform}(0, 1)$
4. if $u \leq \min(1, \text{MHR}(\Theta^{t-1}, \Theta^c))$, then $\Theta^t = \Theta^c$ (Keep Θ^c) Otherwise $\Theta^t = \Theta^{t-1}$
5. Set $t = t + 1$, return step 1

Special case proposals

Random walks ($\theta^c = \theta^o + \epsilon$)

$$\text{Special cases for } \epsilon: \epsilon \sim \begin{cases} \text{Normal}(0, \sigma^2) \\ t_2 \text{ df} \\ \text{Uniform}(-b, b) \end{cases} \quad \text{where } q(\theta^c|\theta^o) = q(\theta^o|\theta^c) \Rightarrow \text{MHR}(\theta^o, \theta^c) = \frac{p(\theta^c)}{p(\theta^o)}$$

Independence proposals

$$\text{MHR}(\theta^o, \theta^c) = \frac{p(\theta^c)q(\theta^o)}{p(\theta^o)q(\theta^c)}$$

Multidimensional Θ : One-at-a-time updates

let $\Theta_U^{t+1}, \Theta_0^t \subset \Theta$ be the set of updated or yet to be updated $\theta \in \Theta$. Then:

$$\text{MHR}(\theta_i^t, \theta_i^c) = \frac{p(\Theta_U^{t+1}, \theta_i^c, \Theta_0^t \setminus \{\theta_i^c\})q(\theta_i^t|\Theta_U^{t+1}, \theta_i^c, \Theta_0^t \setminus \{\theta_i^c\})}{p(\Theta_U^{t+1}, \Theta_0^t)q(\theta_i^c|\Theta_U^{t+1}, \Theta_0^t)}$$

15 MCMC Diagnostics

Notations

- **Burn-in Length B**

The idea burn-in length B is the point where the chain is converged (aka. future values will be in the limiting distribution)

- **Trace plots/Sample paths**

A time series plot of chains of generated values

Multiple Chains

The reason to plot multiple chains for a single random values is that a random variable θ may have multiple modes θ_1, \dots , therefore we would like to capture as many as possible. A plot of such chains does not prove convergence, but instead can indicate a lack of convergence.

Let's say variables $\theta_1, \dots, \theta_K$ converge at $B = \{B_1, \dots, B_K\}$ respectively, then we say any sample length $N \geq \max B$ would indicate a "good mixing"

A more quantitative assessment of a good mixing/lack of convergence is the **Brooks-Gelman-Rubin** (BGR) statistics, and is defined as below:

$$R = \frac{\text{Width of 80\% credible interval of all chains combined}}{\text{Average width of 80\% credible interval for each chain}}$$

If the chain has converged, $R \approx 1$. If the chains have yet been converged then $R > 1$

Variance of Monte Carlo Estimate $\hat{\mathbb{E}}[\theta]$

Notation

- $\gamma_k = \text{Cov}(\theta^i, \theta^{i+k})$: autocovariance of lag $k \geq 0$ of the chain
- $\sigma^2 = \gamma_0$: variance of the chain
- $\rho_k = \gamma_k / \sigma^2$: autocorrelation of lag k
- τ_n^2 / n : Variance of $\hat{\mathbb{E}}[\theta]$

$$\begin{aligned} \text{Var}[\hat{\mathbb{E}}[\theta]] &= \frac{\tau_n^2}{n} \\ \tau_2^2 &= \sigma^2 \left[1 + 2 \sum_{k=1}^n \frac{n-k}{n} \rho_k \right] \quad (\text{Inefficiency factor / Integrated autocorrelation}) \end{aligned}$$

As $n \rightarrow \infty$:

$$\tau_2^2 = \sigma^2 \left[1 + 2 \sum_{k=1}^n \frac{n-k}{n} \rho_k \right] \approx \sigma^2 \left(1 + 2 \sum_{k=1}^{\infty} \rho_k \right) \Rightarrow \text{Var}[\hat{\mathbb{E}}[\theta]] \approx \frac{\sigma^2 (1 + 2 \sum_{k=1}^{\infty} \rho_k)}{n}$$

Effective Sample Size, ESS

$$n_{eff} = \frac{n}{1 + 2 \sum_{k=1}^n \rho_k} \Rightarrow \text{Var} \left[\hat{\mathbb{E}}[\theta] \right] \approx \frac{\sigma^2}{n_{eff}}$$

Autocorrelation Plots

$$\hat{\rho}_k = \frac{\hat{\gamma}_k}{\hat{\sigma}^2} = \frac{\sum_{i=1}^{n-1} (\theta^i - \bar{\theta}) (\theta^{i+k} - \bar{\theta})}{\sum_{i=1}^n (\theta^i - \bar{\theta})^2}, \quad \bar{\theta} \text{ is the sample average}$$

Central Limit Theorem for $\hat{\mathbb{E}}[\theta]$

$$\hat{\mathbb{E}}[\theta] \sim \text{Asymptotically Normal} \left(\mathbb{E}[\theta], \frac{\sigma^2}{n_{eff}} \right) \Rightarrow \frac{\hat{\mathbb{E}}[\theta] - \mathbb{E}[\theta]}{\sigma/\sqrt{n_{eff}}} \sim \text{Normal}(0, 1)$$

Estimating $\text{Var} \left[\hat{\mathbb{E}}[\theta] \right]$

Batching is a popular method for estimating $\text{Var} \left[\hat{\mathbb{E}}[\theta] \right]$, whe algorithm is as below:

1. Partition chain of length n into m batches of $T \in [10, 30]$ successive values
2. For each batch calculate the mean of the values $\bar{\theta}_i$
3. Estimate of $\text{Var} \left[\hat{\mathbb{E}}[\theta] \right]$ is then calculated:

$$\frac{\hat{\tau}^2}{n} = \frac{\frac{T}{m-1} \sum_{i=1}^m (\bar{\theta}_i - \bar{\theta})^2}{n}$$

Improve Performance

- Changing proposal distribution
- Reparameterising model
- Blocking

Gibbs Sampler

Can be seen as a special case of a Metropolis-Hasting sampler with two features:

- The proposal distributions are conditional distributions for θ 's
- All candidates values are accepted ($\text{MHR}(\theta_i^t, \theta_i^c) = 1$)

The Algorithm

1. Let $\Theta = (\theta_1, \dots, \theta_q) \sim p(\Theta)$
2. For $i \in 1, \dots, q$ denote the full conditional distribution for θ_i by

$$p(\theta_i | \theta_{-i}) = p(\theta_i | \theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_q)$$

3. Initialise chain $\Theta^0 = (\theta_1^0, \dots, \theta_q^0)$
4. At iteration $t + 1$ given Θ^t , Θ^{t+1} is generated as follows:
 - Generate θ_1^{t+1} from $p(\theta_1 | \theta_2^t, \dots, \theta_q^t)$
 - \vdots
 - Generate θ_i^{t+1} from $p(\theta_i | \theta_1^{t+1}, \dots, \theta_{i-1}^{t+1}, \theta_{i+1}^t, \dots, \theta_q^t)$
 - \vdots
 - Generate θ_q^{t+1} from $p(\theta_q | \theta_1^{t+1}, \dots, \theta_{q-1}^{t+1})$
5. The (joint) transition probability of going from Θ^t to Θ^{t+1} is:

$$\mathcal{K}_{\text{Gibbs}}(\Theta^t, \Theta^{t+1}) = \prod_{i=1}^q p(\theta_i^{t+1} | \theta_j^{t+1}, j < i \cap \theta_j^t, j > i)$$

An the stationary distribution is the target $p(\Theta)$

Gibbs Sampler vs Metropolis-Hasting

Gibbs Sampler

- Strengths
 - Automatically defined proposal dist.
 - Accepts all candidate values
 - views as an "adaptive algorithm"
- Weaknesses
 - Conditional distribution may not be tractable
 - Sampling may be computationally intensive
 - 100% acceptance \neq good mixing

Metropolis Hasting

- Strengths
 - flexible proposal dist., can be fast sample from, fast to evaluate MHR
 - Don't need to know conditional dist.
 - easily block updating
- Weaknesses
 - Hard to find good proposal with good mixing
 - Can take time to tune a proposal