ch 2.
- what is conjugate prior.
- change of variable · { why not uniform
- Jeffreys prior :     $\det(I(\theta))^{\frac{1}{2}}$
        what it connects to data

- Non - Informative    · Reference Prior.
- Expert prior      ·
· KL distance

· { ·Direct integration
· { · with (hyper) prior
why? Integration is hard.
·Bayesian Regression
· Bayesian variable selection

ch 4.
· Laplace approximation
· Independent MC methods

Bayesian Testing

$H_0:$   $\theta$ belongs to some interval.

$\quad \theta \in \Theta_0$

$\pi(\theta_0 | x) = \int_{\Theta_0} \pi(\theta | x)$

• Bayes Factor

$$\frac{\pi(\theta_0 | x)}{\pi(\theta_1 | x)} \left/ \overline{\left( \frac{\pi(\theta_1)}{\pi(\theta_0)} \right)} \right. \quad \text{hidr of w/e (hyper)priors.}$$

"Posterior" odds   ×   "prior odds"

• A subset of Lesbegue measure zero.

e.g. $\Theta_0 = \{ \mu = \mu_0 \}$   zero measure

$\Rightarrow \pi(d\theta) = P_0 \pi_0(d\theta) + (1-P_0) \pi_1(\theta) d\theta$

$\quad \pi_0:$ dist. concentrated on $\Theta_0$
$\quad P_0:$ prior prob.

$p$-values : Frequent: evidence against null hypothesis

$$\pi(\theta_0 | x) = \frac{f(x | \theta_0)}{f(x|\theta_0) + \int \pi_1}$$   $\Theta_0 = \{\theta_0\}$

$$\geqslant \frac{f(x|\theta_0)}{f(x|\theta_0) + \sup_\theta f(x|\theta)}$$

$\beta$ - Bayesian confidence set with level $(1-\alpha)$

$(1-\alpha)$ creditable set.

$\quad P(\theta \in C_x | X = x) = \pi(C_x | x) \geqslant 1-\alpha$

$\qquad\qquad\qquad\qquad\qquad \underset{\text{depend on } x}{\uparrow}$

HPD. the one minimizing the volume

$\quad L_k = \{ \theta; \pi(\theta | x) \geqslant k \} \quad k_\alpha = \sup \{ k; \pi(L_k | x) \geqslant 1-\alpha \}$

why.

· △

---

Bayesian asymptotics

Fisher Information

$\quad \hat{\theta}_n \text{ MLE} \sim N(\theta_0, \frac{1}{I_n})$

Bayesian   $\theta | x_n \sim N(\hat{\theta}_n, \frac{1}{n} I^{-1})$

influence of prior vanish.

Bayes formula    1.1

$i$ countable set

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_k P(B|A_i)P(A_i)}$$

$P(A_i) > 0$
$P(B) > 0$

continous case.

$$f_{x|y}(x|y) = \frac{f_{xy}(x,y)}{f_y(y)} =$$

$$f_{y|x}(y|x) = \frac{f(x,y)}{f(x)} = \frac{f(x|y)f(y)}{f(x)} = \overline{\frac{}{f_R}}$$

- Frequentist: Relatively frequency of $A_i$ among those happening with $B$
- Bayesian: describing the uncertenty after we saw the data.

- How to choose a prior?
  1. cannot depend on data.
  2.

prior and posterior predictive distribution:

$f(x) = \int f(x|\theta) \pi(\theta) d\theta$ is called marginal density.

prior predictor dens

After X observed:
Could observe distribution of Y from the same model

Posterior predictive distribution:

$$f(y|x) = \int f(y, \theta|x) d\theta$$
$$= \int f(y|x, \theta) \cdot \pi(\theta|x) d\theta$$
$$= \int f(y|\theta) \pi(\theta|x) d\theta \quad \mathcal{R} \quad Y \perp X | \theta$$

$\triangle$

Example $n+1$ th observation $\in N$

page3: $\exp\left(-\frac{1}{2}\left(\sum x_i - \mu)^2 - \frac{n\tau^2}{1+n\tau^2}(\bar{x} - \mu)^2\right)\right)$

—

$\begin{cases} \text{frequenst risk} \\ \text{Bresha} \\ \quad \text{weighted} \end{cases}$

Point Estimation

Summarize Baye dist. into number.

Loss function  $L: \Theta^2 \longrightarrow [0, \infty)$
Find the minimize of risk

$$\hat{\theta} = \arg\min_T \rho(T(x), \pi) \quad \frac{\rho(T(x), \pi)}{= \mathbb{E}[L(T(x), \theta|x)]}$$
$$= \int_\Theta L(T(x), \theta) \pi(\theta|x) d\theta$$

$\Rightarrow \exists$ ye

$T_1 \quad I_2$ are good for a particular estimator
but not all  $R(T_1, \theta_1) < R(T_2, \theta_1)$ but $R(T, \theta_2) > +$

$\Rightarrow$ look for estimator which minimizes the max risk
$\sup_\theta R(T, \theta)$

$$R(T, w) = \int_\Theta R(T, \theta) w(\theta) d\theta$$

If $\int w(\theta) d\theta = 1$  $w$ as prior for $\theta$.

$T(x) = \arg\min_T \rho(T(x), w) = \arg\min_T \mathbb{E}[L(T, \theta)|x]$
is well defined for  $f(x) = \int f(x|\theta) w(\theta) d\theta$
Then any minima $T'$ is almost surely $T$

Admissible   No other estimator is better
$$R(T', \theta) \leq R(T, \theta)$$

Bayes estimator are often biased •

$I_{\alpha}$

$I_\sigma$

## 2.1 Conjugate prior

**Def** $P_\Xi = \{ \pi_\xi(\theta), \xi \in \Xi \}$ of prior densities is called conjugate

if $\forall x \in P_\Xi$, $\pi(\theta | x) \in P_\Xi$

Example: Normal $\to$ Normal
Gamma $\to$ Gamma.

- choosing a prior 's hyper parameters is even harder than choosing parameters itself. Similar in the case of Normal mean, it tells how much we want to rely on the $\pi(\xi)$

## 2.2 Non-informative Priors

Problems with const. uniform/flat prior.

- $\Theta$ must be finite which is often not the case
- under c.o.v. the density is NOT invariant.
- ⓔⓖ $\eta = g(\theta)$ be the new parameter

$$\lambda(\gamma) := \pi(g^{-1}(\theta)) |\det Dg^{-1}(\gamma)| \qquad \text{when } g \text{ is NOT}$$

linear, $\lambda(\gamma)$ is NOT const.

**Improper priors** : a measure that is NOT a probability measure.

if $\pi(\theta) f(x|\theta)$ has FINITE mass, then inference can be made easily as before. and this is NOT a easy task to check

**Jeffrey's Prior** $\qquad \pi(\theta) \propto \det(I(\theta))^{\frac{1}{2}}$

eg. $X \sim N(\theta, 1)$ $I(\theta) = 1$ $\qquad \to$ improper on $\mathbb{R}$.
$X \sim N(0, \theta^2)$ $I(\theta) = 2/\theta^2$ $\qquad \sim \frac{1}{\theta}$

- Reasoning behind Jeffreys priors:
since $I(\theta)^{-1}$ is asymptotic variance of MLE.
thus it is a indicator of how much info. from data about $\theta$.

It's also Invariant under change of variables so. Good !

eg. $\gamma = g(\theta)$ $\qquad I_\gamma(\gamma) = Dg^{-1}(\gamma)^T I_\theta(g^{-1}(\gamma)) Dg^{-1}(\gamma)$

$\Rightarrow \pi_\gamma(\gamma) \propto \det I_\gamma(\gamma)^{\frac{1}{2}} = |\det Dg^{-1}| \det I_\theta(g^{-1}(\gamma))$
same results with. c.o.v. formula.

**Reference Priors**
If $X$ (data) has the largest impact then the impact of the prior is minimal.

- **KL distance** $\qquad KL(f, g) = \int f(x) \log \frac{f(x)}{g(x)} dx$

Idea is to maximize

$$I(x, \theta) = \int_X f(x) \int_\Theta \pi(\theta|x) \log \frac{\pi(\theta|x)}{\pi(\theta)} d\theta dx$$

$$= \int_\Theta \pi(\theta) \int_X f(x|\theta) \log \frac{\pi(\theta) f(x|\theta)}{\pi(\theta) f(x)} d\theta dx$$

$$= \int_\Theta \pi(\theta) \int_X f(x|\theta) \log \pi(\theta|x) dx d\theta -$$

$$\int_\Theta \pi(\theta) \log \pi(\theta) d\theta$$

is also unfeasible to find the maximum value.

How to solve

$$I_\infty(x, \theta) = \lim_n I((x_1 \cdots x_n), \theta) \qquad \text{, still infinite}$$

**Approximating.**

we again have Jeffrey's prior in the limit

hyperparameters that $\pi(\vartheta)$ depends on also has a prior distribution.

$$\pi(\xi)\,\pi(\vartheta|\xi)\,f(x|\vartheta)$$

○ The primary interest is posterior $\pi(\vartheta|x)$

**Way 1.** $\pi(\vartheta) = \int \pi(\vartheta|\xi)\,\pi(\xi)\,d\xi$

**Way 2** $\pi(\vartheta|x) = \int \pi(\vartheta,\xi|x)\,d\xi = \int \pi(\vartheta|x,\xi)\,\pi(\xi|x)\,d\xi$

$$* \propto \int \pi(\vartheta|x,\xi)\,\pi(\xi)\,f(x|\xi)\,d\xi$$

This leaves $\pi(\vartheta|x,\xi)$ and $f(x|\xi)$ / $\pi(\xi|x)$ to determine
$\underset{\textcircled{1}}{\qquad}$ $\underset{\textcircled{2}}{\qquad}$

$$\pi(\xi|x) = \frac{\pi(\vartheta,\xi|x)}{\pi(\vartheta|x,\xi)} \quad \text{or} \quad = \int \pi(\vartheta,\xi|x)\,d\vartheta$$

**Conjugate priors** have explicit expression for
$$\pi(\vartheta|x,\xi) \text{ also } f(x|\xi)$$

**1. Normal means**

Setting $f(x|\vartheta) \sim N(\vartheta, 1)$
prior: $\pi(\vartheta|\mu,\tau^2) \sim N(\mu,\tau^2)$ with $\mu$ fixed
○ $\tau^{-2} \sim \text{Gamma}(\gamma,\lambda)$

The posterior $\pi(\vartheta|\mu,\tau^2,x) \sim N\left(\frac{\mu + \tau^2\bar{x}}{1+\tau^2}, \frac{\tau^2}{1+\tau^2}\right)$

Using **way 1** Direct integration

$\pi(\vartheta) \propto \int (\tau^{-2})^{\frac{1}{2}} \exp\left\{-\tau^{-2}\frac{(\vartheta-\mu)^2}{2}\right\} (\tau^{-2})^{\gamma-1}\exp(-\lambda\tau^{-2})\,d\tau^{-2}$

$= \int u^{\gamma-1/2} \exp\left(-(\lambda + \frac{(\vartheta-\mu)^2}{2})u\right)\,du$

$\propto \frac{1}{(\lambda + \frac{(\vartheta-\mu)^2}{2})^{\gamma+1/2}}$  scaled and shifted

t~ dist. w/ $2\gamma$ d.o.f.

$\pi(\vartheta|x) \propto \pi(\vartheta)\,f(x|\vartheta)$ does NOT belong to standard family but we know sth about the posterior mode.
  is closed to $\bar{x}$ if prior and data are in conflict.

**Using Way 2**
$\pi(\vartheta|\tau^2,x)$ is already known
$f(x|\tau^2)$ is also kind of known by marginizing
$f(x|\tau^2) \propto (1+n\tau^2)^{-\frac{1}{2}} \exp\left\{-\frac{n}{2(1+n\tau^2)}(\bar{x}-\mu)^2\right\}$

$\Rightarrow$ marginal posterior of $\tau^2$ has little mass for small values. $\rightarrow$ It favors large values of $\tau^2$
    prior and data are in conflict

**Empirical Bayes for Normal means.**
$\hat{\xi}(x) = \arg\max_\xi f(x|\xi)$
For $f(x|\tau^2) \Rightarrow \hat{\tau}^2 = \max(0, (\bar{x}-\mu)^2 - \frac{1}{n})$
Means: we choose a wider prior when $\bar{x}$ and $\mu$ are in conflict.

$E[\vartheta|x,\hat{\tau}^2] = \begin{cases} \frac{\mu}{\bar{x} + \frac{1}{n\frac{1}{(\mu-\bar{x})^2}}} \end{cases}$

converges to $\mu$ as $n \uparrow \infty$
    $\Rightarrow$ EB methods tends to underestimate
  uncertainty

---

**2. Hierarchical Poisson Model**
likelihood $f(x|\theta_j)$ – Poisson  $\dfrac{e^{-\theta_j}\,\theta_j^{x_j}}{x_j!}$

Prior Gam. $\theta_j \sim \text{Gam}(\gamma,\lambda) = \dfrac{\lambda^\gamma}{\Gamma(\gamma)}\,\theta_j^{\gamma-1}\exp(-\lambda\theta_j)$
超先验分布 $\pi(\gamma,\lambda)$

**Way 1** $\pi(\theta_1 \cdots \theta_J) = \int \dfrac{\lambda^{J\gamma}}{\Gamma(\gamma)^J}\prod\theta_j^{\gamma-1}\exp\{-\lambda\sum\theta_j\}\,\pi(\gamma,\lambda)\,d\lambda\,d\gamma$
This is NOT a standard distribution.

**Way 2** $\pi(\theta_1 \cdots \theta_J|\gamma,\lambda,x) = \prod_{j}^{J} \pi(\theta_j|\gamma,\lambda,x_j)$
$f(x_1 \cdots x_J|\gamma,\lambda) = \prod^{J} f(x_j|\gamma,\lambda)$

$\pi(\theta_j|\gamma,\lambda,x_j) = \dfrac{(\lambda+1)^{\gamma+x_j}}{\Gamma(\gamma+x_j)}\cdot\theta_j^{\gamma+x_j-1}\exp(-(\lambda+1)\theta_j)$
$\underline{\text{Gamma}}$

To solve $f(x_j|\gamma,\lambda) = \int f(x_j|\theta_j)\,\pi(\theta|\gamma,\lambda,x_j)\,d\theta_j$
$\qquad = \int \text{Pois}(\theta_j)\cdot\text{Gamma}\;d\theta_j$
$\qquad = \dfrac{\Gamma(\gamma+x_j)}{x_j!\,\Gamma(\gamma)}\dfrac{\lambda^\gamma}{(\lambda+1)^{\gamma+x_j}}$
$\qquad = \binom{\gamma+x_j}{x_j}\left(\frac{\lambda}{\lambda+1}\right)^\gamma\left(\frac{1}{\lambda+1}\right)^{x_j}$

Negative binom!

$\Rightarrow \pi(\theta_j|x_1\cdots x_n) = \int \pi(\theta_j|\gamma,\lambda,x_j)\,\pi(\gamma,\lambda|x)\,d\gamma\,d\lambda$
$\propto \int \prod_j^{J} f(x_j|\gamma,\lambda)\,\pi(\gamma,\lambda)\,d\gamma\,d\lambda$

Cannot be Integrated in close form

**Empirical Bayes.**
$f(x|\gamma,\lambda) \sim \text{Negative Binomial}$
$\Rightarrow J\gamma\log\left(\frac{\lambda}{1+\lambda}\right) - \sum^J x_j\log(1+\lambda) + \sum\sum\log(\gamma+k)$
$\Rightarrow \frac{\hat{\gamma}}{\hat{\lambda}} = \bar{x}$

$y = \alpha 1 + X_\gamma \beta_\gamma + \varepsilon \qquad \varepsilon \sim N(0, \sigma^2)$

Assumption: $\gamma$ fixed $\alpha$ independent

$\pi(\alpha, \beta_\gamma, \sigma^2) \propto \pi(\sigma^2) \cdot \pi(\beta_\gamma | \sigma^2 \qquad$ since $\alpha$ is flat

$\qquad = \sigma^{-2} \pi(\beta_\gamma | \sigma^2)$

Popular choice of g-prior of zellner

$\pi(\beta_\gamma | \sigma^2) \sim N(0, g\sigma^2 (X_\gamma^T X_\gamma)^{-1})$

Posterior: $\beta_\gamma | \sigma^2, y \sim N(\frac{g}{g+1} \hat{\beta}_\gamma, \frac{g\sigma^2}{\sigma+1} (X_\gamma^T X_\gamma)^{-1})$

$\qquad \alpha | \sigma^2, y \sim N(\hat{\alpha}, \frac{\sigma^2}{n})$

$\hat{\beta}, \hat{\alpha}$ is the MLE.

posterior mean $\qquad \frac{g}{g+1} \hat{\beta}_\gamma + \frac{1}{g+1} \beta_\gamma^0$

a convex combination of prior mean and MLE

shrinkage towards zero

• $g \to \infty \qquad$ non-informative prior

$\qquad \beta(\gamma, 0) \to 0 \qquad$ in the limit we will

always choose empty model.

$x$   $k$ compont.

we have to know $x_i^t = \pi(x_i | x_{-i})$

Every $x_i$ needs to be known.

## MH algorithm.

- Reversible distribution:

$$\int_A \pi(x) P(x, B) dx = \int_B \pi(x) P(x, A) dx$$

ie $\mathbb{P}[X^t \in A, X^{t+1} \in B] = \mathbb{P}[X^{t+1} \in A, X^t \in B]$

- $\pi$ is also invariant

- $\pi(x) q(x, y) = \tilde{\pi}(y) p(y, x)$   $\forall x, y$.
  reversibly

### MH assumption

- generate a chain which has reversible distribution

可以随意选择吗？     No. Spuriosly $= 1$

usually does NOT hold for $x$

### Algorithm

1. Simulate $X^0$
2. For $t = 1 \ldots$
   ⓐ Generate $Y^t \sim q(X^{t-1}, x) dx$ and $U^t$
   ⓑ $X^t = \begin{cases} Y^t & \text{if } U^t \le a(X^{t-1}, Y^t) \\ X^{t-1} & \text{else} \end{cases}$

---

How to solve the problem of reversibility?

$p(x, y) = q(x, y)$          $p(y, x) = \dfrac{\tilde{\pi}(x) q(x, y)}{\tilde{\pi}(y)}$

                              ⊅ Arbitary Transition Kernel.

$\Rightarrow p(x, y) = \dfrac{\tilde{\pi}(y) q(y, x)}{\tilde{\pi}(x)}$      $p(y, x) = q(y, x)$

找最一         $p(x, y) \le q$
              $p(y, x) \le \ell$       $\forall x \ne y$

         $a(x, y) = \min\left( 1, \dfrac{\pi(y)}{\pi(x)} \dfrac{q(y, x)}{q(x, y)} \right)$

---

## Hamiltonian MC

- Drawbacks of Gibbs sampler / · MH
  step-wise is small.

Assume: we can evaluate the gradient efficiently

How: ⅁ consider a new target $\tilde{\pi}$
$$\tilde{\pi}(x, u) \propto \pi(x) \exp\left\{ -\frac{1}{2} u^T M^{-1} u \right\}$$   $M = diag(m_i)$
$U_i$ is called momentum variable.

Based on a determistics, invertible map $G(x, u)$

the transformation $G(x, u)$ is given by ODE

$$\frac{dx_i}{dt} = \frac{\partial H(x, u)}{\partial u_i} = \frac{u_i}{m_i}$$

$$\frac{du_i}{dt} = \frac{\partial H(x, u)}{\partial x_i} = \frac{\partial \log \pi(x)}{\partial x_i}$$     $0 \le t_i \le T$

---

## Remarks of MCMC.

- Biasness   $\mathbb{E} h_{N, r} \ne \int h(x) \pi(x) dx$

2. $Var(h_{N, r})$ is complicated as it includes the
   covariance term.

.4 Bayesian Computation

## 1. Laplace Approximation

$$\int h(\theta) q(\theta) d\theta \qquad \text{where} \quad \cdot q \text{ has max} \\ \cdot h \text{ arbitary smooth}$$

- $\log q(\theta) = \log q(\theta_0) - \frac{1}{2}(\theta-\theta_0)^T J (\theta-\theta_0)$

$\Rightarrow \int h(\theta) q(\theta) d\theta \approx h(\theta_0) q(\theta_0) (\det J_0)^{-\frac{1}{2}} (2\pi)^{\frac{p}{2}}$

$q(\frac{1}{N})$

Bayes Factor and BIC

$B_{12} = \dfrac{f_c \quad |M_1)}{f_c \quad |M_2)} = \int \pi \prod f_c(|\theta) d\theta_i$

$\approx \pi_i \hat{\theta}_i \prod f(x_i|\hat{\theta}_i) (\det n I(\theta))^{-\frac{1}{2}} (2\pi)^{\frac{p_i}{2}}$

$\Rightarrow \log f_c x_n |M_i) \approx \sum \log \quad - \frac{p_i}{2} \log n + O(1)$

BIC thus

## 2. Independent Monte Carlo Methods

on the basis of drawing independ. $U^t$

Basic idea is that we can use quantile function
and $(0,1)$ to generate RV $F^{-1}(u) = \inf\{x,\quad\}$
following certain distribution

But this is often intractable in high-dim setting.

- Simulate with a diff. dist $\gamma$ (Proposal)
- Connect to target $\pi$

$$r(x) = \frac{\pi(x)}{M \cdot \gamma(x)} \le 2 \quad \text{envelope}$$

Proof:
$$\mathbb{P}[x \in dx] = \mathbb{P}[\gamma \in dx | U \le \frac{\pi(\gamma)}{M\gamma(\gamma)}]$$
$$\propto \mathbb{P}[\gamma \in dx] \, \mathbb{P}[U \le \# | \gamma \in dx]$$
$$= \frac{\pi(x)}{M\gamma(x)} \gamma(x) dx \propto \pi(x) dx$$

Requirement know:
$\pi$ up to a normalizing const.

\# Rejection is high unless $\pi$ is close to $\gamma$
however, It is hard to find a really close one

## Importance Sampling

Goal: $\mathbb{E}_\pi(h(x)) = \int h(x) \pi(x)$

Instead of rejection $\Rightarrow$ weight them with an appropriate
weighting function.

$$\frac{1}{N} \sum_{t=1}^N h(\gamma^t) w(\gamma^t) \quad \text{where} \quad w(x) = \frac{\pi(x)}{\gamma(x)}$$

No need to bound this

These are methods to generate Independent

MCMC generate

$$x^t = G(x^{t-1}, U^t)$$
$\uparrow$ can be more in $\mathbb{R}^d$

$$P(x^t \in A | \longrightarrow) = \mathbb{P}[x^t \in A | x^{t-1}] = \mathbb{P}[x^{t-1}, A]$$

Determined by G
$$P(x, A) = P(G(x, u) \in A) = \mathbb{P}[u \in \{u; G(x,u) \in A\}]$$
How to specify G?

Some properties about MC.
Positivity: All state can be reached $\pi(A) > 0$

- Invariance / Stationary properties
$$\pi(A) = \int \pi(x) P(x, A) dx \quad \forall A$$