

$$P(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}, \quad \lambda > 0$$

$$\cancel{\Gamma(\lambda + n)} = \beta^n$$

$$\Gamma(\alpha) = \int_0^\infty \lambda^{\alpha-1} \cdot e^{-\lambda} d\lambda$$

Posterior for  $\lambda = \theta$

$$P(\lambda | X^{(n)}) = \frac{\lambda^{\alpha+n-1} \exp\{-\lambda(\beta + \sum_{i=1}^n x_i)\}}{\int_0^\infty \lambda^{\alpha+n-1} \exp\{-\lambda(\beta + \sum_{i=1}^n x_i)\} d\lambda}$$

$$\Rightarrow P(\lambda | X^{(n)}) = \frac{(\beta + \sum_{i=1}^n x_i)^{\alpha+n}}{\Gamma(\alpha+n)} \lambda^{\alpha+n-1} \exp\{-\lambda(\beta + \sum_{i=1}^n x_i)\}$$

Posterior is Gamma( $\alpha+n, \beta + \sum_{i=1}^n x_i$ )

$$\hat{\lambda}_{\text{prior}} = \frac{\alpha+n}{\beta + \sum_{i=1}^n x_i} \left| \frac{n}{n \cdot \bar{x}} \right|$$

$$\text{posterior mode} : \frac{\alpha+n-1}{\beta + \sum_{i=1}^n x_i}$$

$$\Gamma(\alpha+n) = \int_0^\infty \lambda^{\alpha+n-1} e^{-\lambda} d\lambda$$

$$\int_0^\infty \lambda^{\alpha+n-1} e^{-\lambda} \underbrace{e^{-\lambda(\beta + \sum_{i=1}^n x_i)}}_{\cancel{\lambda} - m} d\lambda$$

$$\lambda = \frac{m}{\beta + \sum_{i=1}^n x_i}$$

$$\pi^{\alpha+n-1} = m^{\alpha+n-1} \cdot (\beta + \sum_{i=1}^n x_i)^{-(\alpha+n-1)}$$

$$(\beta + \sum_{i=1}^n x_i)^{-\cancel{(\alpha+n-1)}} \cdot \frac{+}{\cancel{\beta + \sum_{i=1}^n x_i}} \Gamma(\alpha+n)$$

# Lec 3

MATH 459

Recall:  $X^{(n)} = \{X_1, \dots, X_n\}$ ,  $X_i \stackrel{iid}{\sim} f(x_i|\theta)$

Likelihood:  $L(\theta) \equiv L(\theta; X^{(n)}) = \prod_{i=1}^n f(x_i|\theta)$

(log likelihood:  $\log L(\theta) \equiv l(\theta)$ )

Posterior:  $P(\theta|X^{(n)}) = \frac{P(\theta) \cdot f(X^{(n)}|\theta)}{\int_{\Theta} P(\theta) f(X^{(n)}|\theta) d\theta}$

Question: If the  $p(\theta)$  is a density, then is  $P(\theta)f(X^{(n)}|\theta)$  also a density?

Where does this ~~from~~? posterior come from?

- Imagine  $\theta$  is a random variable

obviously  $X$  is a random variable

- Can not assume  $\theta, X$  are indep

- Let  $g(\theta, X)$  be the joint density

$$f(X|\theta) = \frac{g(\theta, X)}{m(x)}, P(\theta) > 0$$

the posterior is  $P(\theta|X) = \frac{g(\theta, X)}{m(x)}, m(x) > 0$

- in general, ~~max~~  $m(x)$  unknown.

- Can find  $m(x) = \int_{\Theta} g(\theta, x) d\theta$

- and we know

$$g(\theta, x) = f(x|\theta) \cdot P(\theta)$$

$$\text{thus } P(\theta|X) = \frac{g(\theta, X)}{m(x)} = \frac{P(\theta) \cdot f(x|\theta)}{\int_{\Theta} P(\theta) f(x|\theta) d\theta}$$

2 things we learn:

① denominator of posterior is marginal density for data

$$m(x) = \frac{P(\theta) f(x|\theta)}{\int_{\Theta} P(x^n|\theta) P(\theta) d\theta}$$

Basic Marginal Likelihood Identity  
of data

②  $m(x)$  : normalizing constant  $C$   
so that  $\int_{\Theta} P(\theta|x) d\theta = 1$

Interpretation:  
posterior  $\propto$  prior  $\times$  likelihood

Multi-parameter Model

$\theta$  is a vector

$$\theta = (\theta_1, \dots, \theta_d)$$

Then  $x^{(n)}$  has  $x_i \stackrel{iid}{\sim} f(x|\theta)$

$$\theta \in \Theta \subseteq \mathbb{R}^d$$

Joint prior density  $P(\theta_1, \dots, \theta_d)$

$$\int_{\Theta} \int_{\Theta} \dots \int_{\Theta} P(\theta_1, \dots, \theta_d) d\theta_1 \dots d\theta_d = 1$$

not always true

$$P(\theta_1, \dots, \theta_d) = P(\theta_1) \dots P(\theta_d)$$

- come up with a joint density

- for a subset of parameters of interest say  $\theta$   
 $P(\theta_1 | \theta_2, \dots, \theta_d)$

The likelihood equations:

$$\nabla_{\theta} l(\theta_1, \dots, \theta_d) = 0$$

$$\begin{pmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_d} \end{pmatrix}$$

$$m(x|X^n) = \frac{P(\theta|x^n) f(x|\theta)}{P(\theta|x)}$$

$$m(x|X^n) \cdot p(\theta|x, X^n) = p(\theta|X^n) f(x|\theta, X^n)$$

$$m(x|X^n) = \frac{p(\theta|X^n) f(x|\theta, X^n)}{p(\theta|x, X^n)}$$

$$\frac{m(x|X^n)}{\int m dx} = \frac{\frac{p(X^n|\theta) \cdot P(\theta)}{\int dx^n}}{\frac{p(\theta|X^n) \cdot P(X|\theta)}{\int dx dx^n}} = \int \frac{p(X^n|\theta) P(\theta)}{dx dx^n}$$

$\begin{array}{c} x^n \\ \downarrow \\ x \end{array}$

Suppose a Bayesian wants to inference for  $\theta_1, \dots, \theta_d$ .

— Need marginal posterior for  $\theta_1$ .

$$P(\theta_1|X^{(n)}) = \int \int \dots \int p(\theta_1, \dots, \theta_d | X^{(n)}) d\theta_2 \dots d\theta_d.$$

Bayesian estimators.

(1) mean of marg post

$$\int_{\theta_1} \theta_1 P(\theta_1 | X^{(n)}) d\theta_1$$

(2) median of marg post

$$\hat{\theta}_{\text{median}}: \int_{-\infty}^{\hat{\theta}_{\text{median}}} P(\theta_1 | X^{(n)}) d\theta_1 = 0.5$$

$$(3) E(\theta_1 | X^{(n)}) = \int \theta_1 P(\theta_1 | X^{(n)}) d\theta_1$$

$$\left( \begin{array}{c} \vdots \\ \int \theta_d \dots d\theta_d \end{array} \right)$$

Ex. Normal model

$$X_i \sim Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

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

$$L(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i-\mu)^2}{2\sigma^2}\right\}$$

$$\ell(\mu, \sigma^2) = -\sum_{i=1}^n \frac{(x_i-\mu)^2}{2\sigma^2}$$

$$\underbrace{\ell(\mu, \sigma^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2}$$

likelihood equations

$$\frac{\partial \ell(\mu, \sigma^2)}{\partial \mu} = -\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$\frac{\partial \ell(\mu, \sigma^2)}{\partial \sigma^2} = \frac{1}{-2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 - \frac{n}{2\sigma^2} = 0$$

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

The solution  $\hat{\theta} = (\hat{\mu}, \hat{\sigma}^2)$

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$E(\hat{\sigma}^2) \neq \sigma^2 = \frac{n}{n-1} \sigma^2$$

Joint position

$$P(\mu, \sigma^2 | X^{(n)}) \propto P(\mu, \sigma^2) L(\mu, \sigma^2)$$

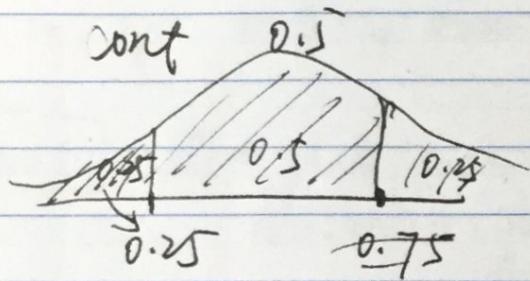
# MATH 459

$$\hat{\theta} = g(X^{(n)})$$

$P(\theta) = 1/\theta \rightarrow \text{improper prior}$  (not a density)

$\downarrow$

may still be possible posterior can be normalized.



$$z_n = \frac{n\bar{X} - n\mu}{\sigma\sqrt{n}} = \frac{n(\bar{X} - E(\bar{X}))}{\sqrt{n} \cdot \sqrt{\text{Var}(\bar{X})}}$$

Sampled

$$\frac{\sqrt{n} \cdot (\bar{X} - E(\bar{X}))}{\sqrt{\text{Var}(\bar{X})}} \xrightarrow{\text{d}} N(0, 1)$$

Sampled

$$\Rightarrow \frac{(\bar{X} - E(\bar{X}))}{\sqrt{\text{Var}(\bar{X})}} \xrightarrow{\text{d}} \frac{[\bar{X} - E(\bar{X})] \cdot \sqrt{n}}{\sqrt{\text{Var}(\bar{X})}}$$

$$\Rightarrow \frac{\bar{X} - M}{S/\sqrt{n}}$$

$$E[\bar{X}] = E[X]$$

$$n \text{Var} \bar{X} = \text{Var} X$$

What is a good decision rule?

Math  
Topic

Ideal case:  $\exists d \in D$  st.

$R(\theta, d)$  is uniformly smallest  $\forall \theta \in \Theta$

Not practical.

Admissibility: Given two decision rules  $d, d'$   $d$  strictly dominates  $d'$  if  $R(\theta, d) \leq R(\theta, d')$

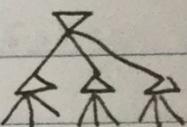
$\forall \theta \in \Theta$  and is strictly less. for at least one  $\theta$

Any rule which is strictly is strictly dominated is inadmissible.

If  $d$  is not strictly dominated. it is admissible

Minimaxity: The max risk of  $d$  is

$$MR(d) = \sup_{\theta \in \Theta} R(\theta, d)$$



$d$  is minimax if  $MR(d) \leq MR(d') \quad \forall d' \in D$

Equivalently,  $d$  must satisfy  $\sup_{\theta \in \Theta} R(\theta, d) = \inf_{d' \in D} \sup_{\theta \in \Theta} R(\theta, d')$

If The supremum and infimum are actually attained.

$$\max_{\theta \in \Theta} R(\theta, d) = \min_{d' \in D} \max_{\theta \in \Theta} R(\theta, d')$$

Unbiasedness

$d$  is unbiased if  $E_{\theta} [L(\theta', d(x))] \geq E_{\theta_{\text{true}}} [L(\theta, d(x))]$

$d$  is able to find out the optimal  $\theta$  from other  $\theta'$

for  $\forall \theta, \theta' \in \Theta$

Recall, in estimator:

$d(x)$  is unbiased for  $\theta$ , if  $E_{\theta} d(x) = \theta$

If  $L(\theta, d) = (\theta - d)^2$  these definition are the same

## Bayes' decision rule

Problem has a seventh component  
corresponding  $\rightarrow$  prior.

Bayes risk of  $d$

$$r(p, d) = \int_{\theta \in \Theta} \underbrace{R(\theta, d)}_{\text{average over data}} p(\theta) d\theta$$

Bayes decision rule w.r.t  $p(\theta)$  minimizes the Bayes risk, s.t.

$$r(p, d) = \inf_{d' \in D} \frac{r(p, d')}{\text{given } d}$$

Why people like this?

1. always (almost) admissible.
2. every admissible rule is either Bayes or a limit of Bayes rules.

# Statistical Decision Problem

① Parameter space

$$\Theta \subseteq \mathbb{R}^d$$

Set of possible true states of nature

② Sample space  $\mathbb{X}$

where data live

typically  $n$  observations

generic element  $x \in \mathbb{X}$  is  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$

③ family of prob. dist on  $\mathbb{X}$

indexed by values  $\theta \in \Theta$

$$\{P_\theta(x) : x \in \mathbb{X}, \theta \in \Theta\}$$

④ Action Space  $A$

Set of all actions available to the experimenter

(a) Hypothesis testing

decide b/w  $H_0, H_1$

$$A = \{a_0, a_1\} \rightarrow \text{reject } H_0$$

Fail to reject

(b) Estimation

estimate  $\theta$  by a func of  $X^{(n)}$ , e.g.

$$\bar{x} = n^{-1} \sum_{i=1}^n x_i, \text{ or } x_1^3 + x_2 \cos \left\{ \sqrt{n} x_3 \right\}$$

⑤ Loss function

$$L : \Theta \times A \mapsto \mathbb{R}$$

links action to unknown Parameter. If we take action

$a$  when true state of nature is  $\theta \in \Theta$ , incur loss  $L(\theta, a)$

⑥ Set of des. decision rule  $D$

An element  $d \in D$ ,  $d : \mathbb{X} \mapsto A$

s.t each  $x \in \mathbb{X}$  is assoc. with  $d(x) \in A$

EX: Hypothesis testing

Adopt rule. Fail to reject  $H_0$ .

if  $\bar{x} \leq 3.6$ , otherwise reject  $H_0$ .

$$d(x) = \begin{cases} a_0 & \text{if } \bar{x} \leq 3.6 \\ a_1 & \text{if } \bar{x} > 3.6 \end{cases}$$

(\*)

### Risk function

For  $\theta \in \Theta$  the risk of decision  $d$  based on random  $X$

$R(\theta, d) = E_{\theta}[L(\theta, d(x))] \leftarrow \text{given } \theta \text{ average over data. [given } \theta\text{]}$

$$= \int_X L(\theta, d(x)) f(x; \theta) dx \quad \begin{matrix} \rightarrow \text{frequentist view} \\ \text{over } \theta \text{ (fp)} \end{matrix}$$

$$\sum_X L(\theta, d(x)) f(x; \theta) \quad \text{Bayesian avg (RL)} \xrightarrow{H(p,d)} E\left(E_{\theta}[L(\theta, d(x))]\right) \quad \text{over } X$$

### Key notion

Different decision rules should be composed in terms of their risk as function  $\theta$

$\partial d \times \partial \theta$

Common loss Functions for estimation.

(a)  $L(\theta, a) = (\theta - a)^2$

(b)  $L(\theta, a) = |\theta - a|$  absolute error

(c)  $L(\theta, a) = \begin{cases} 0 & \text{if } |\theta - a| \leq \delta \\ 1 & \text{if } |\theta - a| > \delta \end{cases}$

mean square error

Common loss func in Testing.

$$L(\theta, a) = \begin{cases} 1 & \text{if } \theta \in H_0 \text{ and } a = a_1 \\ 1 & \text{if } \theta \in H_1 \text{ and } a = a_0 \\ 0 & \text{otherwise} \end{cases}$$

Risk associated with these losses.

Type I and Type II error

Hence  $R(\theta, d) = \begin{cases} \Pr_{\theta}\{d(X) = a_1\} & \text{if } \theta \in H_0 \\ \Pr_{\theta}\{d(X) = a_0\} & \text{if } \theta \in H_1 \end{cases}$

## Bayes Risk

$$R(P_d, d) = \int_{\theta \in \Theta} R(\theta, d) P(\theta) d\theta$$

↑  
given  $\theta$

## Bayes decision rule

$$r(P, d) = \inf_{d' \in D} r(P, d')$$

Bayes rule of  $d$  satisfies this

Why Bayes is good.

① almost all Bayes rules admissible

② every admissible rule is either Bayes or a limit of Bayes rules

③ a decision rule  $d$  is minimax if (a) it is a Bayes for some prior  $P$

$$(b) \max_{\theta} R(\theta, d) \leq R(P_d, d)$$

$P$ : prior, s.t.  $d$  is minimax  
minimize the worst case cost

type 1 error

type 2 error

## Bayes Rules

Risk function:

$$R(\theta, d) = E_{\theta} L(\theta, d)$$

$$= \int_X L(\theta, d(x)) f(x; \theta) dx$$

## Bayes Risk

$$R(P, d) = \int_{\theta} R(\theta, d) P(\theta) d\theta$$

$$= \int_{\theta} \int_X L(\theta, d(x)) f(x; \theta) P(\theta) dx d\theta$$

$$= \int_{\theta} \int_X L(\theta, d(x)) f(x; \theta) P(\theta | x) dx d\theta$$

$$\text{and } (f(x) = \int f(x|\theta) P(\theta) d\theta)$$

$f(x; \theta) \rightarrow$  for frequentist concept

$P(\theta | x) \rightarrow$  get Bayesian posterior using frequentist concept

key for Bayesian aspect

$$= \int_X f(x) \underbrace{\left\{ \int_{\theta} L(\theta, d(x)) P(\theta | x) d\theta \right\}}_{\text{w.r.t. } d} dx$$

To minimize  $R(P, d)$  sufficient to minimize

$$\int_{\theta} L(\theta, d(x)) P(\theta | x) d\theta$$

"expected posterior loss"



## Bayes Point Estimators

(Let  $L(\theta, d) = (\theta - d)^2$ )

for observed  $X = x$

choose  $d(x)$  to minimize

$$\int_{\theta} (\theta - d)^2 P(\theta | x) d\theta$$

Differentiate w.r.t  $d$

$$= \int_{\theta} (\theta - d) P(\theta | x) d\theta = 0$$

Hint  $\int_{\theta} P(\theta | x) d\theta = 1$

$$\Rightarrow \int_{\theta} \theta P(\theta | d) d\theta = \int_{\theta} d P(\theta | x) d\theta$$

||      ↓      ||

$$\boxed{\text{Posterior mean} = d}$$

more generally if  $L(\theta, d) = (\theta - d)^2$

consider estimator

$$R(\theta, d) = \text{Mean squared error.}$$

(Now suppose  $L(\theta, d) = |\theta - d|$ )

Bayes rule minimizes

$$\int_{-\infty}^d (d - \theta) P(\theta | x) d\theta + \int_d^{\infty} (\theta - d) P(\theta | x) d\theta$$

$$= \int_{\theta} |\theta - d| P(\theta | x) d\theta$$

Differentiate w.r.t  $d$ , and set to 0

$$\Rightarrow \int_{-\infty}^d P(\theta | x) \frac{d\theta}{d} = \int_d^{\infty} P(\theta | x) d\theta$$

~~∴~~ [when  $d$  is the posterior median.]

$$\int_{-\infty}^{\infty} L(\theta, \delta) f_{\theta|x}(v|x) dv = \int_{-\infty}^{\delta-\frac{\delta}{2}} 1 f_{\theta|x}(v|x) dv + \int_{\delta+\frac{\delta}{2}}^{\delta} 1 f_{\theta|x}(v|x) dv$$

Bayes interval estimator

Suppose  $L(\theta, d) = \begin{cases} 0 & \text{if } |\theta - d| \leq \delta \\ 1 & \text{if } |\theta - d| > \delta \end{cases}$

$\delta$  is a tolerance level.

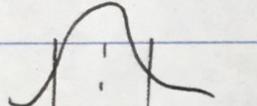
Problem: Find the 'best' interval

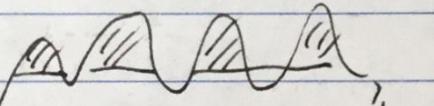
with ~~respect~~ prespecified length  $2\delta$

$$(d - \delta, d + \delta) \Rightarrow \text{Loss} = 0$$

'Best' means the interval maximizes  
the posterior probability that  
 $\theta \in (d - \delta, d + \delta)$

HPDI (Highest posterior density interval)

If   
easy to construct

If tolerance is   
HPP set union of disjoint  
interval

simplify problem.

# Lecture 10:

## Background: ① Prior Independence

If  $P(\theta) = P(\theta_1)P(\theta_2)\cdots P(\theta_d)$

Does not imply a posteriori independence  
also need  $L(\theta) = L(\theta_1)L(\theta_2)\cdots L(\theta_d)$

## ② Fisher information

$$I(\theta) = \text{Var}(D_2) = -E(D_2)$$

here:  $I(\theta) = E_{\theta}[\nabla_{\theta} \log f(x; \theta) \nabla_{\theta} \log f(x; \theta)^T]$

$$\begin{aligned} I_{ij} &= E_{\theta}\left[\left(\frac{\partial}{\partial \theta_i} \log f(x; \theta)\right)\left(\frac{\partial}{\partial \theta_j} \log f(x; \theta)\right)^T\right] \\ &= -E_{\theta}\left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(x; \theta)\right] \end{aligned}$$

## ③ Orthogonality

$\theta_i, \theta_j$  are orthogonal if  $I_{ij} = 0$

## Select Prior by Formal Rules

### Noninformative Prior & Objective prior

"maximize divergence between pri & post"  
afunc that Jose M. Bernardo  
use flat prior (not invariant)  $\Rightarrow$  variance is large?

## Subjective, Bayesian

- ① choose a prior
- ② apply it to the likelihood  
to derive the posterior

$P(\eta) \geq 1$  (noninformative)

$P(\eta) | \eta = \log \frac{\theta}{1-\theta}$  is informative  
and also not invariant under  
one-to-one reparameterization

•  $I(\eta) = I(\theta) \cdot \left(\frac{d\theta}{d\eta}\right)^2$

Invariant  $\sqrt{I(\eta)} = \sqrt{I(\theta)} \cdot \left|\frac{d\theta}{d\eta}\right|$   
by change-of-variable formula,

this shows Jeffreys prior

$\pi_j(\theta) = \sqrt{I(\theta)}$  is invariant to  
a change of variable:  $\eta \leftrightarrow \theta$

ex:  $X \sim \exp(\lambda)$

$$I(\theta) = -E[D_2] = \lambda^{-2}$$

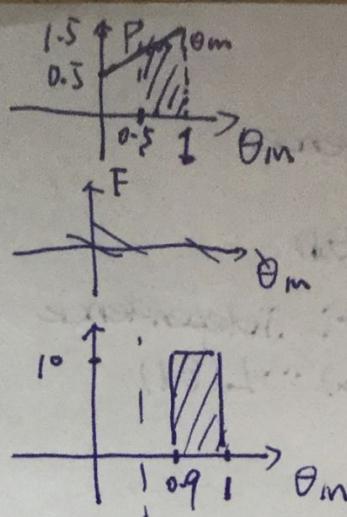
$$\pi_j(\lambda) \propto \frac{1}{\sqrt{I(\lambda)}} = \lambda^{-1}$$

$$P_m(\theta_m) = \theta_m + 0.5, \theta_m \in (0, 1)$$

$$P_{\text{rm}}(\theta_m > 0.5) = \int_{0.5}^1 P_m(\theta) d\theta_m = 0.625$$

$$P_N(\theta_s) = 10, \theta \in (0.9, 1)$$

$$P_{\text{rm}}(\theta_m > 0.5) = \int_{0.5}^1 P_m(\theta) d\theta_m = 0.1$$



prior dist

posterior dist

posterior confidence  
interval

$$P(\theta_m | y_m = 0) \propto P(\theta_m) \cdot P(y_m | \theta_m)$$

$$\propto \frac{\theta_m}{2} - \theta_m^2 + \frac{1}{2}$$

$$P_{\theta_m | Y_m}(\theta_m > 0.5) = \int_{0.5}^1 P_m(\theta_m | y_m) d\theta_m = 0.35$$

give up after one try

prior allow us to incorporate common sense

Lecture 9

## Lecture 12

Approximate Bayesian Inference = Exact ~~analytic~~  
often not practical  
principle

### Basics of Parametric Bayesian Asymptotics

#### Consistency of Posterior

$$\pi(f(x|\theta_0)) \Rightarrow \pi(\cdot | X^{(n)}) \rightarrow 0$$

if  $U$  near  $\theta_0 \in \Theta$

$$\pi(U | X^{(n)}) \rightarrow 1 \text{ as } n \rightarrow \infty$$

Some of the principle

violated by frequentist  
procedures:

① admissibility

② unbiasedness

③ probability matching

#### Asymptotic Normality of MLE: 4法

and  $\frac{\theta - \hat{\theta}_n}{\sqrt{n}} \rightarrow N(0, I^{-1}(\theta_0))$

MLE  $\rightarrow$  MVN

#### (S.S.) sufficiency principle

$$\theta = (\mu, \sigma^2) \quad T(x) = (\bar{x}, s^2)$$

- ① inference on  $\theta$  should be based  
only on  $T(x)$

completeness: unique unbiased S.S  
MLE

# Lecture 13

## Approximate Bayesian Inference II

numerical integration

$$\text{inverse} \Rightarrow \text{elementary function}$$

$$P(X \leq f(x)) \quad \text{key}$$

$$P(X < x) = P(X \leq F^{-1}(u))$$

$$= F(x) \quad X = F_x^{-1}(u)$$

some special elementary function

Integration Approximation (Type

(1) asymptotic expansion

(2) deterministic numerical approximation

(quadrature method) converge order  $O(n^{-2})$  not good for high dim

(3) Monte Carlo integration (utilizing randomness) conv order  $O(n^{-1/2})$  large sample for high dimension

MC approx for mean:  $\hat{\mu} = E[h(x)]$

$$\hat{\mu}_{mc} = \frac{1}{n} \sum_{i=1}^n h(x_i)$$

$$\hat{\sigma}^2_{(f)mc} = \frac{1}{n-1} \sum_{i=1}^n [h(x_i) - \hat{\mu}_{mc}]^2$$

Magic of MH, for every given  $q$ , we can construct a MH kernel  $k$  so that  $f$  is its stationary distribution. (more easier than A-R method)

Comment:

(1) closely relate to A-R method (Rejection Sampling)  
but the proposal dist changes over time

(2) the acceptance rate is the average of all iterations

$$\bar{p} = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T p(x^{(t)}, y_t)$$

$$= \int p(x, y) f(x) q(y|x) dy dx$$

(3) independent of normalizing constants ( $M$ )

(4) in the symmetric case,  $\bar{p}$  can be reduced to a ratio independent of  $q$ . ex.  $q(x|y) = q(y|x)$

Detailed Balance:  $f(x) k(y|x) = f(y) k(x,y)$  exist a  $f(x)$

Independent MH  
genetic MH:  $q$  depend on current state of the chain.  
If we require  $q$  independent of  $\dots$ , ie.  $q(y|x) = g(y)$   
we have special case

Given initial  $x^{(t)}$

1. Generate  $y_t \sim g(y)$

2. Take

$$x^{(t+1)} = \begin{cases} y_t & p = \min\left(\frac{f(y_t)}{f(x^{(t)})}, 1\right) \\ x^{(t)} & \text{otherwise. } \frac{g(x^{(t)})}{g(y_t)}, 1 \end{cases}$$

generalization of A-R  
Sampling

# Comparison of Rejection with M-H

- ① AR is sampled is iid, MH depend on  $X^{(t)}$  (not matter how, always)
- ② MH involves repeated occurrence of the same value Reject  $Y_t$   
 $\Rightarrow$  repeat  $X_t$  at time  $t+1$
- ③ no  $M$  needed

## Select of proposal.

- ① Random walk M-H (Generic approach)  
 the proposal is random walk,  
 but the MC is not (A-R involved)

(RWMH) | (long chain, half time revisiting)

Given initial  $X^{(t)}$

1. Gen  $Y_t \sim g(y - X^{(t)})$

2. Take  $X^{t+1} = \begin{cases} Y_t & U_p^t = \min\{1, f(Y_t)/f(X^{(t)})\} \\ x_t & \text{otherwise} \end{cases}$

$\bar{p}$  does not depend on  $g$ , but different choice of  $g \Rightarrow$  different ranges for  $Y_t$  and  $\Rightarrow$  diff  $\bar{p}$

- ② Langevin algorithm
- $$Y_t = X^{(t)} + \frac{\delta^2}{2} \nabla \log f(x^{(t)}) + \delta \epsilon_t,$$
- $$\epsilon_t \sim g(t)$$
- (require a lot of tuning in practice)

# Lecture 14

MCMC method

source of randomness

① physical process

is not fast

need storage

② Pseudo-random number generators

~~Linear Congruential Generators~~  
multiplicative

LCG/MCG

multiple recursive generator  
MRG

Mersenne twister RNG, default in R

Definition, for any continuous CDF F

$U \sim \text{Unif}(0,1)$ , then

$$X = F^{-1}(U)$$

proof.  $P(X < x) = P(X < F^{-1}(U))$

$$= P(U < F(x))$$

$$= F(x)$$

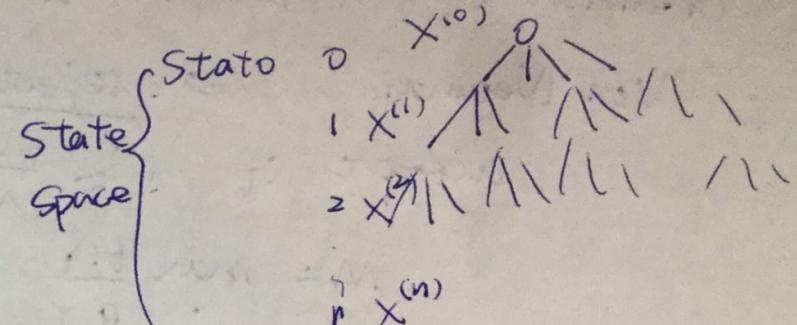
Ergodic Theorem

problem no high dimension  
high cost

revolutionary idea: Markov Chain

Monte Carlo Markov Chain

MCMC: a random process



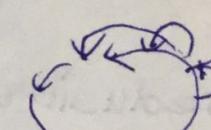
$$X^{(t+1)} | X^{(0)}, X^{(1)}, X^{(2)}, \dots, X^t$$

$$\sim k(x^{(t)}, \underline{x^{(t+1)}})$$

• Stationary: if  $x^{(t)} \sim f$   
then  $x^{(t+1)} \sim f$

$$\int_X k(x, y) f(x) = f(y)$$

• Recurrent



the stationary dist is also a limiting dist which means the limiting dist of  $x^{(t)}$  is  $f$  for almost all initial  $x^{(0)}$   $\Rightarrow$  ergodicity  
 $\Rightarrow$  eventually produce  $f$

# lecture 15

Metropolis-Hastings, Gibbs sampling  
reversible jump, convergence diagnostics

Generating iid samples accept-reject method

Von Neumann's idea: Rejection Sampling

accept if  $U \leq \frac{f(y)}{Mg(y)}$

$$M = \max_y \frac{f(y)}{g(y)}$$

Simple Random Walk

$$\begin{aligned} x^{(t+1)} &= x^{(t)} + t_t \\ k(x^{(t)}, x^{(t+1)}) &\equiv N(x^{(t)}, 1) \\ f(x^{(t+1)} | x^{(t)}) & \end{aligned}$$

irreducible: regardless of  $x^{(0)}$   
positive prob  $\rightarrow$  reach any part  
of the state space

an irreducible MC  $\rightarrow$  is current if  
return to arbitrary part of the state infinite time

The stationary of ~~of~~ is also a limiting dist  $\Rightarrow$  ergodicity theorem

Basic idea of MCMC

Given target dist ~~of~~  $f$

①  $\Rightarrow$  build a Markov kernel

$$K(x, y) = P(X^{(t+1)} = y | X^{(t)} = x)$$

with stationary dist  $\pi$  the same  
as target  $f$

②  $\rightarrow$  generate a MC  $\{X^{(t)}\}$

using the kernel  $K$  so that  
the limiting dist of  $\{X^{(t)}\}$  is  $f$   
and  $T^{-1} \sum_{t=1}^T h(x^{(t)}) \rightarrow E_f[h(x)]$

Idea of Metropolis-Hastings

Given initial  $x^{(t)}$ , and candidate  
density  $q(y|x^{(t)})$   $\rightarrow$  candidate dist

1. Generate  $y_t \sim q(y|x^{(t)})$

2. Set  $x^{(t+1)} = \begin{cases} y_t & \text{with Prob} \\ p(x^{(t)}, y_t) \\ (R) x^{(t)} & \text{otherwise} \end{cases}$

$$\text{where } p(x, y) = \min \left\{ \frac{f(y) q(x|y)}{f(x) q(y|x)}, 1 \right\}$$

$$\frac{1}{T} \sum_{t=1}^T h(x^{(t)}) \rightarrow E_f[h(x)]$$

Gibbs sampling

$$\theta = (\theta_1, \dots, \theta_d)^T \text{ (Posterior of } \theta)$$

$$p(\theta_j | \theta_{(-j)}^{t-1}, y) \quad \begin{matrix} \uparrow \\ \text{expect } j \end{matrix}$$

$$\theta_{(-j)}^t = (\theta_1^t, \dots, \theta_{j-1}^t, \theta_{j+1}^{t-1}, \dots, \theta_d^{t-1})^T$$

ex. from BDA

~~say~~  $(y_1, y_2)$  (obsver)  
 $\theta = (\theta_1, \theta_2)^T$ ,  $\text{cov} \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$

$$\Rightarrow \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} | y \sim N \left( \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right) \text{ (easier)}$$

$$\Rightarrow \theta_1 | \theta_2, y \sim N(y_1 + \rho(\theta_2 - y_2), 1 - \rho^2) \text{ (Gibbs Sample, time consuming)}$$
$$\theta_2 | \theta_1, y \sim N(y_2 + \rho(\theta_1 - y_1), 1 - \rho^2)$$

Mixing of Chain

Some application tuning method

R package

MCMCpack for MCMC

tstan  $\xrightarrow{\text{output}}$  coda mcmc  
object  
of boa

Mixing of Chain: Mixing time

$$\text{Auto correlation: } P_k = \frac{\sum_{i=1}^{n-k} (X_i - \bar{X})(X_{i+k} - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

expect the  $k$ -th to be  $\downarrow$  as  $k \uparrow$

Any Markov chain will have autocorrelation current value depend on the previous one

$$\Rightarrow \beta | \delta, X, Y \sim N(\hat{\beta}, (\mathbf{X}^T \mathbf{X})^{-1} \delta^2)$$

Lec 17

Linear Regression model

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_{k-1} X_{ik-1} + \varepsilon_i$$

$$Y = X\beta + \varepsilon$$

$$E(Y|X) = \beta_0 + \sum_{j=1}^{k-1} \beta_j X_j + E(\varepsilon|X), \quad \text{Var}(\varepsilon|X) = \delta^2 I_{n \times n}$$

$$\text{Least square estimator: } \min_{\beta} (Y - X\beta)^T (Y - X\beta)$$

$$\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T Y, \quad \delta^2 = \frac{(Y - X\hat{\beta})^T (Y - X\hat{\beta})}{n - k}$$

provided exists.

$$\text{maximize } \underbrace{N(0, \delta^2 I_{n \times n})}_{\text{Ls problem.}} \Rightarrow \text{Ls problem.}$$

$$\begin{aligned} f(x, y | \psi, \theta) & , L(\psi, \theta) \\ \text{prior } \psi, \theta & , \end{aligned}$$

$$\Rightarrow P(\beta, \delta^2 | X, Y) = P(\beta | \delta^2, X, Y) P(\delta^2 | X, Y)$$

start by find the posterior of  $\beta$ , conditional on  $\delta$ , then find marginal distribution of  $\delta^2$ .

$$P(\beta | \delta^2, X, Y) = \frac{P(\beta, \delta^2 | X, Y)}{P(\beta | \delta^2, X, Y)}$$

$$\Rightarrow \delta^2 \sim \text{Inv-}\chi^2(n - k, s^2)$$

$$s^2 = \frac{1}{n - k} (Y - X\hat{\beta})^T (Y - X\hat{\beta})$$

In practice, we use MCMC by drawing from  $\delta$  and then drawing  $\beta | \delta$

The joint posterior  $P(\beta, \delta^2 | X, Y)$  is

proper provided that  $(X^T X)$  is invertible

$$\Rightarrow \begin{aligned} 1. \quad n & \geq k \\ 2. \quad \text{rank}(X) & = k \end{aligned}$$

Posterior Predictive Density

given:  $\tilde{y} | \beta, \sigma^2, \underline{x}^* \sim N(\underline{x}^* \beta, \sigma^2)$

$$P(\tilde{y} | \underline{y}_{\text{train}}) = \int P(\tilde{y} | \beta, \sigma^2) P(\beta, \sigma^2 | \underline{y}) d\beta d\sigma^2$$

What MCMC regression does.. Bayesian Fit  
Simulates from posterior using  
Gibbs Sampling

Sample from  $P(\beta, \sigma^2 | X, Y)$

1. Compute  $\hat{\beta}$  and  $(X^T X)^{-1}$
2. Compute  $S^2$
3. Draw  $\sigma^2$  from inverse  $\chi^2$  dist
4. Draw  $\beta$  from multivariate normal dist.

If we know  $P(\sigma^2 | \beta, X, Y)$  and  $P(\beta | \sigma^2, X, Y)$

we can sample  $P(\beta, \sigma^2 | X, Y)$  using Gibbs sampling

Initialize  $\beta^{(1)}, \sigma^{2(1)}$

for  $t = 1:T$

$$\beta_{t+1} \sim P(\beta^{(t)} | t^{(t)}, X, Y)$$

$$\sigma^{2(t+1)} \sim P(\sigma^2 | \beta_{t+1}, X, Y)$$

END.

The puffin data from Learn Bayes package

(lm)  $\rightarrow$  Frequentist Fit

('blinReg')  $\rightarrow$  Bayesian Fit