

# Lecture 7: Point estimation

Ma Xuejun

School of Mathematical Sciences

Soochow University

<https://xuejunma.github.io>



# Outline

- 1 Estimator
- 2 MoM
- 3 MLE
- 4 Bayes Estimators
- 5 EM Algorithm
- 6 Methods of Evaluating Estimators
  - Mean Square Error
  - Best Unbiased Estimators

# Review

- Population and Sample
- Popular models: Normal model; exponential family; Bayesian
- statistics



# Estimators

## Estimator

let  $X_1, \dots, X_n \sim f(x; \theta)$ . An estimator

$$\hat{\theta} = \hat{\theta}_\theta = W(X_1, \dots, X_n)$$

is a function of the data.

- An estimator is a statistic. It is a **random variable**.
- When we have an observation  $x_1, x_2, \dots, x_n$ , the corresponding result  $W(x_1, x_2, \dots, x_n)$  is called an estimate of  $\theta$ .
- Note: estimator  $\hat{\theta}$  is a random variable, and estimate is a realization of this random variable.
- Multiple estimator for  $\theta$ .

# Method of Moments(MoM)

Suppose the unknown parameter  $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ . Define a sequence of the moment functions

Data moment

$$m_1 = \frac{1}{n} \sum_{i=1}^n X_i$$

$$m_2 = \frac{1}{n} \sum_{i=1}^n X_i^2$$

$$m_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

Theoretical moment

$$\mu_1(\theta) = E[X_i]$$

$$\mu_2(\theta) = E[X_i^2]$$

$$\vdots \quad \quad \quad \vdots$$

$$\mu_k(\theta) = E[X_i^k]$$

Let the estimator  $\hat{\theta}$  satisfies that

$$m_j = \mu_j(\hat{\theta}), j = 1, \dots, k.$$

Then there are  $k$  equations. Solve the  $k$  equations to get  $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$ , where each of them is a function of  $m_1, m_2, \dots, m_k$ , the moments from data.

# MOM:Example

**Example 7.2.1** Let  $X_i \sim N(\mu, \sigma^2)$ . Find the estimators of  $\mu$  and  $\sigma^2$  with MOM method.

Note that the theoretical moments are

$$E[X] = \mu, \quad E[X^2] = \mu^2 + \sigma^2.$$

The MOM estimator should satisfy that

$$\frac{1}{n} \sum_{i=1}^n X_i = \hat{\mu}, \quad \sum_{i=1}^n X_i^2 = \hat{\mu}^2 + \hat{\sigma}^2.$$

Therefore, the estimators are

$$\hat{\mu} = \bar{X}_n \quad \hat{\sigma}^2 = \sum_{i=1}^n X_i^2 - \bar{X}_n^2 = \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

## MOM:Example 2

**Example 7.2.2** Let  $X_i \sim \text{Binomial}(k, p)$  random variable with  $\theta = (k, p)$ . Find the estimators of  $\mu$  and  $p$  with MOM. Note that the theoretical moments are

$$E[X] = kp, \quad E[X^2] = (kp)^2 + kp(1 - p).$$

The MOM estimator should satisfy that

$$\frac{1}{n} \sum_{i=1}^n X_i = \hat{k}\hat{p}, \quad \frac{1}{n} \sum_{i=1}^n X_i^2 = (\hat{k}\hat{p})^2 + \hat{k}\hat{p}(1 - \hat{p}).$$

Therefore, the estimators are

$$\hat{p} = \bar{X}_n / \hat{k} \quad \hat{k} = \frac{\bar{X}_n^2}{\bar{X}_n^2 - \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2}$$

- Easy to calculate
- May not fit the parameter set ( $\hat{k}$  may not be an integer here)

# Maximum Likelihood Estimator(MLE)

- Consider a sample  $X_1, \dots, X_n \sim f(x; \theta), \theta \in \Theta$
- Observe  $x_1, \dots, x_n$
- Which  $\theta$  maximizes the possibility that we have this observation?

$$\hat{\theta} = \arg \max_{\theta \in \Theta} f_{\theta}(x_1, \dots, x_n)$$

- Similar as what did for I.I.D.Normal Model, we find the joint density:

$$L(\theta; x_1, \dots, x_n) = f(x_1, \dots, x_n; \theta)$$

which is also a function of  $\theta$ . We call  $L(\theta)$  as the likelihood function.

- The estimate is called Maximum Likelihood Estimate.



### Definition 7.2.4

For each sample point  $x$ , let  $\hat{\theta}(x)$  be a parameter value at which  $L(\theta|x)$  attains its maximum as a function of  $\theta$ , with  $x$  held fixed. Then  $\hat{\theta}(x)$  is a maximum likelihood estimator of  $\theta$ .

- **Remark**

Sometime maximizing  $l(\theta|x) = \log(L(\theta|x))$  is much easier than maximizing  $L(\theta|x)$ .

**Remark**  $X_1, \dots, X_n$  in most cases do not have to be **identically distributed**.

# MLE: Procedure

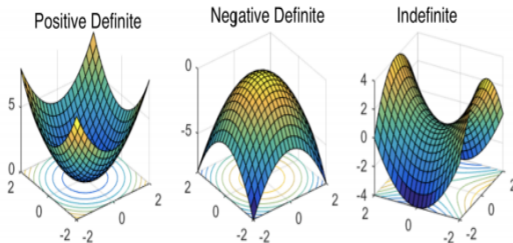
The procedure to figure out the maximiser of  $L(\theta)$ ;

- Find the likelihood function  $L(\theta)$  (actually, the joint density function)
- Find the log-likelihood function  $l(\theta) = \log L(\theta)$
- Compute the gradient vector of  $l(\theta)$  with respect to  $\theta$ , denoted by  $\nabla l(\theta)(x_1, \dots, x_n; \theta)$ . For univariate case, it is the gradient only.
- Solve  $\nabla l(\theta)(x_1, \dots, x_n; \theta) = 0$ , with respect to  $\theta \in \Theta$ , call this solution  $\tilde{\theta}_n$ , check whether  $H(\tilde{\theta}_n)$  is negative definite. If yes, then  $\tilde{\theta}_n = \hat{\theta}_n$ . Here,  $H(\theta)$  is the Hessian matrix for  $l(\theta)$ , defined as

$$\begin{pmatrix} \frac{\partial(\theta)}{\partial \theta_1^2} & \frac{\partial(\theta)}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial(\theta)}{\partial \theta_1 \partial \theta_k} \\ \frac{\partial(\theta)}{\partial \theta_2 \partial \theta_1} & \frac{\partial(\theta)}{\partial \theta_2^2} & \cdots & \frac{\partial(\theta)}{\partial \theta_2 \partial \theta_k} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial(\theta)}{\partial \theta_k \partial \theta_1} & \frac{\partial(\theta)}{\partial \theta_k \partial \theta_2} & \cdots & \frac{\partial(\theta)}{\partial \theta_k^2} \end{pmatrix}$$

# Hessian Matrix

- When  $k = 1$ , the Hessian Matrix is the second derivative of  $l(\theta)$ .
- The matrix can be seen as the second derivative for a multivariate function
- Negative definite indicates **local maximum**
- Since there is only one solution, the local maximum is also the global maximum (if  $\Theta$  is properly defined)



# MLE:Remarks

- It may be hard to figure out the solution for  $\nabla l(\theta) = 0$  analytically. Computationally, some methods can be applied, such as Newton method. Details in applied mathematics / computing mathematics.
- The general Hessian matrix  $H(\theta)$  may be complicated, yet  $H(\tilde{\theta}_n)$  is usually easier since  $\tilde{\theta}_n$  satisfies  $\nabla l(\tilde{\theta}_n) = 0$ .
- If there are multiple solutions, then we should check the values to figure out the maximal one, check the Hessian matrix at this maximal solution, and also check  $l(\theta)$  when some  $\theta$  goes to infinity.
- If  $\Theta$  is not properly defined, it is possible that the maximal is achieved at the bound of  $\Theta$ , not the maximum point. Or, the solution  $\tilde{\theta}_n$  does not exist in  $\Theta$ .

# MLE: Example

**Example** Let  $X_1, X_2, \dots, X_n \sim \text{Poisson}(\lambda)$ . Find the MLE of  $\lambda$ , given observations  $x_1, \dots, x_n$ .

The likelihood function is

$$L(\lambda) = f(x_1, \dots, x_n; \lambda) = \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} = e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i} \frac{1}{\prod_{i=1}^n x_i!}$$

The log-likelihood function is

$$l(\lambda) = \log L(\lambda) = -n\lambda + \left(\sum_{i=1}^n x_i\right) \log \lambda - \sum_{i=1}^n \log x_i!$$

The gradient vector is a derivative:

$$\frac{d}{d\lambda} l(\lambda) = -n + \frac{1}{\lambda} \left(\sum_{i=1}^n x_i\right)$$

Set it to be 0, and solve the equation, we have  $\tilde{\lambda}_n = \bar{x}_n$ . Note that

$$\frac{d^2}{d\lambda^2} l(\lambda) \Big|_{\lambda=\tilde{\lambda}} = -\frac{1}{\lambda^2} \left(\sum_{i=1}^n x_i\right) \Big|_{\lambda=\tilde{\lambda}} = -n/\bar{X}_n < 0$$

So, the MLE is  $\hat{\lambda}_n = \bar{X}_n$ .

## MLE: Example 2

**Example** Let  $X_1, X_2, \dots, X_n \sim \text{Exp}(\lambda)$ , Find the MLE of  $\lambda$ , given observations  $x_1, x_2, \dots, x_n$ .  
The likelihood function is

$$L(\lambda) = f(x_1, \dots, x_n; \lambda) = \prod_{i=1}^n \lambda e^{-\lambda x_i} = e^{-\lambda \sum_{i=1}^n x_i} \lambda^n$$

The log-likelihood function is

$$l(\lambda) = \log L(\lambda) = n \log \lambda - \lambda \left( \sum_{i=1}^n x_i \right)$$

The gradient vector is a derivative:

$$\frac{d}{d\lambda} l(\lambda) = n/\lambda - \sum_{i=1}^n x_i.$$

Set it to be 0, and solve the equation, we have  $\tilde{\lambda}_n = 1/\bar{x}_n$ . Check the second derivative at  $\tilde{\lambda}_n = \bar{x}_n$ . Note that

$$\frac{d^2}{d\lambda^2} l(\lambda) \big|_{\lambda=\tilde{\lambda}} = -\frac{n}{\tilde{\lambda}^2} \big|_{\lambda=\tilde{\lambda}} < 0.$$

So, the MLE is  $\hat{\lambda}_n = 1/\bar{X}_n$ .

# Invariance of MLE

## Theorem 7.2.10 (Invariance Property of MLE)

If  $\hat{\theta}$  is the MLE of  $\theta$ , then for any function  $\tau(\theta)$ , the MLE of  $\tau(\theta)$  is  $\tau(\hat{\theta})$ .

- For a function of  $\theta$ , we do not need to calculate the derivative with respect to this new function, but introduce in  $\hat{\theta}_n$  onto the function directly.
- Introduce the method as **Profile likelihood**. Assuming we are interested in only part of the parameters  $\eta$ , where  $\theta = (\eta, \xi)$ . Then the profile likelihood is defined by

$$L(\eta) = \max_{\xi} L(\eta, \xi)$$

Maximizing  $L(\eta)$  gives the same MLE for  $\hat{\eta}_n$  with the one we get from  $\hat{\theta} = (\hat{\eta}_n, \hat{\xi}_n)$ .

# Bayes Estimators

Let  $\mathbf{X} \sim f(\mathbf{x}|\theta)$ ,  $\theta \sim \pi(\theta)$ , where  $\pi(\theta)$  is the prior distribution of  $\theta$ . Then after observing  $\mathbf{X} = \mathbf{x}$ , the posterior distribution of  $\theta$  is

$$\pi(\theta|\mathbf{x}) = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{m(\mathbf{x})}$$

where  $m(\mathbf{x})$  is the marginal distribution of  $\mathbf{X}$ ,

$$m(\mathbf{x}) = \int f(\mathbf{x}|\theta)\pi(\theta)d\theta.$$



## Example 7.2.14

Let  $X_1, X_2, \dots, X_n$  be i.i.d. Bernoulli( $p$ ). Then  $Y = \sum_{i=1}^n X_i$  is Bin( $n, p$ ).

We assume the prior distribution on  $p$  is Beta( $\alpha, \beta$ ). Then

$$\begin{aligned} f(y, p) &= \left[ \binom{n}{y} p^y (1-p)^{n-y} \right] \left[ \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} \right] \\ &= \binom{n}{y} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{y+\alpha-1} (1-p)^{n-y+\beta-1} \\ f(y) &= \int_0^1 f(y, p) dp = \binom{n}{y} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot \frac{\Gamma(y + \alpha)\Gamma(n - y + \beta)}{\Gamma(n + \alpha + \beta)} \end{aligned}$$

Therefore  $f(p|y) = \frac{\Gamma(n+\alpha+\beta)}{\Gamma(y+\alpha)\Gamma(n-y+\beta)} p^{y+\alpha-1} (1-p)^{n-y+\beta-1}$ , which is Beta( $y + \alpha, n - y + \beta$ ). Bayes estimator of  $p$  is the mean of  $f(p|y)$ , which is

$$\hat{p}_B = \frac{y + \alpha}{\alpha + \beta + n}.$$

### Definition 7.2.6 (Conjugate Family)

Let  $\mathcal{F}$  denote the class of pdf or pmf  $f(x|\theta)$  (indexed by  $\theta$ ). A class  $\Pi$  of prior distribution is a conjugate family of  $\mathcal{F}$  if the posterior distribution is in the class  $\Pi$  for all  $f \in \mathcal{F}$ , all prior in  $\Pi$ , and all  $x \in \mathcal{X}$

# Bayesian Analysis

$$P(\theta) = \pi(\theta|\mathbf{y}) = \frac{f(\mathbf{y}|\theta)\pi(\theta)}{\int_{\theta} f(\mathbf{y}|\theta)\pi(\theta)d\theta}^1 \quad (12.3')$$

- Estimating the normalising constant
- Markov chain Monte Carlo (MCMC): Monte Carlo integration and Markov chain sampling

we estimated unknown parameters using the methods:

- Maximum likelihood : Newton - Raphson(NR)
- MCMC
- NR比MCMC收敛的快。
- 对于多峰分布，NR只能找到一个，而MCMC可以找到多个。
- NW对于固定的似然初始值，它的路径是一样的；而MCMC即使初始值相同，它的路径也是随机的。

---

<sup>1</sup>Annette J. Dobson and Adrian G. Barnett (2008). An Introduction to Generalized Linear Models Third Edition. Chapman & Hall/CRC

- **Target density**  $P(\theta)$  is not always achievable because it may have a complex, or even unknown, form.
- Markov chains provide a method of drawing samples from target densities (regardless of their complexity).
- Using these conditional steps, we build up a chain of samples  $(\theta^1, \dots, \theta^M)$  after specifying a starting value  $\theta^0$

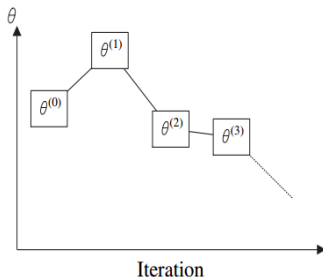


Figure 13.3 A simple example of a Markov chain.

Markov property:

$$P(\theta^i = a | \theta^{i-1}, \theta^{i-2}, \dots, \theta^0) = P(\theta^i = a | \theta^{i-1})$$

An algorithm for creating a Markov chain for a target probability density  $P(\theta)$  is:

- 1 Choose an initial value  $\theta^0$ . The restriction on the initial value is that it needs to be within the distribution of  $P(\cdot)$ , so that  $P(\theta) > 0$
- 2 Create a new sample using  $\theta^1 \sim \pi(\theta^1 | \theta^0, \mathbf{y})$
- 3 Repeat step 2 M times, each time increasing both indices by 1.

**Transitional density:**  $\pi(\theta^{i+1} | \theta^i)$ , Normal distribution.

# The Metropolis – Hastings sampler

- The Metropolis – Hastings sampler works by randomly **proposing a new value**  $\theta^*$
- If this proposed value is **accepted** (according to a criterion below),  $\theta^{i+1} = \theta^*$
- If this proposed value is **rejected** (according to a criterion below),  $\theta^{i+1} = \theta^i$
- Another proposal is made and the chain progresses by assessing this new proposal

- $\theta^* = \theta^i + Q$ ,  $Q$  is called the **proposal density**,  $N(0,1)$  or  $U[-1,1]$
- The acceptance criterion is:

$$\theta^{i+1} = \begin{cases} \theta^* & \text{if } U < \alpha \\ \theta^i & \text{otherwise} \end{cases}$$

where  $U \sim U[-1, 1]$  and

$$\alpha = \min\left\{\frac{\pi(\theta^*|\mathbf{y})}{\pi(\theta^i|\mathbf{y})} \cdot \frac{Q(\theta^i|\theta^*)}{Q(\theta^*|\theta^i)}, 1\right\}$$

where  $P(\theta|\mathbf{y})$  is the probability of  $\theta$  given the data  $\mathbf{y}$  (the likelihood)

If the proposal density is symmetric ( $Q(a|b) = Q(b|a)$ ), then

$$\alpha = \min\left\{\frac{\pi(\theta^*|\mathbf{y})}{\pi(\theta^i|\mathbf{y})}, 1\right\}$$

## 例1: The Metropolis - Hastings sampler

设  $Y_1, Y_2, \dots, Y_n \sim^{iid} N(\mu, \sigma^2)$ ,  $(\mu, \sigma^2)$  的先验分布为  $\pi(\mu, \sigma^2) \propto \frac{1}{\sigma^2}$ 。计算  $E(\mu|Y = y)$  和  $E(\sigma^2|Y = y)$ 。

解:  $(\mu, \sigma^2)$  后验分布为

$$\pi(\mu, \sigma^2|Y = y) \propto \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2}+1} \exp\left\{-\frac{\sum_{i=1}^n (y_i - \mu)^2}{2\sigma^2}\right\}$$

- 《统计模拟及其R实现》P203
- 为了比较方法, 假设  $Y_1, Y_2, \dots, Y_n \sim^{iid} N(2, 4^2)$



# The Gibbs sampler

- The Gibbs sampler is another way of generating a Markov chain.
- It splits the parameters into a number of components and then updates each one in turn.
- For the beetle mortality example, a Gibbs sampler to update the two unknown parameters would be:
  - 1 Assign an initial value to the two unknowns:  $\beta_1^0$  and  $\beta_2^0$
  - 2 (a) Generate  $\beta_2^1 \sim \pi(\beta_2 | \mathbf{y}, \beta_1^0)$   
(b) Generate  $\beta_1^1 \sim \pi(\beta_1 | \mathbf{y}, \beta_2^0)$
  - 3 Repeat the step 2  $M$  times, each time increasing the sample indices by 1.

上面举例比较简单，复杂的见《统计模拟及其R实现》P199 例8.4

## 例1续: The Gibbs sampler

设  $Y_1, Y_2, \dots, Y_n \sim^{iid} N(\mu, \sigma^2)$ ,  $(\mu, \sigma^2)$  的先验分布为  $\pi(\mu, \sigma^2) \propto \frac{1}{\sigma^2}$ 。计算  $E(\mu|Y = y)$  和  $E(\sigma^2|Y = y)$ 。

解:  $(\mu, \sigma^2)$  后验分布为

$$\pi(\mu, \sigma^2|Y = y) \propto \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2}+1} \exp\left\{-\frac{\sum_{i=1}^n (y_i - \mu)^2}{2\sigma^2}\right\}$$

则:

$$\pi(\mu|\sigma^2, y) \propto \exp\left\{-\frac{\sum_{i=1}^n (y_i - \mu)^2}{2\sigma^2}\right\} \propto N(\bar{y}, \frac{\sigma^2}{n}) \quad (1)$$

$$\pi(\sigma^2|\mu, y) \propto \left(\frac{1}{\sigma^2}\right) \exp\left\{-\frac{\sum_{i=1}^n (y_i - \mu)^2}{2\sigma^2}\right\} \propto IG\left(\frac{n}{2}, \frac{\sum_{i=1}^n (y_i - \mu)^2}{2}\right) \quad (2)$$

# 考虑Logit模型

例：考虑54个老人智力得分。

Table 7.8 *Symptoms of senility ( $s=1$  if symptoms are present and  $s=0$  otherwise) and WAIS scores ( $x$ ) for  $N=54$  people.*

$x$	$s$	$x$	$s$	$x$	$s$	$x$	$s$	$x$	$s$
9	1	7	1	7	0	17	0	13	0
13	1	5	1	16	0	14	0	13	0
6	1	14	1	9	0	19	0	9	0
8	1	13	0	9	0	9	0	15	0
10	1	16	0	11	0	11	0	10	0
4	1	10	0	13	0	14	0	11	0
14	1	12	0	15	0	10	0	12	0
8	1	11	0	13	0	16	0	4	0
11	1	14	0	10	0	10	0	14	0
7	1	15	0	11	0	16	0	20	0
9	1	18	0	6	0	14	0		

注：中科大张伟平《计算统计讲义》

考虑Logit模型:

$$Y_i \sim \text{Bin}(1, \pi_i), \quad \log \frac{\pi_i}{1 - \pi_i} = \beta_0 + \beta_1 x_i, \quad i = 1, 2, \dots, 54$$

则似然函数为:

$$\begin{aligned} f(\mathbf{y}|\beta_0, \beta_1) &= \prod_{i=1}^n \left( \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right)^{y_i} \left( \frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \right)^{1 - y_i} \\ &= \exp \left\{ \sum_{i=1}^n [(\beta_0 + \beta_1 x_i) y_i - \log(1 + e^{\beta_0 + \beta_1 x_i})] \right\} \end{aligned}$$

考虑 $\beta_0, \beta_1$ 的先验分布为独立的正态分布:

$$\beta_j \sim N(\mu_j, \sigma_j^2)$$

从而后验分布为：

$$\begin{aligned}
 f(\beta_0, \beta_1 | \mathbf{y}) &\propto f(\mathbf{y} | \beta_0, \beta_1) \pi(\beta_0, \beta_1) \\
 &\propto \exp \left\{ \sum_{i=1}^n [(\beta_0 + \beta_1 x_i) y_i - \log(1 + e^{\beta_0 + \beta_1 x_i})] \right. \\
 &\quad \left. - \frac{(\beta_0 - \mu_0)^2}{\sigma_0^2} - \frac{(\beta_1 - \mu_1)^2}{\sigma_1^2} \right\}
 \end{aligned}$$

# Diagnostics of chain convergence

- Assumption: the sample densities for the unknown parameters were good estimates of the target densities.
- If this assumption is incorrect, then inferences could be invalid.
- We can only make valid inference when a chain has **converged** to the target density.
  - Chain history
  - Chain autocorrelation
  - Multiple chains

# Chain history

A chain that has converged should show a **reasonable degree of randomness** between iterations, signifying that the Markov chain has found an area of high likelihood and is integrating over the target density (known as mixing).

- $\text{logit}(\pi_i) = \beta_1 + \beta_2 x_1$
- $\text{logit}(\pi_i) = \beta_1 + \beta_2 (x_1 - \bar{x})$

**Note:** By centring the dose covariate we have greatly **improved the convergence** because centring **reduces the correlation between the parameter estimates**  $\beta_1$  and  $\beta_2$

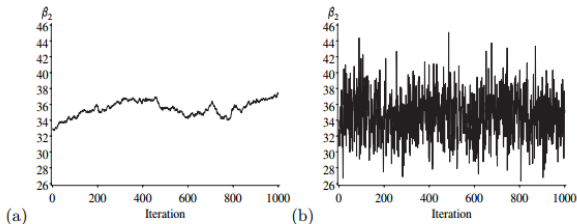


Figure 13.11 Example of a chain showing (a) poor convergence and (b) reasonable convergence (first 1,000 iterations using Gibbs sampling). Estimate for  $\beta_2$  using the logit link using two different parameterizations.



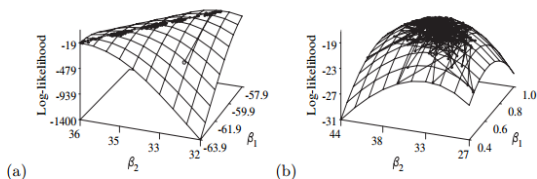


Figure 13.12 Three-dimensional plots of the log-likelihood and 200 Gibbs samples for the beetle mortality data using the logit link function and (a) uncentered dose or (b) centered dose. The initial value is shown as an open circle and subsequent estimates as closed circles.

## 13.5.2 Chain autocorrelation

- Autocorrelation is a useful diagnostic because it summarizes the dependence between neighbouring samples.
- Ideally we would like neighbouring samples to be completely independent, as this would be the most efficient chain possible.
- In practice we usually accept some autocorrelation, but large values (greater than 0.4) can be problematic.
- **Autocorrelation function (ACF)**

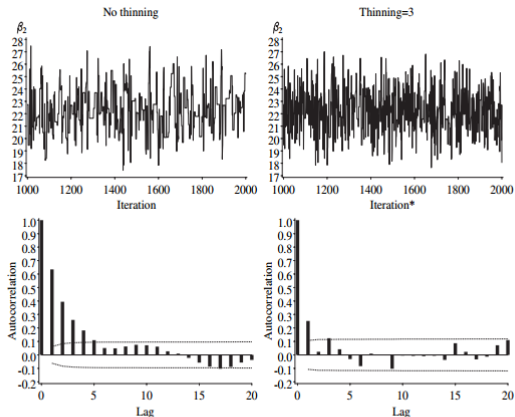


Figure 13.13 *Reduction in autocorrelation of Metropolis-Hastings samples after thinning. Chain history (top row) and ACF (bottom row) for the estimate of  $\beta_2$  from the beetle mortality example using the extreme value model and a centered dose.*

# Multiple chains

## Advantage:

- Using multiple chains is a good way to assess convergence.
- If we start multiple chains at widely varying starting values and each chain converges to the same solution, this would increase our confidence in this solution.
- This method is particularly good for assessing the influence of initial values.

**Drawback:** It may be difficult to generate suitably varied starting values, particularly for **complex problems** with **many unknown parameters** and **multi-dimensional likelihoods**.

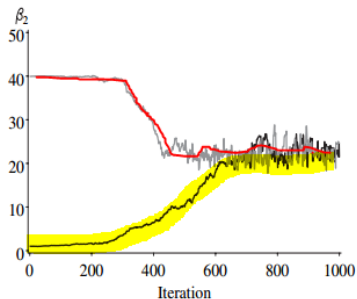


Figure 13.14 *Two chains with different starting values. Estimates of  $\beta_2$  using Metropolis–Hastings sampling for the extreme value model using the beetle mortality data.*

# EM Algorithm

- Consider a mixture model:

$$Z_i \sim \text{Bernoulli}(\theta), \quad X_i|Z_i \sim N(Z_i, \sigma^2).$$

However, our observations are  $X_i$ 's only.

- Since  $Z_i$ 's are missing. We call it as **missing data** problem
- The MLE is difficult to solve analytically for this problem.
- One generally used procedure is EM algorithm.

- Now there are two unknown things: the missing data  $z$  and the parameters  $\theta$
- If we **know**  $z$ , the MLE is found by

$$\arg \max_{\theta} \log f(x, z|\theta).$$

- If we **know**  $\theta$ , then we have the conditional distribution for  $z$ , which is

$$f(z|x, \theta) = \frac{f(x, z|\theta)}{f(x|\theta)} = \frac{f(x, z|\theta)}{\int f(x, z|\theta) dz}$$

- Combine these two things, find  $\theta$  as

$$\arg \max_{\theta} Q(\theta, \theta') = \arg \max_{\theta} E_Z \left[ \log f(x, z|\theta) | x; \theta' \right]$$

where the expectation is taken with respect to the conditional distribution of  $Z$ ,  $f(z|x, \theta')$ .

# Procedure

- ① Set initial value  $\theta^0, N = 1$ .
- ② Expectation step:
  - With  $\theta^{N-1}$ , find  $f(z|x; \theta^{N-1})$
  - Compute  $Q(\theta, \theta^{N-1}) = E_Z \left[ \log f(x, z; \theta) | x; \theta^{N-1} \right]$ , where  $Z \sim f(z|x; \theta^{N-1})$ .
- ③ Maximization step: Find  $\theta^N$  as

$$\theta^N = \arg \max_{\theta} Q(\theta, \theta^{N-1})$$

- ④ Repeat Steps 2-3 until  $\|\theta^{N-1} - \theta^N\| \leq \epsilon$ , where  $\epsilon$  is a pre-set threshold. Or, stop the algorithm when  $N$  is large enough.

Remark:

- Very popular method, since very good for complicated models
- Seems different in different applications
- However, maybe trapped by local maxima



# The EM Algorithm: Example

Suppose we observe  $Z_{1:n}$  and  $Y_{1:n}$  both independent random variables and independent of each other. In particular,  $Y_i \sim \text{Poisson}(\tau_i)$  and  $Z_i \sim \text{Poisson}(\tau_i)$ , where  $\theta = (\beta, \tau_1, \dots, \tau_n) \in R_+^{n+1}$  are the parameters.

- If we have the **full data**, then the joint density is

$$f_Y(y|\theta) = \prod_{i=1}^n \frac{(\beta\tau_i)^{y_i}}{y_i!} e^{-\beta\tau_i} \quad f_Z(z|\theta) = \prod_{i=1}^n \frac{\tau_i^{z_i}}{z_i!} e^{-\tau_i}$$

It is straightforward to find the MLEs;

$$\hat{\beta}_n = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n \hat{z}_i}, \quad \hat{\tau}_i = \frac{y_i + z_i}{\hat{\beta}_n + 1}, i = 1, \dots, n.$$

- Now, if  $z_1$  was missing, we have the marginal data likelihood of the observations:

$$\begin{aligned} f(y, z_{2:n}; \theta) &= \frac{(\beta\tau_1)^{y_1}}{y_1!} e^{-\beta\tau_1} \left( \prod_{i=2}^n \frac{(\beta\tau_i)^{y_i}}{y_i!} e^{-\beta\tau_i} \frac{\tau_i^{z_i}}{z_i!} e^{-\tau_i} \right) \sum_{z_1=0}^{\infty} \frac{\tau_1^{z_1}}{z_1!} e^{-\tau_1} \\ &= \frac{(\beta\tau_1)^{y_1}}{y_1!} e^{-\beta\tau_1} \left( \prod_{i=2}^n \frac{(\beta\tau_i)^{y_i}}{y_i!} e^{-\beta\tau_i} \frac{\tau_i^{Z_i}}{z_i!} e^{-\tau_i} \right) \end{aligned}$$

Now the Q function is

$$\begin{aligned}
 & Q(\theta, \theta') \\
 &= \sum_{z_1=0}^{\infty} \log \left[ \prod_{i=1}^n \frac{(\beta \tau_i)^{y_i}}{y_i!} e^{-\beta \tau_i} \frac{\tau_i^{z_i}}{z_i!} e^{-\tau_i} \right] \frac{(\tau'_1)^{z_1}}{z_1!} e^{-\tau'_1} \\
 &= \sum_{i=1}^n \left( -\beta \tau_i + y_i [\log \beta + \log \tau_i] - \log y_i! \right) + \sum_{i=2}^n \left[ -\tau_i + z_i \log \tau_i - \log z_i! \right] \\
 &\quad + \sum_{z_1=0}^{\infty} \log \left( -\tau_1 + z_1 \log \tau_1 - \log z_1! \right) \frac{(\tau'_1)^{z_1}}{z_1!} e^{-\tau'_1} \\
 &= \sum_{i=1}^n \left( -\beta \tau_i + y_i [\log \beta + \log \tau_i] \right) + \sum_{i=2}^n \left[ -\tau_i + z_i \log \tau_i \right] \\
 &\quad + \sum_{z_1=0}^{\infty} \log \left( -\tau_1 + z_1 \log \tau_1 - \log z_1! \right) \frac{(\tau'_1)^{z_1}}{z_1!} e^{-\tau'_1} + \sum_{i=1}^n -\log y_i \\
 &\quad + \sum_{i=2}^n -\log z_i!
 \end{aligned}$$

Note that

$$\sum_{z_1=0}^{\infty} \log \left( -\tau_1 + z_1 \right) \frac{(\tau_1')^{z_1}}{z_1!} e^{-\tau_1'} = -\tau_1 + \tau_1' \log \tau_1$$

and the last term does not depend on  $\theta$ , which can be denoted as  $C$ . So we have

$$\begin{aligned} Q(\theta, \theta') \\ &= \sum_{i=1}^n \left( -\beta \tau_i + y_i [\log \beta + \log \tau_i] \right) + \sum_{i=2}^n \left[ -\tau_i + z_i \log \tau_i \right] - \tau_1 + \tau_1' \log \tau_1 + C \end{aligned}$$

Maximizing  $Q(\theta, \theta')$  w.r.t.  $\theta$ , the solution is

$$\beta = \frac{\sum_{i=1}^n y_i}{\tau_1' + \sum_{i=2}^n z_i} \quad \tau_1 = \frac{\tau_1' + y_1}{\beta + 1}, \quad \tau_i = \frac{y_1 + z_i}{\beta + 1}, i = 2, \dots, n$$

# Methods of Evaluating Estimators

The **mean square error (MSE)** of an estimator  $W$  of a parameter  $\theta$  is the function of  $\theta$  defined by

$$E_{\theta}(W - \theta)^2.$$

$\text{Bias}_{\theta}W = E_{\theta}W - \theta$ . If  $\text{Bias}_{\theta}W = 0$ , then  $W$  is unbiased.

- **Example 7.3.3 (Normal MSE)** Let  $X_1, X_2, \dots, X_n$  be i.i.d.  $N(\mu, \sigma^2)$ . Then statistics  $\bar{X}$  and  $S^2$  are both unbiased.

$$\begin{aligned} \text{MSE}(\bar{X}) &= E(\bar{X} - \mu)^2 = \text{Var}(\bar{X}) = \sigma^2/n \\ E(S^2 - \sigma^2)^2 &= \text{Var}(S^2) = \frac{2\sigma^4}{n-1} \end{aligned}$$

---

2

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1) \Rightarrow \text{Var}\left(\frac{(n-1)S^2}{\sigma^2}\right) = 2(n-1)$$

- **Example 7.3.4** Maximum Likelihood estimator of  $\sigma^2$  is

$$\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n-1}{n} S^2.$$

$$\begin{aligned} \text{Var} \left( \frac{n-1}{n} S^2 \right) &= \frac{(n-1)^2}{n^2} \cdot \frac{2\sigma^4}{n-1} = \frac{2(n-1)}{n^2} \sigma^4 \\ \text{MSE} \left( \frac{n-1}{n} S^2 \right) &= \left( \frac{n-1}{n} ES^2 - \sigma^2 \right)^2 + \frac{2(n-1)}{n^2} \sigma^4 \\ &= \sigma^4 \left( \frac{n-1}{n} - 1 \right)^2 + \frac{2(n-1)}{n^2} \sigma^4 \\ &= \sigma^4 \frac{2n-1}{n^2} \end{aligned}$$

Since

$$\frac{2n-1}{n^2} < \frac{2}{n-1},$$

So in this case MLE has smaller MSE than the unbiased estimator  $S^2$ .

**Remark** While MSE is a reasonable measurement for location parameters, it may not be a good to compare estimators of scale parameters with MSE.

**Example 7.3.5 (MSE of binomial Bayes Estimator)**  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ .

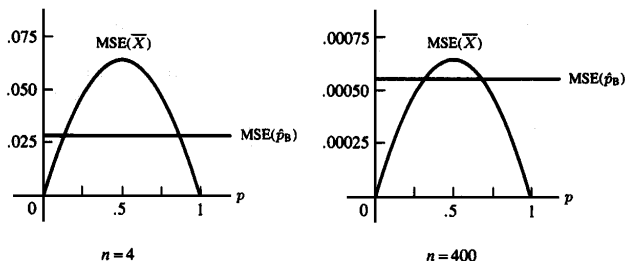
- Let  $\hat{p} = \frac{X_1 + \dots + X_n}{n}$ .  $E_p(\hat{p} - p)^2 = \text{Var}_p(\bar{X}) = \frac{p(1-p)}{n}$ .
- Let  $\hat{p}_B = \frac{Y + \alpha}{\alpha + \beta + n}$  be the Bayes estimator. Here  $Y = \sum_{i=1}^n X_i$

$$\begin{aligned} \text{MSE}(\hat{p}) &= \text{Var}_p(\hat{p}_B) + (\text{Bias}_p(\hat{p}_B))^2 \\ &= \text{Var}\left(\frac{Y + \alpha}{\alpha + \beta + n}\right) + \left(E_p\left(\frac{Y + \alpha}{\alpha + \beta + n}\right) - p\right)^2 \\ &= \frac{np(1-p)}{(\alpha + \beta + n)^2} + \left(\frac{np + \alpha}{\alpha + \beta + n} - p\right)^2 \end{aligned}$$

In the absence of good prior information about  $p$ , we might choose  $\alpha$  and  $\beta$  to make the MSE of  $\hat{p}_B$  constant. Choose  $\alpha = \beta = \sqrt{n/4}$  gives

$$\hat{p}_B = \frac{Y + \sqrt{n/4}}{n + \sqrt{n}}, \quad E(\hat{p}_B - p)^2 = \frac{n}{4(n + \sqrt{n})^2}$$

**Figure 7.3.1** Comparison of  $MSE(\hat{p})$  and  $MSE(\hat{p}_B)$  for sample size  $n = 4$  and  $n = 400$  in Example 7.3.5



- For small  $n$ ,  $\hat{p}_B$  is the better choice (unless there is a strong belief that  $p$  is near 0 or 1)
- For large  $n$ ,  $\hat{p}$  is the better choice (unless there is a strong belief that  $p$  is close to  $\frac{1}{2}$ )



# Best Unbiased Estimators

As we have discussed, there is usually no "best MSE" estimator. However, if we restrict our choice from [unbiased estimators](#), then there exists best estimator in this class.

**Definition 7.3.7** An estimator  $W^*$  is a best unbiased estimator of  $\tau(\theta)$  if it satisfies  $E_{\theta}W^* = \tau(\theta)$  for all  $\theta$  and, for any other estimator  $W$  with  $E_{\theta}W = \tau(\theta)$ , we have

$$\text{Var}_{\theta}W^* \leq \text{Var}_{\theta}W \text{ for all } \theta.$$

$W^*$  is also called a uniform minimum variance unbiased estimator (UMVUE) of  $\tau(\theta)$ .

\* Finding UMVUE is not easy.

- **Theorem 7.3.9 (Cramér-Rao Inequality)** Let  $X_1, X_2, \dots, X_n$  be a sample with pdf  $f(\mathbf{x}|\theta)$ , and let  $W(\mathbf{X})$  be any estimator satisfying

$$\frac{d}{d\theta} E_{\theta} W(\mathbf{X}) = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} [W(\mathbf{x}) f(\mathbf{x}|\theta)] d\mathbf{x}$$

and  $\text{Var}_{\theta} W(\mathbf{X}) < \infty$ . Then

$$\text{Var}_{\theta} W(\mathbf{X}) \geq \frac{\left( \frac{d}{d\theta} E_{\theta} W(\mathbf{X}) \right)^2}{E_{\theta} \left( \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X}|\theta) \right)^2 \right)}$$

In particular, if  $W(X)$  is an unbiased estimator of  $\theta$ , then

$$\text{Var}_{\theta} W(\mathbf{X}) \geq \frac{1}{E_{\theta} \left( \left( \frac{\partial}{\partial \theta} \log f(\mathbf{X}|\theta) \right)^2 \right)}$$

- **Corollary 7.3.10 (Cramér-Rao Inequality, iid case)** If  $X_1, X_2, \dots, X_n$  are i.i.d.  $f(x|\theta)$ , and the condition of Theorem 7.3.9 are satisfied, then

$$\text{Var}_{\theta} (W(\mathbf{X})) \geq \frac{\left( \frac{d}{d\theta} E_{\theta} W(\mathbf{X}) \right)^2}{n E_{\theta} \left( \left( \frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2 \right)}$$

To evaluate  $E_{\theta} \left( \left( \frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2 \right)$ , we have the following Lemma.

- **Lemma 7.3.11** If  $f(x|\theta)$  satisfies

$$\frac{d}{d\theta} E_{\theta} \left( \frac{\partial}{\partial \theta} \log f(x|\theta) \right) = \int \frac{\partial}{\partial \theta} \left[ \left( \frac{\partial}{\partial \theta} \log f(x|\theta) \right) f(x|\theta) \right] dx$$

(true for an exponential family), then

$$E_{\theta} \left( \left( \frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2 \right) = -E_{\theta} \left( \frac{\partial^2}{\partial \theta^2} \log f(X|\theta) \right).$$

- **Example 7.3.12**  $\bar{X}$  is UMVUE for  $\lambda$  if  $X_1, \dots, X_n$  are i.i.d. Poisson( $\lambda$ ). From [Theorem 7.3.9](#), we have for any unbiased estimator  $W(\mathbf{X})$  of  $\lambda$ .

$$\text{Var}_\lambda W(\mathbf{X}) \geq \frac{1}{-nE_\lambda \left[ \frac{\partial^2}{\partial \lambda^2} \log f(\mathbf{x}|\lambda) \right]} \quad (8.1)$$

$$\log f(\mathbf{x}|\lambda) = \log \left[ e^{-\lambda} \frac{\lambda^x}{x!} \right] = -\lambda + x \log \lambda - \log x!$$

$$\frac{\partial^2}{\partial \lambda^2} \log f(\mathbf{x}|\lambda) = -x \frac{1}{\lambda^2}.$$

Therefore,  $-E_\lambda \left[ \frac{\partial^2}{\partial \lambda^2} \log f(\mathbf{x}|\lambda) \right] = \frac{1}{\lambda^2} E_\lambda X = \frac{1}{\lambda}.$

(8.1) Becomes  $\text{Var}_\lambda(W(\mathbf{X})) \geq \frac{\lambda}{n}.$

But  $\text{Var}_\lambda(\bar{X}) = \frac{\lambda}{n}.$

- **Example 7.3.13 (Unbiased Estimator for Scale Parameter)** Let  $X_1, \dots, X_n$  be i.i.d. with pdf  $f(x|\theta) = \frac{1}{\theta}$ ,  $0 < x < \theta$ . Since  $\frac{\partial}{\partial \lambda} \log f(x|\theta) = -\frac{1}{\theta}$ , we have

$$E_{\theta} \left[ \frac{\partial}{\partial \lambda} \log f(x|\theta) \right] = \frac{1}{\theta^2}$$

So if  $W$  is unbiased for  $\theta$ , then

$$\text{Var}_{\theta}(W) \geq \frac{\sigma^2}{n}.$$

- On the other hand,  $Y = \max(Y_1, \dots, Y_n)$  is a sufficient statistic.  $f_Y(y|\theta) = ny^{n-1}/\theta^n$ ,  $0 < y < \theta$ . So

$$E_{\theta} Y = \int_0^{\theta} y \cdot \frac{ny^{n-1}}{\theta^n} dy = \frac{n}{n+1} \theta,$$

showing that  $\frac{n+1}{n} Y$  is an unbiased estimator of  $\theta$ .

$$\begin{aligned}\text{Var}_{\theta} \left( \frac{n+1}{n} Y \right) &= \left( \frac{n+1}{n} \right)^2 \text{Var}_{\theta}(Y) \\&= \left( \frac{n+1}{n} \right)^2 [E_{\theta} Y^2 - (EY)^2] \\&= \left( \frac{n+1}{n} \right)^2 \left[ \frac{n}{n+2} \theta^2 - \left( \frac{n}{n+1} \theta \right)^2 \right] \\&= \frac{1}{n(n+2)} \theta^2,\end{aligned}$$

which is uniformly smaller than  $\theta^2/n$ . Cramér-Rao lower bound Theorem is not applicable to this pdf since

$$\frac{d}{d\theta} \int_0^{\theta} h(x) f(x|\theta) dx \neq \int_0^{\theta} h(x) \frac{\partial}{\partial \theta} f(x|\theta) dx.$$

- **Example 7.3.14 (Normal Variance Bound)** Let  $X_1, \dots, X_n$  be i.i.d.  $N(\mu, \sigma^2)$ . The conditions of Cramér-Rao Theorem are satisfied. Let  $W$  be an unbiased estimator of  $\sigma^2$ , then

$$\text{Var}(W|\mu, \sigma^2) \geq 2\sigma^4/n.$$

In Example 7.3.3 we see that  $\text{Var}(S^2|\mu, \sigma^2) \geq \frac{2\sigma^4}{n-1}$ . So  $S^2$  does not attain the Cramér-Rao lower bound.

- **Corollary 7.3.15 (Attainment)** Let  $X_1, \dots, X_n$  be i.i.d.  $f(x|\theta)$ , where  $f(x|\theta)$  satisfies the conditions of the Cramér-Rao Theorem. Let  $L(\theta|\mathbf{x})$  denote the likelihood function. If  $W(\mathbf{X})$  is any unbiased estimator of  $\tau(\theta)$ , then  $W(\mathbf{X})$  attains the Cramér-Rao lower bound if and only if

$$a(\theta)[W(\mathbf{x}) - \tau(\theta)] = \frac{\partial}{\partial \theta} \log L(\mathbf{x}|\theta)$$

for some function  $a(\theta)$ .

- **Proof** The Cramér-Rao inequality, can be written as

$$\begin{aligned} & \left[ \text{Cov}_{\theta} \left( W(\mathbf{X}), \frac{\partial}{\partial \theta} \log \prod_{i=1}^n f(X_i|\theta) \right) \right]^2 \\ & \leq \text{Var}_{\theta} W(\mathbf{X}) \cdot \text{Var}_{\theta} \left( \frac{\partial}{\partial \theta} \log L(\mathbf{X}) \right) \end{aligned}$$

Using the condition for "=" in Cauchy-Schwarz inequality, we obtain the expression (8.1).



- **Example 7.3.16 (Continuation of Example 7.3.14)**

$$L(\mathbf{x}|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 / \sigma^2\right)$$

and hence

$$\frac{\partial}{\partial \sigma^2} \log L(\mathbf{x}|\mu, \sigma^2) = \frac{n}{2\sigma^4} \left( \sum_{i=1}^n \frac{(x_i - \mu)^2}{n} - \sigma^2 \right)$$

Taking  $a(\sigma^2) = \frac{n}{2\sigma^4}$  shows that the best unbiased estimator of  $\sigma^2$  is  $\sum_{i=1}^n (x_i - \mu)^2 / n$ , which is calculable only if  $\mu$  is known.

- So the question of finding best unbiased estimator are still unsolved for many common pdf's.