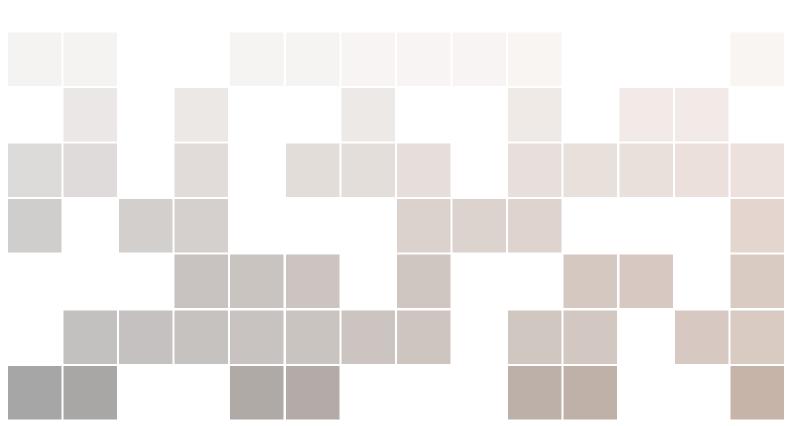


STAT 514 Lecture Notes

Dr. David Hunter



Copyright © 2014 John Smith

PUBLISHED BY PUBLISHER

BOOK-WEBSITE.COM

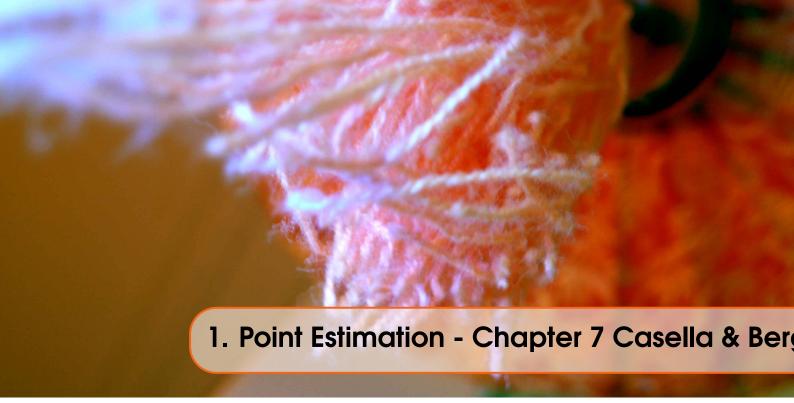
Licensed under the Creative Commons Attribution-NonCommercial 3.0 Unported License (the "License"). You may not use this file except in compliance with the License. You may obtain a copy of the License at http://creativecommons.org/licenses/by-nc/3.0. Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

First printing, March 2013



1	Point Estimation - Chapter 7 Casella & Berger	. 5
1.1	Introduction	5
1.2	Mean Squared Error	6
1.3	Best Unbias Estimator	7
1.4	Lost Function Optimality	9
1.5	Practice Problems	10
2	Principles of Data Reduction	13
2.1	Sufficiency Principle	13
2.2 2.2.1	Exponential Families Distribution Families	14 15
2.3 2.3.1	Complete Statistic Sufficiency and Unbiasness	16
2.4	Methods of Finding Estimators	18
2.4.1 2.4.2	Method of Moments Maximum Likelihood Estimation	
2.5	Major/Minorization Mini/Maximization	20
2.5.1	How EM Algorhythms Work	20
2.6 2.6.1	Bayes Estimators Model Based Clustering	21 22
3	Hypothesis Testing	23
	Bibliography	25
	Rooks	25

Article	S															25
Index		 	 	 												27



1.1 Introduction

In the simplest case, we have n observations of data that we believe follow the same distribution.

$$X_1,\ldots,X_n \stackrel{iid}{\sim} f_{\theta}(x)$$

where $f_{\theta}(x)$ is a density function involving a parameter θ . Our goal is to learn something about θ , which could be real or vector valued.

Definition 1.1.1 — Estimator. An *estimator* of θ is any function $W(X_1, ..., X_n)$ of the data. That is, an estimator is a *statistic*.

Note:

- 1. W(X) may not depend on θ .
- 2. W(X) should resemble or "be close" to θ .
- 3. An estimator is *random*.
- 4. $W(X_1,...,X_n)$ is the estimator, $W(x_1,...,x_n)$ is the fixed estimate.
- **Example 1.1** Suppose we have n observations from an exponential distribution,

$$X_1, \dots, X_n \stackrel{iid}{\sim} f_{\theta}(x) = \frac{1}{\theta} \exp\left\{-\frac{x}{\theta}\right\} \mathbb{1}\{x > 0\}$$

for some $\theta > 0$. The **likelihood function** is equivalent to the joint density function, expressed as a function of θ rather than the data:

$$L(\theta) = \prod_{i=1}^{n} \frac{1}{\theta} \exp\left\{-\frac{x}{\theta}\right\} = \frac{1}{\theta^{n}} \exp\left\{-\frac{1}{\theta} \sum_{i=1}^{n} x_{i}\right\}$$

This function represents the *likelihood* of observing the data we observed assuming the parameter was a particular value of θ . If we can maximize this function, we can determine the $\hat{\theta}$ for which the likelihood of observing \boldsymbol{X} was the highest. This might tell us something about the true value of θ .

To maximize $L(\theta)$, we want to take the derivative, set it equal to 0, and solve for θ . However, in many cases taking the derivative of the likelihood function will be very hard, if not impossible.

We can use the fact that taking the logarithm does not change the location of extrema. The **log-likelihood function** in this case is

$$\ell(\theta) = \log L(\theta) = -n\log \theta - \frac{1}{\theta} \sum_{i=1}^{n} x_i$$

Take the derivative with respect to the parameter and set equal to 0:

$$\ell'(\theta) = -\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^{n} x_i \stackrel{\text{set}}{=} 0$$

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

Here $\hat{\theta}$ is an estimator (the sample mean). Since it maximizes $L(\theta)$, we call it the **maximum** likelihood estimator (MLE).

1.2 Mean Squared Error



Read Casella & Berger Chapter 7.3 - Methods of Evaluating Estimation

Definition 1.2.1 — Mean Squared Error. If W(X) is an estimator of θ , then the mean squared error (MSE) is defined as

$$E_{\theta}\left[(W(\boldsymbol{X})-\theta)^2\right].$$

Definition 1.2.2 — Unbiased estimator. If W(X) is an estimator of θ , we say that W(X) is **unbiased** if

$$E_{\theta}[W(\mathbf{X})] = \theta \quad \forall \theta.$$

Furthermore, the **bias** of W(X) is

$$E_{\boldsymbol{\theta}}[W(\boldsymbol{X})] - \boldsymbol{\theta}$$
.

Example 1.2 For $\theta > 0$, let

$$X_1,\ldots,X_n \stackrel{iid}{\sim} f_{\theta}(x) = \theta x^{-2} \mathbb{1}\{x > \theta\}$$

Find the MLE of θ .

$$L(\theta) = \theta^n \prod_{i=1}^n x_i^{-2} \prod_{i=1}^n \mathbb{1}\{x > \theta\}$$

$$\hat{\theta} = \text{minimum of } x_i$$

Theorem 1.2.1 $MSE(W) = bias^2 + Var(W)$

Proof:

$$E[(W(X) - \theta)^{2}] = E[(W - E[W] + E[W] - \theta)^{2}]$$

$$= E[(W - E[W])^{2}] + E[(E[W] - \theta)^{2}] + 2E[(W - E[W])(E[W] - \theta)]$$

$$= Var(W) + bias^{2}(W) + 0$$

Best Unbias Estimator

What does best mean? Answer: Minimum variance.

Recall: Given an esitmator W(X) for θ ,

$$MSE(\theta) = E(W(X) - \theta)^2$$

Definition 1.3.1 — Best Unbiased Estimator. An estimator W* is a best unbiased estima**tor*** of $\tau(\theta)$ if it satisfies

$$E_{\theta}(W*) = \tau(\theta)$$

for all θ and, for any other estimator W with

$$E_{\theta}(W) = \tau(\theta)$$

we have

$$\operatorname{Var}_{\theta}(W*) \leq \operatorname{Var}_{\theta}(W)$$

for all θ . W* is also called a *uniform minimum variance unbiased estimator* (UMVUE) of

Main Result: Under some assumptions we can establish a lower bound on Var(W(X)).

Theorem 1.3.1 — Cremer-Rao Inequality. Also: Information Inequality.

$$\operatorname{Var}(W(\underline{X})) \ge \frac{(\Psi'(\theta))^2}{I(\theta)}$$

where, $\Psi(\theta) = E_{\theta}(W(\underline{X}))$

and, the Fisher/Expected information, $I(\theta) = E\left(\left(\frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})\right)^{2}\right)$.

Proof: Follows from Cauchy-Schwarz Inequality

$$Cov(W(\underline{x}), \frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})))^{2} \leq Var(W(\underline{X})) * Var(\frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X}))$$

Assumptions:

- 1. $I(\theta) = E_{\theta}\left(\left(\frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})^2\right)\right) = Var(\left(\frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})^2\right))$ is well defined and $I(\theta) > 0$.
- 2. $E\left(\left(\frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})\right)^{2}\right) = 0$ Thus, $I(\theta)$ is just varience. 3. $E\left(\left(W(\underline{X})\frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})\right)^{2}\right) = \Psi'(\theta)$

$$Cov(W(\underline{x}), \frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})))^{2} \leq Var(W(\underline{X})) * Var(\frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})) \simeq Var(W(\underline{X})) \geq \frac{(\Psi'(\theta))^{2}}{I(\theta)}$$

Exercise 1.1 Let $\underline{X} \sim Poi(\theta)$, $Y = \sum x_i$, and $Y \sim Poi(n\theta)$. What is $I(\theta)$?

$$I(\theta) = E\left(\left(\frac{\delta}{\delta\theta}\log f_{\theta}(\underline{X})^{2}\right)\right)$$

$$= E\left(\left(\frac{\delta}{\delta\theta}\log \frac{\theta^{\Sigma}e^{-\theta}}{x!}\right)^{2}\right)$$

$$= E\left(\left(\frac{\delta}{\delta\theta}\log \frac{\theta^{\Sigma}ie^{-\theta n}}{\sum x_{i}!}\right)^{2}\right)$$

$$= E\left(\frac{\delta}{\delta\theta}(-n\theta + \sum x_{i}\log \theta - \sum \log x_{i}!\right)^{2}\right)$$

$$= E\left(\left(-n + \frac{\sum x_{i}}{\theta}\right)^{2}\right)$$

$$= E\left(n^{2} - 2(\frac{\sum x_{i}}{\theta})(n) + (\frac{\sum x_{i}}{\theta})^{2}\right)$$

$$= n^{2} - \frac{2n * E(\sum x_{i})}{\theta} + \frac{E((\sum x_{i})^{2})}{\theta^{2}}$$

$$= n^{2} - \frac{2n * E(Y)}{\theta} + \frac{E(Y)^{2}}{\theta}$$

$$= \frac{n}{\theta}$$

Note: $I(\theta)$ is equal to information in the whole sample, but sometimes it's just one sample (based on context).

If we assume (as in any exponential family):

$$E_{\theta}\left(\frac{\delta^{2}}{\delta\theta^{2}}\log f_{\theta}(\underline{x})\right) = \frac{\delta^{2}}{\delta\theta^{2}}\int \log f_{\theta}(\underline{x})f_{\theta}(\underline{x})dx$$

then, the observed information is

$$I(\theta) = -E_{\theta} \left(\frac{\delta^2}{\delta \theta^2} \log f_{\theta}(\underline{x}) \right)$$

■ Example 1.3 Let $X_i ... X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ Find the information based on σ^2 . Write

$$\begin{split} \ell(\mu,\sigma^2) &= \log \left(\prod (2\pi\sigma^2)^{\frac{1}{2}} \exp\left\{ \frac{-(x-\mu)^2}{2\sigma^2} \right\} \right) \\ &= \frac{-n}{2} \log(2\pi\sigma^2) + \sum \left(\frac{-(x_i-\mu)^2}{2\sigma^2} \right) \end{split}$$

If we try to use the expected information:

$$I(\sigma^2) = E\left(\frac{\delta}{\delta\theta}\left(\frac{-n}{2}\log(2\pi\sigma^2) + \sum\left(\frac{-(x_i - \mu)^2}{2\sigma^2}\right)\right)^2\right)$$
 which is a mess.

However, using the observed information:

ı

$$I(\sigma^2) = \frac{-n}{2\sigma^4} + \frac{1}{\sigma^6} E\left(\sum (X_i - \mu)^2\right)$$
$$= \frac{-n}{2\sigma^4} + \frac{n\sigma^2}{\sigma^6}$$
$$= \frac{n}{2\sigma^4}$$

Example 1.4 Continued from previous example.

Define
$$S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$$

Note:
$$\frac{n-1}{\sigma^2}S^2 \sim \chi_{n-1}^2$$

Define $S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$ Note: $\frac{n-1}{\sigma^2} S^2 \sim \chi_{n-1}^2$ Does S^2 acheive the C-R lower bound?

$$Var(S^{2}) \ge \frac{\Psi'(\theta)}{I(\theta)}$$
$$\frac{2\sigma^{4}}{n-1} \ge \frac{1}{I(\theta)}$$
$$\ge \frac{2\sigma^{4}}{n}$$

Reread 7.3 and 6.2 in Casella and Berger

What would give equality? When $W(\underline{X})$ is a linear function of $\frac{\delta}{\delta\theta} \log f_{\theta}(\underline{X})$ which leads

Corollary 1.3.2 — Attainment. Let X_1, \ldots, X_n be iid $f(x|\theta)$, where $f(x|\theta)$ satisfies the conditions of the Cramer-Rao Theorem. Let $L(\theta|x) = \prod f(x_i|\theta)$ denote the likelihood function. If $W(X) = W(X_1, \dots, X_n)$ is any unbiased estimator of $\tau(\theta)$, then W(X) attains the Cramer-Rao Lower Bound iff

$$a(\theta)(W(x) - \tau(\theta)) = \frac{\delta}{\delta\theta} \log L(\theta|x)$$

for some function $a(\theta)$.

1.4 Lost Function Optimality

Definition 1.4.1 — Loss. $L(\theta, W(x))$ assigns a nonnegative real value called the **loss** to our decision to estimate θ by W(\underline{X}). General context: Decision Theory.

Typically $L(\theta, \theta) = 0$ because nothing is lost if your decision is exactly correct.

■ Example 1.5

$$L(\theta, W(\underline{X})) = (\theta - W(X)^2)$$
 square error loss
$$= |\theta - W(X)|$$
 absolute error loss
$$= \frac{W(X)}{\theta} - 1 - \log \frac{W(X)}{\theta}$$
 Stein's loss

Definition 1.4.2 — Risk. Risk of estimating θ by $W(\underline{X})$ is

$$R(\theta, W) = E(\theta, W(\underline{X}))$$

Exercise 1.2 If $X_1, ..., X_n$ are iid with mean μ and varience σ^2 what is

$$E\left(\sum_{i=1}^n (X_i - \bar{X})^2\right)?$$

Note:

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

2.
$$Var(X) = E(X^2) - E(X)^2$$

3.
$$\operatorname{Var}(\frac{X_i}{n}) = \frac{\sigma^2}{n^2}$$

4.
$$\operatorname{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

Thus.

$$E(\sum X_i^2) = (\sigma^2 + \mu^2)n$$

$$-nE(\bar{X}^2) = -n(\mu^2 + \frac{\sigma^2}{n})$$

We may conclude, $\frac{1}{n-1}\sum (X_i - \bar{X})^2$ is an unbiased estimator of σ^2 called S^2 .

Theorem 1.4.1 — Cauchy-Scwarz Inequality.

$$|| < x, y > ||^2 \le || < x, x > || * || < y, y > ||$$

In terms of E(X) if A & B have $\mu = 0$:

$$(E(AB))^2 \le E(A^2) * E(B^2)$$

Or in terms of covariance:

$$Cov^2(AB) \le Var(A) * Var(B)$$

Proof: Let $D = B - \frac{E(AB)}{E(A^2)}A$, given $D^2 \ge 0$.

$$\begin{split} E(D^2) &= E(B^2 - 2(B)(\frac{E(AB)}{E(A^2)}A) + (\frac{E(AB)}{E(A^2)}A)^2) \\ &= E(B^2) - 2\left(\frac{E(AB)^2}{E(A^2)}\right) + \frac{E(AB)}{E(A^2)}E(A^2) \\ &= E(B^2) - \frac{E(AB)^2}{E(A^2)} \ge 0 \end{split}$$

1.5 Practice Problems

Chapter 7 Exercises (pp. 355 to 367): 7.2, 7.6, 7.7, 7.10, 7.12, 7.13, 7.14, 7.19, 7.20, 7.21, 7.22, 7.23, 7.25, 7.30, 7.37, 7.38, 7.44, 7.47, 7.48, 7.49, 7.52

Problem 1.1 (C&B 7.2) Let $X_1, ..., X_n$ be a random sample from a Gamma(α, β) population. (a) Find the MLE of β , assuming α is known.

(b) If α and β are both unknown, there is no explicit formula for the MLEs of α and β , but the maximum can be found numberically. The result in part (a) can be used to reduce the problem to the maximization of a univariate function. Find the MLEs for α and β for the data in Exercise 7.10(c).

Solution.

(a) Gamma distribution pdf: $\frac{\beta^{\alpha}}{\Gamma(\alpha)}x^{\alpha-1}e^{-\beta x}$

$$L(\beta) = \frac{\beta^{n\alpha}}{\Gamma(\alpha)^n} [\prod x]^{\alpha - 1} e^{-\beta \sum x}$$

$$\ell(\beta) = n\alpha \log \beta - n \log \Gamma(\alpha) + (\alpha - 1) \log(\prod x_i) - \beta \sum x_i$$

$$\frac{\delta\ell}{\delta\beta} = \frac{n\alpha}{\beta} - \sum x_i$$

$$\hat{\beta} = \frac{n\alpha}{\sum x_i}$$

- R The density of Gamma can also be in the form with $\frac{1}{\beta}$ which would mean $\hat{\beta}$ would be $\frac{\sum x_i}{n\alpha}$.
- (b) The new likelihood with $\hat{\beta}$ plugged in is:

$$L(\beta) = \frac{\left(\frac{n\alpha}{\sum x_i}\right)^{n\alpha}}{\Gamma(\alpha)^n} \left[\prod x\right]^{\alpha - 1} e^{-\left(\frac{n\alpha}{\sum x_i}\right)\sum x}$$

To solve...use a computer.

Problem 1.2 (C&B 7.6) Let $X_1, ..., X_n$ be a random sample from the pdf

$$f(x|\theta) = \theta x^{-2}, 0 < \theta \le x < \infty$$

- (a) What is a sufficient statistic for θ ?
- (b) Find the MLE of θ .
- (c) Find the method of moments estimator of θ ?

Solution.

(a) Use Factorization Theorem.

$$L(\theta) = \theta^n \prod_i (x_i^2) \prod_i I\{\theta \le x < \infty\}$$
 where $\prod_i I\{\theta \le x < \infty = x_{(1)}\}$

Thus the sufficient statistic is the minimum x.

- (b) $L(\theta|x) = \theta^n \prod (x_i^2) \prod I\{\theta \le x < \infty\}$. θ^n is increasing in θ . The second term does not involve θ . So to maximize $L(\theta|x)$, we want to make θ as large as possible. But because of the indicator function, $L(\theta|x) = 0$ if $\theta > x(1)$. Thus, $\hat{\theta} = x_{(1)}$.
- (c) $E(X) = \int_{\theta}^{\infty} \theta x^{-1} dx = \theta \log x|_{\theta}^{\infty} = \infty$. Thus the method of moments estimator of θ does not exist.

Problem 1.3 Let $X_1, ..., X_n$ be iid with one of two pdfs. If $\theta = 0$, then

$$f(x|\theta) = \begin{cases} 1, & \text{if } 0 < x < 1 \\ 0, & \text{otherwise} \end{cases}$$

$$f(x|\theta) = \begin{cases} \frac{1}{2\sqrt{x}}, & \text{if } 0 < x < 1\\ 0, & \text{otherwise} \end{cases}$$

Find the MLE of θ .

Solution.

$$L(0|x) = \begin{cases} 1, & \text{if } 0 < x_i < 1 \\ 0, & \text{otherwise} \end{cases}$$

while if $\theta = 1$, then

$$L(1|x) = \begin{cases} \prod \frac{1}{2\sqrt{x}}, & \text{if } 0 < x_i < 1\\ 0, & \text{otherwise} \end{cases}$$

Thus,

$$\hat{\theta}_{MLE} = \begin{cases} 1, & \text{if } \prod \frac{1}{2\sqrt{x}} \\ 0, & \text{otherwise} \end{cases}$$



2.1 Sufficiency Principle

Definition 2.1.1 sufficient T(X) is **sufficient** for θ if the distribution of $X|T(\underline{X})$ does not depend on θ .

Theorem 2.1.1 Factorization Theorem If we have $\underline{X} \sim f_{\theta}(\underline{x})$ then T is sufficient iff we can write f_{θ} as

$$g(T(x), \theta)h(x)$$

for some g and h.

Example 2.1 $X_1, \dots, X_n \stackrel{iid}{\sim} \operatorname{Bern}(\theta)$ pdf: $\theta^x (1 - \theta)^{1 - x}$ Joint pdf = $f_{\theta}(x) = \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}$ So, $T(\underline{x}) = \sum x_i$

Q: What is the connection to Section 7.3?

A:

$$Var_{\theta}(W(\underline{X})) = Var[E(W(X)|T(X))] + E[Var(W(X)|T(X))]$$

$$\geq Var[E(W(X)|T(X))]$$

Note: $E\theta \{E_{\theta}[W(X)|T(X)]\} = E_{\theta}(W(X))$

 $E_{\theta}[W(X)|T(X)]$ does not depend on θ . So this is a legitimate estimator since T(X) is sufficient and its varience is smaller than any estimator with same mean.

General Idea of Sufficiency: If T(X) is sufficient for θ , then all information in X about θ is captured in T(X).

More technically, conditioning T(X), the remaining randomness does not depend on θ .



Know Thm 6.2.6 - Factorization Theorem

Definition 2.1.2 minimal sufficient T(X) is **minimal sufficient** if

• is a function of any sufficient statistic

Note: "T(X) is a function S(X)" means that S(x) = S(y) implies T(x) = T(y).

■ **Example 2.2** $X_1, ..., X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$. Find possible $T(\mu, \sigma^2)$.

$$f(\underline{x}) = \prod (2\pi\sigma^2)^{\frac{1}{2}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$
$$= (2\pi\sigma^2)^{\frac{-n}{2}} \exp\left\{\frac{-\sum x_i^2 + 2\mu\sum x_i - n\mu^2}{2\sigma^2}\right\}$$

Therefore our pair of sufficient statistics are $\sum x_i^2$ and $\sum x_i$. But are the minimal?

Theorem 2.1.2 6.2.13 $T(\underline{x}) = T(\underline{x})$ iff $\frac{f_{\theta}(\underline{x})}{f_{\theta}(\underline{y})}$ does not depend on θ , implying $T(\underline{x})$ is minimal sufficient statistic.

Exponential Families

Definition 2.2.1 Exponential Family An exponential family of densities is the set of densi-

$$f_{\underline{\theta}}(\underline{x}) = h(\underline{x})c(\underline{\theta}) \exp \left\{ \sum w_i(\underline{\theta})t_i(\underline{x}) \right\}$$

As θ takes values in some set Θ , this equation may be rewritten as

$$f_{\underline{\eta}}(\underline{x}) = \exp\left\{\underline{\eta}^T \underline{T}(\underline{x}) - A(\underline{\eta})\right\} h(\underline{x}) I_B(\underline{x})$$

The latter equation is called the canonical/natural parameterization.

Example 2.3 Deriving the canonical parameterization for a Normal distribution.

$$\begin{split} f_{(\mu,\sigma^2)} &= k \frac{1}{\sqrt{2\sigma^2}} \exp\left\{ \frac{-1}{2\sigma^2} (x - \mu)^2 \right\} \\ &= k \exp\left\{ \frac{1}{2\sigma^2} (-x^2) + \frac{\mu}{\sigma^2} (x) - \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log((2\sigma^2)) \right\} \\ &= k \exp\left\{ \eta_1 (-x^2) + \eta_2 (x) - \frac{\eta_2^2}{4\eta_1} - \frac{1}{2} \log \eta_1 \right\} \end{split}$$

$$\begin{split} &\eta_1 = \frac{1}{2\sigma^2} \Leftrightarrow \sigma^2 = \frac{1}{2\eta_1} \\ &\eta_2 = \frac{\mu}{2\sigma^2} \Leftrightarrow \mu = \frac{\eta_2}{2\eta_1} \\ &A(\underline{\eta}) = \frac{\eta_2^2}{4\eta_1} - \frac{1}{2}\log\eta_1 \end{split}$$

Cool fact: The partial derivatives of $A(\eta)$ gives the expectations of T. The second partials give the covarience of T.

Observe $T(x) = (-x^2, x)$

$$E(-x^2) = \frac{\delta}{\delta \eta_1} A(\underline{\eta})$$

$$= \frac{-\eta_2^2}{4\eta_1^2} - \frac{1}{2\eta_1}$$

$$= \frac{-\mu^2 * 4\sigma^4}{4\sigma^4} - \frac{2\sigma^2}{2}$$

$$= -\mu^2 - \sigma^2$$

$$E(x) = \frac{\delta}{\delta \eta_2} A(\underline{\eta})$$
$$= \frac{\eta_2}{2\eta_1}$$
$$= \mu$$

$$Var(x^{2}) = \frac{\delta^{2}}{\delta^{2}\eta_{1}}$$

$$= \frac{\eta_{2}^{2}}{2\eta_{1}^{3}} + \frac{1}{2\eta_{2}}$$

$$= 4\mu^{2}\sigma^{2} + 2\sigma^{4}$$

2.2.1 Distribution Families

Poisson

$$f_{\lambda} = e^{-\lambda} \frac{\lambda^{x}}{x!} I\{x \in \mathbb{N}\}$$

$$= \exp\{(\log \lambda)x - \lambda\}$$

$$\eta = \log \lambda$$

$$T(\underline{x}) = x$$

$$A(\eta) = \lambda = e^{\log \lambda} = e^{\eta}$$

Beta

$$\begin{split} f_{\alpha,\beta}(x) &= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) + \Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} I \{ 0 < 1x < 1 \} \\ &= \exp \{ \log \Gamma(\alpha+\beta) - \log \Gamma(\alpha) - \log \Gamma(\beta) + (\alpha-1) \log x + (\beta-1) \log (1-x) I \{ 0 < x < 1 \} \} \\ \eta_1 &= \alpha - 1 \\ \eta_2 &= \beta - 1 \\ T_1(x) &= \log(x) \\ T_2(x) &= \log(1-x) - x\beta + \alpha \log(\beta) - \log(\Gamma(\alpha)) \end{split}$$

Gamma

$$\begin{split} f_{\alpha,\beta}(x) &= \frac{\beta^{\alpha} x^{\alpha - 1}}{\Gamma(\alpha)} \exp\left\{-x\beta\right\} I\{x > 0\} \\ &= \exp\left\{(\alpha - 1)\log(x) - \beta(x) + (\eta_1 + 1)\log(-\eta_2) + \log(\Gamma(\eta_1 + 1))\right\} \\ \eta_1 &= \alpha - 1 \\ \eta_2 &= \beta \\ T_1(x) &= \log(x) \\ T_2(x) &= x \\ A(\eta) &= (\eta_1 + 1)\log(-\eta_2) + \log(\Gamma(\eta_1 + 1)) \end{split}$$

Binomial

$$f_p(x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} I\{0...,n\}$$

$$= \binom{n}{x} \exp\left\{x(\log\frac{p}{1-p}) - (-n\log(1-p))\right\}$$

$$\eta = \log\frac{p}{1-p}$$

$$T(x) = x$$

$$A(\eta) = -n\log(1-p)$$

Why is $E(T(x)) = \frac{\delta}{\delta \eta} A(\eta)$?

Notice: $e^{A(\eta)}$ must be the integral of $\exp\left\{\frac{\eta}{T}\underline{T}(\underline{x}) - A(\underline{\eta})\right\}h(\underline{x})$.

Differentiating on both sides and using the fact that for an expoonential family we may switch integral and derivative, we obtain:

$$\frac{\delta}{\delta \eta} A(\eta) e^{A(\eta)} = e^{A(\eta)} \int \exp\left\{\eta^T T(x)\right\} h(x) [T(x)] e^{-A(\eta)}$$
(2.1)

$$=E(T(x)) (2.2)$$

(2.3)

Recall, $f_{\eta}(x)$ must always integrate to 1!



Read Casella & Berger Chapter 7.2

2.3 Complete Statistic

Recall, minimal statistics contain no redundant or unnessicary information. What might redundant information look like? From T(x) we may construct g_1 and g_2 such that $E_{\theta}(g_1(T)) = E_{\theta}(g_2(T))$ for all θ which we want to avoid.

Definition 2.3.1 complete statistic (T) is obtained only if g=0 gives $E_{\theta}(g(T)) = o$ for all θ .

Theorem 2.3.1 Complete and sufficient inplies minimal statistics.

Theorem 2.3.2 If $f_{\theta}(\underline{x}) \exp \{ \eta^T T(x) = A(\eta) \} h(x)$ is a cononical exponential family of full rank then T is a complete and sufficient statistic.

Definition 2.3.2 Full Rank Full rank means

Parameter space H contains an open set

• T does not satisfy any linear constraint (i.e. the T are linearly independent).

Definition 2.3.3 Open Set Let k be dimention of η . To say that H contains an **open set** means there exists $\begin{cases} \varepsilon > 0 \\ \eta_0 \in H \end{cases}$ such that $B(\eta_0; \varepsilon) \subset H$

Theorem 2.3.3 Basu's Theorem Any complete and sufficient statistic (minimum sufficient) is independent of any ancillary statistic (that is, any statistic whose distribution does not depend on θ).

2.3.1 Sufficiency and Unbiasness

Suppose $E_{\theta}(W(X)) = \theta$ for all θ , thus W(X) is unbiased. Then $Var[E(W(X)|T)] \leq Var(W(X))$. But what if T is not sufficient? Wel...then you won't get an estimator.

■ Example 2.4 Suppose that $X \sim Unif(0, \theta)$. Find the UMVUE of $\sin \theta$.

Solution.

pdf:
$$f(x) = \frac{1}{b-a} = \frac{1}{\theta}$$

$$E(X) = \frac{\theta}{2}$$

Thus 2X is unbiased for θ . But is $\sin(2X)$ unbiased for $\sin(\theta)$? NO!

What about cos(2X)?

$$E(\cos(2X)) = \int_{1}^{\theta} \cos(2X) \frac{1}{\theta} dx$$
$$= \frac{1}{\theta} \left(\frac{\sin(2X)}{2} \right) \Big|_{0}^{\theta}$$
$$= \frac{\sin(2\theta)}{2\theta}$$

No! Blarg. Let's bust out some calculus.

$$\int_{0}^{\theta} \mathbf{estimator}(\mathbf{x}) dx = \theta \sin(\theta)$$

$$\mathbf{estimator} = \sin \theta + \theta \cos(\theta)$$

$$= x \sin(x) + x \cos(x)$$

Theorem 2.3.4 Rae-Blackwell [7.3.17] Let W(X) be an extimator with $E_{\theta}(W(x)) = \tau(\theta)$. If T(X) is a sufficient statistic, then an alternative, unbiased estimator of $\tau(\theta)$ whose varience is

uniformly not worse than that of W(X) is

$$\phi(X) = E_{\theta}(W(X)|T(X))$$

Proof: Based on conditional arguments.

$$E(W) = E(E(W(X)|T(X)))$$

Var(W)=Var(Cond Exp)+E(Cond Var)

Theorem 2.3.5 7.3.19 If W(X) is UMVUE for $\tau(\theta)$ then it is the unquie UMVUE, aka the best unbiased estimator for all theta.

Theorem 2.3.6 7.3.20 If $E_{\theta}(W(X)) = \tau(\theta)$, then W(X) is UMVUE for $\tau(\theta)$ iff W(X) is uncorrelated with any S(X) such that $E_{\theta}(S(X)) = 0$.

Theorem 2.3.7 7.2.23 If T(X) is a complete and sufficient statistic, then $\phi(T)$ is the unique UMVUE of its expectation (see Rao-Blackwell).

Definition 2.3.4 Jenson's Inequality If h(x) is a convex function then

$$h(E(x)) \le E(h(x))$$

for any random variable x.

2.4 Methods of Finding Estimators

Setup: $x_1, \ldots, x_n \stackrel{iid}{\sim} f_{\theta}, \theta \in \Theta$

Goal: Estimate θ , or some function $\tau(\theta)$

Definition 2.4.1 — Point Estimator. A point estimator $T(x_1, \ldots, x_n)$ is a function of the random sample. T should not depend on θ .



Estimator is defined above, while an estimate is a fixed value. $(T(x_1 = x_1, \dots, x_n = x_n))$

2.4.1 **Method of Moments**

This method is one of the oldest, simplest estimators developed by Karl Pearson.

$$\begin{aligned} \textbf{Construct:} & \begin{cases} \frac{1}{n} \sum X_i = E_{\theta}(x) = \int_{\theta} x f_{\theta}(x) dx \\ \frac{1}{n} \sum X_i^2 = E_{\theta}(x^2) = \int_{\theta} x^2 f_{\theta}(x) dx \\ \vdots \end{cases} \\ & \text{So, } \hat{\theta}_M M \text{ depends only on } (\frac{1}{n} \sum X_i = E_{\theta}(x), \dots, \frac{1}{n} \sum X_i = E_{\theta}(x^k)). \end{cases}$$

■ Example 2.5
$$x_1, ..., x_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

 μ, σ^2 unknown: $\Theta = \mathbb{R}x\mathbb{R}^+$

$$\frac{1}{n}\sum X_i = \mu \implies \mu_{MM}^2 = \frac{1}{n}\sum X_i \triangleq \bar{X}$$
$$\frac{1}{n}\sum X_i^2 = \sigma^2 + \mu^2 \implies \sigma_M M^2 = \frac{1}{n}\sum (X_i - \bar{X})^2$$

■ Example 2.6 $x_1, ..., x_n \stackrel{iid}{\sim} \text{Bin}(k, p)$ where k, p are unknown: $\Theta : \mathbb{N} \times [0, 1]$.

$$\frac{1}{n} \sum X_i = kp \triangleq \bar{X}$$

$$\frac{1}{n} \sum X_i^2 = kp(1-p) + k^2 p^2 = \bar{X}(1-p) + \bar{X}^2$$

R

Read about Satterthwait approximation

2.4.2 Maximum Likelihood Estimation

MLE makes data the most likely under the model specified. The quality of MLE depends on the quality of the model. Unlike Method of Moments, MLE does not lie outside of Θ .

Invariance Principle: if $\hat{\theta}$ is MLE for θ , then $\tau(\hat{\theta})$ is MLE for $\tau(\theta)$.

Roughly speaking, MLEs enjoy a nice property called **asymptotic efficiency**. They acheive the smallest possible varience in large sample. Similar to Lower Bound of Cramer-Rao as n -> infinity.

Definition 2.4.2 Score Function Suppose that throughout $X_1, \ldots, X_n \sim f_{\theta}(x)$ the **score function** is $\sum \frac{\delta}{\delta \theta} \log f_{\theta}(x_i)$.

Take a first-order Taylor expansion of $\ell(\theta)$ around true value (θ) .

$$\ell'(\hat{\theta}) \simeq \ell'(\theta_0) + (\hat{\theta} - \theta_0)\ell''(\theta_0)$$

$$\hat{\theta} - \theta_0 \simeq [\ell''(\theta_0)^{-1}][\ell'(\hat{\theta}) - \ell'(\theta_0)]$$

$$= [\frac{-1}{n}\ell''(\theta_0)^{-1}][\frac{1}{n}\ell'(\theta_0)]$$

$$\frac{(\hat{\theta} - \theta_0)}{\sqrt{n}} \simeq \frac{\frac{-1}{\sqrt{n}}\ell'(\theta_0)}{\frac{-1}{n}\ell''(\theta_0)}$$

$$\frac{1}{n}\ell'(\theta_0) = \frac{1}{n}\sum \frac{\delta}{\delta\theta} \log f_{\theta}(x)$$

$$\frac{-1}{n}\ell''(\theta_0) = \frac{1}{n}\sum \frac{\delta^2}{\delta\theta^2} \log f_{\theta}(x)|_{\theta = \theta_0}$$

Thus,
$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \sim N(0, [I(\boldsymbol{\theta}_0)]^{-1})$$

Finding the Maxima: $\theta \in \mathbb{R}^k$

Suppose L is differentialbe, then the only possible candidtes for $\hat{\theta_{MLE}}$ are the stationary points of L which satisfy $\frac{\delta L}{\delta \theta_i} = 0, i = 1, \dots, k$, which is a necessary condition for optimality. Suppose L is twice differentiable

$$\mathbf{H} \triangleq \begin{bmatrix} \frac{\delta^2 L}{\delta \theta_1^2} & \frac{\delta^2 L}{\delta \theta_1 \delta \theta_2} & \cdots & \frac{\delta^2 L}{\delta \theta_1 \delta \theta_k} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\delta^2 L}{\delta \theta_k \theta_1} & \cdots & \cdots & \frac{\delta^2 L}{\delta \theta_k^2} \end{bmatrix}$$

2.5 Major/Minorization Mini/Maximization

Ever EM Algorhythm is special case of a MiMax. EM = Expectation Maximization



Read 7.2.4 in Casella & Berger

2.5.1 How EM Algorhythms Work

• Choose a starting parameter, θ .

E-Step Construct the (minimizing) function.

- Maximize this function of θ . The maximizer will be next value of *theta*₀. Call it θ_1 .
- Return to (b) as long as we haven't converged.

Given μ_0 we find after one EM iteration that

$$\mu_1 = \frac{\sum Y_i}{n} + \mu_0 (1 - \frac{\mu}{n})$$
$$= \mu_0 (1 - \frac{\mu}{n}) + \frac{\sum Y_i}{\mu} (\frac{\mu}{n})$$

Thus, μ_1 always takes us $\frac{\mu}{n}$ of the way to the final answer at each iteration. This type of convergence is called **linear convergence**, which is a characteristic of EM algorhythems. The linear rate of convergence is governed by "amount of missingness" and is considered slow compared to other similar methods (e.g. Newton-Raphson is optimization method that enjoys quadradic convergence). However, EM algorhithem tend to trade more iterations for simpler iterations

Suppose we wish to maximize a fuction that is a product of sums or integrals. Taking the log gives a sum of logs of sums/integrals:

$$\sum_{a}^{A} \log \sum_{b}^{B} s_{ab}(\theta)$$

To simplify, ignore the sumation over A and take $f(\theta) = \log \sum s_{ab}(\theta)$. Fix a θ_0 . Define

$$W_{ob} \equiv \frac{S_b(\theta_0)}{\sum_{c_a}^c S(\theta_0)}$$

Goal: Maximize $f(\theta) = \log \sum S_b(\theta)$ with fixed θ_0 .

Claim: Define $Q_0(\theta) \equiv \sum_{b} W_{ob} * \log S_b(\theta)$.

Then $f(\theta) - f(\theta_0) \ge Q_0(\theta) - Q_0(\theta_0)$.

To verify this claim, notice that Jenson's inequality says

$$E(\log(\bullet)) \ge \log E(\bullet)$$

$$\begin{aligned} Q_0(\theta) - Q_b(\theta_0) &= \sum^B W_{ob} \log(\frac{S_b(\theta)}{S_b(\theta_0)}) \\ &\leq \log \sum^B W_{ob} \frac{S_b(\theta)}{S_b(\theta_0)} \\ &\leq \log \sum^B \frac{S_b(\theta_0)}{\sum S_c(\theta_0)} \frac{S_b(\theta)}{S_b(\theta_0)} \\ &f(\theta) - f(\theta_0) = \log \frac{S_b(\theta)}{S_c(\theta_0)} \end{aligned}$$

■ Example 2.7 $X \sim f_{\theta}(x) = \lambda f_{\xi_1}(x) + (1 - \lambda) f_{\xi_2}(x)$, $\theta = (\lambda, \xi)$ Intuition: To generate X according to mixture density flip coin with Heads probability of λ . If H, generate $X \sim f_{\xi_1}(x)$. If T, generage $X \sim f_{\xi_2}$.

Write down the observed data likelihood, $\ell(\lambda, p_1, p_2)$. For a sample size, n, form a mixture of Binom(m, p_1) and Binom(m, p_2).

$$\sum \log \left(\lambda \binom{m}{\xi} p_1^{\xi_1} (1 - p_1)^{m - \xi_1} + (1 - \lambda) \binom{m}{\xi_2} p_2^{\xi_2} (1 - p_2)^{m - \xi_2} \right)$$

Now impliment EM algorithm.

Suppose that "complete data" consists of $x_{obs} & x_{miss}$. Using previous notation,

$$S_{b}(\theta) = P_{\theta} (X_{obs} = x_{obs}, X_{miss} = b)$$

$$\Rightarrow \sum_{b}^{B} S_{b}(\theta) = P_{\theta} (X_{obs} = x_{obs})$$

$$W_{ob} = \frac{P_{\theta_0} \left(X_{obs} = x_{obs}, X_{miss} = b \right)}{P_{\theta} \left(X_{obs} = x_{obs} \right)}$$
$$= P_{\theta_0} \left(X_{miss} = b | X_{obs} = x_{obs} \right)$$

$$Q_0(\theta) = E_{\theta_0}[\log P_{\theta}(X_{obs} = x_{obs}, X_{miss} = b) | X_{obs} = x_{obs}]$$

2.6 Bayes Estimators

Theorem 2.6.1 Bayes' Theorem

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

LHS: posterior density of the parameters given data. This is what we want to know. The RHS inclues

likelihood

- prior density
- marginal density of x

Proposition 2.6.2 Dave's Bayes Way posterior ∝ prior*likelihood

Definition 2.6.1 Conjugate Family If the prior is chosen from a family of distributions in such a way that the posterior is also a member, then we call this family of priors a **conjugate family**.

Definition 2.6.2 Bayes Risk Recall we previously define risk as a function of a parameter and estimator - E(Loss). Now that θ can be random, we may define **Bayes risk** as the average risk.

$$\int R(\theta, \delta) \pi(\theta) d\theta = E(R(\theta, \delta)) = E(L(\theta, \delta(\underline{X})))$$

Definition 2.6.3 Bayes Estimator For a give loss function, we define the **Bayes estimator** as the minimizer of the Rayes Risk.

Note: if δ^* minimizes $E(L(\theta, \delta(\underline{X}))|\underline{X} = x)$ for (almost) all x then δ^* is the Bayes Estimator. We can conclude that under square error loss, the posterior mean is the Bayes Estimator. Wrap up: How to choose prior? Even choosing family of distribution is hard. Considering restraints and conjugacy can help.

- Use data to estimate prio parameters
- Assign priors on our priors "hyper priors". Priorception.
- Use a Jeffery's prior: reparameterize

2.6.1 Model Based Clustering



Definition 3.0.4 hypothesis testing We see that a hypothesis testing procedure is a rule that partitions the sampe space into we will accept (fail to reject) H_0 as true or not.

Let's consider a simple null hypothesis, i.e. a hypothesis of the form:

$$H_0: \theta = \theta_0$$

We will consider three different ways to determine whether $heta_{true} = heta_0$

- 1. Wald Test $theta \theta_0$ consider if this is large relative to the distribution. $\hat{\theta}$ should have if H_0 true. Recall, $\sqrt{n}(\hat{\theta} \theta_0) \sim N(0, I^{-1}(\theta_0))$.
- 2. **Likelihood Ratio Test** Consider $\ell(\hat{\theta}) \ell(\theta_0)$ and determine whether this is large relative to χ^2 distribution that it will approximately have under H_0 .
- 3. **Rao Score Test** Consider $\nabla \ell(\theta_0)$ and determine if it is too far from 0 relative to it's true approximate distribution if H_0 true.



Read 8.2.1 and 8.2.2

Theorem 3.0.3 If T(X) is sufficient for the family of distributions from which X is drawn and $\delta(X)$ is an unbiased estimater



Books Articles



Paragraphs of Text, 5