

# **CHAPTER 5.**

## **BAYES METHODS AND ELEMENTARY DECISION THEORY**

1. **Elementary Decision Theory**
2. **Structure of the risk body: the finite case**
3. **The finite case: relations between Bayes minimax, admissibility**
4. **Posterior distributions**
5. **Finding Bayes rules**
6. **Finding minimax rules**
7. **Admissibility and Inadmissibility**
8. **Asymptotic theory of Bayes estimators**
9. **Nonparametric Bayes estimation**

## 1. ELEMENTARY DECISION THEORY

### Notation 1.1.

Let

- $\Theta \equiv$  the states of nature.
- $\mathbf{A} \equiv$  the action space
- $\mathbf{X} \equiv$  the sample space of a random variable  $X$  with distribution  $P_\theta$
- $P : \mathbf{X} \times \Theta \rightarrow [0, 1]$ ;  $P_\theta(X = x) =$  probability of observing  $X = x$  when  $\theta$  is true.
- $L : \Theta \times \mathbf{A} \rightarrow R^+$ ; a loss function.
- $d : \mathbf{A} \times \mathbf{X} \rightarrow [0, 1]$ ; a decision function function.  
 $d(a, x) = d(a | x) =$  probability of action  $a$  when  $X = x$  is observed.
- $\mathbf{D} = \{\text{all decision functions}\}$ .

**Definition 1.2.** The *risk* for rule  $d \in \mathbf{D}$  when  $\theta \in \Theta$  is true is

$$\begin{aligned} R(\theta, d) &= E_\theta L(\theta, d(\cdot | X)) \\ &= \int_{\mathbf{X}} \int_{\mathbf{A}} L(\theta, a) d(da | X = x) P_\theta(dx) \\ &= \sum_{i=1}^k L(\theta, a_i) \left\{ \sum_{j=1}^m d(a_i, x_j) P_\theta(X = x_j) \right\} \quad \text{in the discrete case.} \end{aligned}$$

We will call  $R : \Theta \times \mathbf{D} \rightarrow R^+$  the *risk function*.

**Definition 1.3.** A decision rule  $d$  is *inadmissible* if there is a rule  $d'$  such that  $R(\theta, d') \leq R(\theta, d)$  for all  $\theta$  and  $R(\theta, d') < R(\theta, d)$  for some  $\theta$ . A decision rule  $d$  is *admissible* if it is not inadmissible.

**Example 1.4. Hypothesis testing.**  $\mathbf{A} = \{0, 1\}$ ,  $\Theta = \Theta_0 \cup \Theta_1$

$$L(\theta, 0) = l_0 1_{\Theta_1}(\theta), \quad L(\theta, 1) = l_1 1_{\Theta_0}(\theta).$$

Since  $\mathbf{A}$  contains just two points,  $d(1 | x) = 1 - d(0 | x) \equiv \phi(x)$ ; and then

$$P_\theta(\text{accept } H_0) = E_\theta d(0 | X) = E_\theta(1 - \phi(X)),$$

$$P_\theta(\text{reject } H_0) = E_\theta d(1 | X) = E_\theta \phi(X),$$

and

$$R(\theta, d) = l_1 E_\theta \phi(X) 1_{\Theta_0}(\theta) + l_0 E_\theta (1 - \phi(X)) 1_{\Theta_1}(\theta).$$

The classical Neyman - Pearson hypothesis testing philosophy bounds  $l_1 \alpha \equiv \sup_{\theta \in \Theta_0} R(\theta, d)$  and tries to minimize

$$I_0(1 - \beta_d(\theta)) \equiv R(\theta, d)$$

for each  $\theta \in \Theta_1$ .

**B. Estimation.**  $\mathbf{A} = \Theta$ ; typically  $\Theta = R^s$  for some  $s$  (but sometimes it is useful to consider  $\Theta = \mathbf{T}$  where  $(T, |\cdot|)$  is a metric space). A typical loss function (often used for mathematical simplicity more than anything else) is

$$L(\theta, a) = K|\theta - a|^2 \quad \text{for some } K.$$

Then the risk of a rule  $d$  is

$$\begin{aligned} R(\theta, d) &= K E_\theta |\theta - d(\cdot | X)|^2 \\ &= K E_\theta |\theta - d(X)|^2 \quad \text{for a non-randomized } d \\ &= K \{ \text{Var}_\theta d(X) + [\text{bias}_\theta(d(X))]^2 \} \quad \text{when } s = 1. \end{aligned}$$

**Definition 1.5.** A decision rule  $d_M$  is *minimax* if

$$\inf_{d \in \mathbf{D}} \sup_{\theta \in \Theta} R(\theta, d) = \sup_{\theta \in \Theta} R(\theta, d_M).$$

**Definition 1.6.** A probability distribution  $\Lambda$  over  $\Theta$  is called a *prior distribution*.

**Definition 1.7.** For any given prior  $\Lambda$  and  $d \in \mathbf{D}$ , the *Bayes risk of  $d$  with respect to  $\Lambda$*  is

$$\begin{aligned} \mathcal{R}(\Lambda, d) &= \int_{\Theta} R(\theta, d) d\Lambda(\theta) \\ &= \sum_{i=1}^l R(\theta_i, d) \lambda_i \quad \text{if } \Theta = \{\theta_1, \dots, \theta_l\} \end{aligned}$$

**Definition 1.8.** A *Bayes decision rule with respect to  $\Lambda$* ,  $d_\Lambda$ , is any rule satisfying

$$\mathcal{R}(\Lambda, d_\Lambda) = \inf_{d \in \mathbf{D}} \mathcal{R}(\Lambda, d) = \text{Bayes risk}.$$

**Example 1.9.**  $\Theta = \{1, 2\}$ .

Urn 1: 10 red balls, 20 blue balls, 70 green balls.

Urn 2: 40 red balls, 40 blue balls, 20 green balls.

One ball is drawn from one of the two urns. Problem: decide which urn the ball came from if the losses  $L(\theta, a)$  are given by

$\theta/a$	1	2
1	0	10
2	6	0

Let  $d = (d_R, d_B, d_G)$  with  $d_x =$  probability of choosing urn 1 if color  $X = x$  is observed. Then,

$$\begin{aligned} R(1, d) &= 10 P_1(\text{action } 2) \\ &= 10\{.1(1 - d_R) + .2(1 - d_B) + .7(1 - d_G)\} \end{aligned}$$

and, similarly,

$$R(2, d) = 6\{.4d_R + .4d_B + .2d_G\}.$$

If the prior distribution on the urns is given by  $\lambda_1 = \lambda$  and  $\lambda_2 = 1 - \lambda$ , then the Bayes risk is

$$\mathbb{R}(\Lambda, d) = 10\lambda + (2.4 - 3.4\lambda)d_R + (2.4 - 4.4\lambda)d_B + (1.2 - 8.2\lambda)d_G.$$

This is minimized by choosing  $d_x = 1$  if its coefficient is negative, 0 if its coefficient is positive. For example, if  $\lambda = 1/2$ , then the Bayes risk equals

$$5 + .7d_R + .2d_B - 2.9d_G,$$

which is minimized by  $d_R = d_B = 0$ ,  $d_G = 1$ ; i.e. the Bayes rule  $d_\Lambda$  with respect to  $\lambda = 1/2$  is  $d_\Lambda = (0, 0, 1)$ . Note that the Bayes rule is in fact a non-randomized rule. This give us the Bayes risk for  $\lambda = 1/2$  as  $\mathbb{R}(1/2, d_B) = 2.1$ .

I claim that the minimax rule is  $d_M = (0, 9/22, 1)$ , which is a randomization of the two nonrandom rules  $d_2 = (0, 0, 1)$  and  $d_7 = (0, 1, 1)$ . This is easily confirmed by computing  $R(1, d_M) = 263/110 = R(2, d_M)$ .

Here is a table giving all of the nonrandomized rules:

$X$	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$	$d_7$	$d_8$
$R$	0	0	0	1	1	1	0	1
$B$	0	0	1	0	1	0	1	1
$G$	0	1	0	0	0	1	1	1
$R(1, d)$	10	3	8	9	7	2	1	0
$R(2, d)$	0	1.2	2.4	2.4	4.8	3.6	3.6	6

## 2. STRUCTURE OF THE RISK BODY: THE FINITE CASE

Suppose that

$$\Theta = \{\theta_1, \dots, \theta_l\}, \quad \mathbf{X} = \{x_1, \dots, x_m\}, \quad \mathbf{A} = \{a_1, \dots, a_k\}.$$

**Definition 2.1.**  $A \subset R^l$  is *convex* if, for all  $\lambda \in [0, 1]$   $x, y \in A$ ,  $\lambda x + (1 - \lambda)y \in A$ .

**Lemma 2.2.** Let  $\lambda = (\lambda_1, \dots, \lambda_t)$  be a probability distribution,  $\lambda_i \geq 0$ ,  $\sum_{i=1}^t \lambda_i = 1$ ). Then for any  $d_1, \dots, d_t \in \mathbf{D}$ ,  $\sum_{i=1}^t \lambda_i d_i \in \mathbf{D}$ .

**Proof.** set  $d(a_l, x_j) = \sum_{i=1}^t \lambda_i d_i(a_l, x_j)$ . Then  $d(a_l, x_j) \geq 0$  and

$$\sum_{l=1}^k d(a_l, x_j) = \sum_{i=1}^t \lambda_i \sum_{l=1}^k d_i(a_l, x_j) = 1.$$

□

We call

$$\mathbf{R} = \{(R(\theta_1, d), \dots, R(\theta_l, d)) \in R^{+l} : d \in \mathbf{D}\}$$

the *risk body*.

**Theorem 2.3.** The risk body is convex.

**Proof.** Let  $d_1, d_2 \in \mathbf{D}$  denote the rules corresponding to  $R_i = (R(\theta_1, d_i), \dots, R(\theta_l, d_i))$ ,  $i = 1, 2$ , and set  $d = \lambda d_1 + (1 - \lambda)d_2 \in \mathbf{D}$  by the lemma. Then

$$\begin{aligned} & \lambda R(\theta_l, d_1) + (1 - \lambda)R(\theta_l, d_2) \\ &= \sum_{i=1}^k \sum_{j=1}^m L(\theta_l, a_i) \{\lambda d_1(a_i, x_j) + (1 - \lambda)d_2(a_i, x_j)\} p_{\theta_l}(x_j) \\ &= R(\theta_l, d), \end{aligned}$$

so  $d \in \mathbf{D}$  has risk point  $\lambda R_1 + (1 - \lambda)R_2 \in \mathbf{R}$ . □

**Theorem 2.4.** Every  $d \in \mathbf{D}$  may be expressed as a convex combination of the nonrandomized rules.

**Proof.** Any  $d \in \mathbf{D}$  may be represented as a  $k \times m$  matrix of  $\geq 0$  real numbers whose columns add to 1. The nonrandomized rules are those whose entries are 0's and 1's. Set  $d(a_i, x_j) = d_{ij}$ ,  $d = (d_{ij})$ ,  $1 \leq i \leq k$ ,  $1 \leq j \leq m$ .

There are  $k^m$  nonrandom decision rules: call them  $\delta^{(r)} = (\delta_{ij}^{(r)})$ ,  $1 \leq r \leq k^m$ . Given  $d \in \mathbf{D}$ , we want to find  $\lambda_l \geq 0$ ,  $\sum_{l=1}^{k^m} \lambda_l = 1$ , so that

$$\sum_{l=1}^{k^m} \lambda_l \delta^{(r)} = d .$$

Claim:  $\lambda_r = \prod_{j=1}^m \prod_{i=1}^k d_{ij}^{\delta_{ij}^{(r)}} = \prod_{j=1}^m \left\{ \sum_{i=1}^k d_{ij} \delta_{ij}^{(r)} \right\}$  work.

**Proof.** Induction on  $m$  . □

### 3. THE FINITE CASE: RELATIONS BETWEEN BAYES, MINIMAX, ADMISSIBILITY

This section continues our examination of the special, but illuminating, case of a finite set  $\Theta$ . In this case we can prove a number of results about Bayes and minimax rules and connections between them which carry over to more general problems with some appropriate reformulation.

**Theorem 3.1.** A minimax rule always exists.

**Proof.** If  $0 \notin \mathbf{R}$ , then for some  $c > 0$  the cube with all sides of length  $c$  in the positive quadrant does not intersect  $\mathbf{R}$ . Let  $c$  increase until the square intersects  $\mathbf{R}$ ; call this number  $c_0$ . Any decision rule with a risk point which intersects the  $c_0$ -square is minimax. (Note the the equation for the boundary of a cube is  $\max_i R(\theta_i, d) = \text{constant}$ .  $\square$ )

**Theorem 3.2.** If  $\lambda = (\lambda_1, \dots, \lambda_l)$  is a prior, then Bayes decision rules have risk points on the hyperplane  $\{x \in \mathbf{R}^l : \sum_i \lambda_i x_i = c^*\}$  where

$$c^* = \inf\{c \geq 0 : \text{the plane determined by } \sum_i \lambda_i x_i = c \text{ intersects } \mathbf{R}\}.$$

**Proof.**  $\mathcal{R}(\Lambda, d) = \sum \lambda_i R(\theta_i, d)$  is the Bayes risk and we want to minimize it.  $\square$

**Lemma 3.3.** An admissible rule is a Bayes rule for some prior  $\lambda$ .

**Proof.** Picture! □

Note that *not every* Bayes rule is admissible.

**Theorem 3.4.** Suppose that  $d_0$  is Bayes w.r.t.  $\lambda = (\lambda_1, \dots, \lambda_l)$  and  $\lambda_i > 0$ ,  $i = 1, \dots, l$ . Then  $d_0$  is admissible.

**Proof.** Suppose that  $d_0$  is *not* admissible; then there is a rule  $d$  better than  $d_0$ : i.e.

$$R(\theta_i, d) < R(\theta_i, d_0) \quad \text{for all } i$$

with  $<$  for some  $i$ . Since  $\lambda_i > 0$  for  $i = 1, \dots, l$ , contradicting the fact that  $d_0$  is Bayes with respect to  $\lambda$ . □

**Theorem 3.5.** If  $d_B \in \mathbf{D}$  is Bayes for  $\lambda$  and it has *constant risk*, then  $d_B$  is minimax.

**Proof.** Let  $r_0$  be the constant risk. Assume  $d_B$  is *not* minimax and let  $d_M$  be the minimax rule (which exists). Then

$$(a) \quad R(\theta_i, d_M) \leq \max_{1 \leq j \leq l} R(\theta_j, d_M) < \max_{1 \leq j \leq l} R(\theta_j, d_B) = r_0$$

and

$$(b) \quad \min_{d \in \mathbf{D}} \mathbf{R}(\lambda, d) = \mathbf{R}(\lambda, d_B) = \sum \lambda_i R(\theta_i, d_B) = r_0.$$

But (a) yields

$$(c) \quad \mathbf{R}(\lambda, d_M) = \sum \lambda_i R(\theta_i, d_M) < \sum \lambda_i r_0 = r_0$$

which contradicts (b). □



#### 4. POSTERIOR DISTRIBUTIONS

Now we will examine Bayes rules more generally.

**Definition 4.1.** If  $\theta \sim \Lambda$ , a prior distribution over  $\Theta$ , then the conditional distribution of  $\theta$  given  $X$ ,  $P(\theta \in A | X)$ ,  $A \in \mathbf{B}(\Theta)$ , is called the *posterior* distribution of  $\theta$ . Write  $\Lambda(\theta | x) = P(\theta \leq \theta | X = x)$  for the conditional distribution function if  $\Theta$  is Euclidean.

If  $\Lambda$  has density  $\lambda$  with respect to  $\nu$  and  $P_\theta$  has density  $p_\theta$  with respect to  $\mu$ , then

$$\lambda(\theta | x) = \frac{p(x | \theta)\lambda(\theta)}{\int_{\Theta} p(x | \theta)\lambda(\theta) d\nu(\theta)} = \frac{p(x | \theta)\lambda(\theta)}{p(x)}$$

is the *posterior density* of  $\theta$  given  $X = x$ .

**Definition 4.2.** Suppose  $X$  has density  $p(x | \theta)$ , and  $\lambda = \lambda(\cdot)$  is a prior density,  $\lambda \in \mathbf{P}_\Lambda$ . If the posterior density  $\lambda(\cdot | x)$  has the same form (i.e.  $\lambda(\cdot | x) \in \mathbf{P}_\Lambda$  for a.e.  $x$ ), then  $\lambda$  is said to be a *conjugate prior* for  $p(\cdot | \theta)$ .

#### Examples 4.3.

**Example 4.3.1.**  $(X | \theta = \theta) \sim \text{Poisson}(\theta)$ ,  $\theta \sim \Gamma(\alpha, \beta)$ . Then

$$\begin{aligned} \lambda(\theta | x) &= \frac{p(x | \theta)\lambda(\theta)}{p(x)} \propto p(x | \theta)\lambda(\theta) \\ &= e^{-\theta} \frac{\theta^x}{x!} \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta} \\ &= \frac{\beta^\alpha}{x! \Gamma(\alpha)} \theta^{\alpha+x-1} e^{-(\beta+1)\theta}, \end{aligned}$$

so  $(\theta | X = x) \sim \Gamma(\alpha + x, \beta + 1)$ , and gamma is a conjugate prior for Poisson.

**Example 4.3.2.** (Example 1.3, Lehmann, TPE, page 243).  $(X | \theta = \theta) \sim N(\theta, \sigma^2)$ ,  $\theta \sim N(\mu, \tau^2)$ . Then

$$(\theta | X = x) \sim N\left(\frac{\mu/\tau^2 + x/\sigma^2}{1/\tau^2 + 1/\sigma^2}, \frac{1}{1/\tau^2 + 1/\sigma^2}\right).$$

Note that if  $X$  is replaced by  $X_1, \dots, X_n$  iid  $N(\theta, \sigma^2)$ , then by sufficiency,

$$\begin{aligned} (\theta | \underline{X} = \underline{x}) &= (\theta | \bar{X} = \bar{x}) \\ &\sim N\left(\frac{n/\sigma^2}{1/\tau^2 + n/\sigma^2} \bar{x} + \frac{1/\tau^2}{1/\tau^2 + n/\sigma^2} \mu, \frac{1}{1/\tau^2 + n/\sigma^2}\right). \end{aligned}$$

**Example 4.3.3.** If  $(X | \theta = \theta) \sim \text{Binomial}(n, \theta)$ ,  $\theta \sim \text{B}(\alpha, \beta)$ , then

$$\begin{aligned} \lambda(\theta | x) &= \frac{p(x | \theta) \lambda(\theta)}{p(x)} \propto p(x | \theta) \lambda(\theta) \\ &= \binom{n}{x} \theta^x (1 - \theta)^{n-x} \theta^{\alpha-1} (1 - \theta)^{\beta-1} B_{\alpha, \beta} \\ &= \theta^{\alpha+x-1} (1 - \theta)^{\beta+n-x-1} B_{\alpha, \beta} \binom{n}{x}, \end{aligned}$$

so  $(\theta | X = x) \sim B(\alpha + x, \beta + n - x)$ .

**Example 4.3.4.** Multinomial / Dirichlet. Suppose that  $(X | \theta = \theta) \sim \text{Multinomial}_k(n, \underline{\theta})$  and  $\theta \sim \text{Dirichlet}(\underline{\alpha})$ ; i.e.

$$\lambda(\theta) = \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \dots \theta_k^{\alpha_k-1}$$

for  $\theta_i \geq 0$ ,  $\sum \theta_i = 1$ ,  $\alpha_i > 0$ . Then  $(\underline{\theta} | \underline{X} = \underline{x}) \sim \text{Dirichlet}(\underline{\alpha} + \underline{x})$ ,

$$E(\theta) = \frac{\underline{\alpha}}{\sum \alpha_i},$$

and

$$d_{\Lambda}(\underline{X}) = \frac{\underline{\alpha} + \underline{X}}{\sum \alpha_i + n} = \frac{\sum \alpha_i}{\sum \alpha_i + n} \frac{\underline{\alpha}}{\sum \alpha_i} + \frac{n}{\sum \alpha_i + n} \frac{\underline{X}}{n}.$$

## 5. FINDING BAYES RULES

For convex loss function we can restrict attention to non-randomized rules; see corollary 1.6.2, TPE, page 51. The following theorem gives a useful recipe for finding Bayes rules.

**Theorem 5.1.** Suppose that  $\boldsymbol{\theta} \sim \Lambda$ ,  $(X | \boldsymbol{\theta} = \theta) \sim P_\theta$ , and  $L(\theta, a) \geq 0$  for all  $\theta \in \Theta$ ,  $a \in \mathbf{A}$ . If

- (i) There exists a rule  $d_0$  with finite risk, and
- (ii) For a.e.  $x$  there exists  $d_\Lambda(x)$  minimizing

$$E\{L(\boldsymbol{\theta}, d(\cdot | x)) | X = x\} = \int_{\Theta} L(\theta, d(\cdot | x)) \lambda(\theta | x) d\nu(\theta),$$

then  $d_\Lambda(\cdot | x)$  is a Bayes rule.

**Proof.** For any rule  $d$  with finite risk

$$(a) \quad E\{L(\boldsymbol{\theta}, d(X)) | X\} \geq E\{L(\boldsymbol{\theta}, d_\Lambda(X)) | X\} \quad \text{a.s.}$$

by (ii). Hence

$$(b) \quad \begin{aligned} \mathcal{R}(\Lambda, d) &= EE\{L(\boldsymbol{\theta}, d(X)) | X\} = E L(\boldsymbol{\theta}, d(X)) \\ &\geq EE\{L(\boldsymbol{\theta}, d_\Lambda(X)) | X\} = E L(\boldsymbol{\theta}, d_\Lambda(X)) \\ &\geq \mathcal{R}(\Lambda, d_\Lambda). \end{aligned}$$

**Corollary 5.2.** (Estimation with weighted squared - error loss). If  $\Theta = \mathbf{A}$  and  $L(\theta, a) = K(\theta)|\theta - a|^2$ , then

$$d_\Lambda(x) = \frac{E\{K(\boldsymbol{\theta})\boldsymbol{\theta} | X = x\}}{E\{K(\boldsymbol{\theta}) | X = x\}} = \frac{\int \theta K(\theta) d\Lambda(\theta | x)}{\int K(\theta) d\Lambda(\theta | x)}.$$

When  $K(\theta) \equiv 1$ , then

$$d_\Lambda(x) = \int \theta d\Lambda(\theta | x) = E\{\boldsymbol{\theta} | X = x\} \equiv \text{the posterior mean}$$

**Proof.** For an arbitrary (nonrandomized) rule  $d \in \mathbf{D}$

$$\begin{aligned} \int K(\theta) |\theta - d(x)|^2 d\Lambda(\theta | x) &= \int K(\theta) |\theta - d_\Lambda(x) + d_\Lambda(x) - d(x)|^2 d\Lambda(\theta | x) \\ &= \int K(\theta) |\theta - d_\Lambda(x)|^2 d\Lambda(\theta | x) \\ &\quad + 2(d_\Lambda(x) - d(x)) \int K(\theta) \{\theta - d_\Lambda(x)\} d\Lambda(\theta | x) \end{aligned}$$

$$\begin{aligned}
 &+ (d_{\Lambda}(x) - d(x))^2 \int K(\theta) d\Lambda(\theta | x) \\
 &\geq \int K(\theta) |\theta - d_{\Lambda}(x)|^2 d\Lambda(\theta | x)
 \end{aligned}$$

with equality if  $d(x) = d_{\Lambda}(x)$ . □

**Corollary 5.3.** (Estimation with  $L_1$ -loss). If  $\Theta = \mathbf{A}$  and  $L(\theta, a) = |\theta - a|$ , then

$$d_{\Lambda}(x) = \text{any median of } \Lambda(\theta | x).$$

**Corollary 5.4.** (Testing with 0-1 loss). If  $\mathbf{A} = \{0, 1\}$ ,  $\Theta = \Theta_0 + \Theta_1$ ,  $L(\theta, a_i) = l_i 1_{\Theta_i^c}(\theta)$ ,  $i = 0, 1$ , then any rule of the form

$$d_{\Lambda}(x) = \begin{cases} 1 & P(\boldsymbol{\theta} \in \Theta_1 | X = x) > (l_1/l_0)P(\boldsymbol{\theta} \in \Theta_0 | X = x) \\ \gamma(x) & \text{if } P(\boldsymbol{\theta} \in \Theta_1 | X = x) = (l_1/l_0)P(\boldsymbol{\theta} \in \Theta_0 | X = x) \\ 0 & P(\boldsymbol{\theta} \in \Theta_1 | X = x) < (l_1/l_0)P(\boldsymbol{\theta} \in \Theta_0 | X = x) \end{cases}$$

is Bayes with respect to  $\Lambda$ . Note that this reduces to a test of the Neyman-Pearson form when  $\Theta_i = \{\theta_i\}$ ,  $i = 0, 1$ .

**Proof.** Let  $\phi(x) = d(1 | x)$ . Then

$$\begin{aligned}
 E\{L(\boldsymbol{\theta}, \phi(x)) | X = x\} &= \int_{\Theta} L(\boldsymbol{\theta}, \phi(x)) d\Lambda(\boldsymbol{\theta} | x) \\
 &= \int_{\Theta} \{l_1 \phi(x) 1_{\Theta_0}(\boldsymbol{\theta}) + l_0(1 - \phi(x)) 1_{\Theta_1}(\boldsymbol{\theta})\} d\Lambda(\boldsymbol{\theta} | x) \\
 &= l_1 \phi(x) P(\boldsymbol{\theta} \in \Theta_0 | X = x) + l_0(1 - \phi(x)) P(\boldsymbol{\theta} \in \Theta_1 | X = x) \\
 &= l_0 P(\boldsymbol{\theta} \in \Theta_1 | X = x) + \phi(x) \{l_1 P(\boldsymbol{\theta} \in \Theta_0 | X = x) - l_0 P(\boldsymbol{\theta} \in \Theta_1 | X = x)\}
 \end{aligned}$$

which is minimized by any rule of the form  $d_{\Lambda}$ . □

**Corollary 5.5.** (Testing with linear loss). If  $\Theta = R$ ,  $\mathbf{A} = \{0, 1\}$ ,  $\Theta_0 = (-\infty, \theta_0]$ ,  $\Theta_1 = (\theta_0, \infty)$  and

$$L(\theta, 0) = (\theta - \theta_0) 1_{(\theta_0, \infty)}(\theta); \quad L(\theta, 1) = (\theta_0 - \theta) 1_{(-\infty, \theta_0]}(\theta);$$

then

$$d_{\Lambda}(x) = \begin{cases} 1 & E(\boldsymbol{\theta} | X = x) > \theta_0 \\ \gamma(x) & \text{if } E(\boldsymbol{\theta} | X = x) = \theta_0 \\ 0 & E(\boldsymbol{\theta} | X = x) < \theta_0 \end{cases}$$

is Bayes with respect to  $\Lambda$ .

**Proof.** Again, by theorem 5.1 it suffices to minimize

$$E\{L(\boldsymbol{\theta}, \phi(x)) | X = x\}$$

$$\begin{aligned}
 &= \int_{\Theta} \{ \phi(x)(\theta_0 - \theta)1_{(-\infty, \theta_0]}(\theta) + (1 - \phi(x))(\theta - \theta_0)1_{(\theta_0, \infty)}(\theta) \} d\Lambda(\theta | x) \\
 &= \int_{\Theta} (\theta - \theta_0)1_{(-\infty, \theta_0]}(\theta) d\Lambda(\theta | x) + (1 - \phi(x)) \{ E(\boldsymbol{\theta} | X = x) - \theta_0 \}
 \end{aligned}$$

which is minimized for each fixed  $x$  by any rule of the form  $d_{\Lambda}$ . □

**More examples 5.6: Estimation with squared error and weighted squared error loss.**

**Example 5.6.1.** Binomial( $n, \theta$ ) with  $\theta \sim \text{Beta}(\alpha, \beta)$ . Thus  $(\theta | X = x) \sim B(\alpha + x, \beta + n - x)$ , so that

$$E(\boldsymbol{\theta}) = \frac{\alpha}{\alpha + \beta}$$

and

$$\begin{aligned}
 d_{\Lambda}(X) &= \frac{\alpha + X}{\alpha + \beta + n} \\
 &= \frac{\alpha + \beta}{\alpha + \beta + n} \frac{\alpha}{\alpha + \beta} + \frac{n}{\alpha + \beta + n} \frac{X}{n}.
 \end{aligned}$$

If the loss is  $L(\theta, a) = (\theta - a)^2 / \{\theta(1 - \theta)\}$  instead of squared - error loss, then

$$\begin{aligned}
 E\{\boldsymbol{\theta}K(\boldsymbol{\theta}) | X = x\} &= \frac{B(\alpha + x, \beta + n - x - 1)}{B(\alpha + x, \beta + n - x)}, \\
 E\{K(\boldsymbol{\theta}) | X = x\} &= \frac{B(\alpha + x - 1, \beta + n - x - 1)}{B(\alpha + x, \beta + n - x)},
 \end{aligned}$$

and hence the Bayes rule with respect to  $\Lambda$  for this loss function is

$$\begin{aligned}
 d_{\Lambda}(X) &= \frac{B(\alpha + X, \beta + n - X - 1)}{B(\alpha + X - 1, \beta + n - X - 1)} = \frac{\alpha + X - 1}{\alpha + \beta + n - 2} \\
 &= \lambda_n(\alpha, \beta) \frac{\alpha - 1}{\alpha + \beta - 2} + (1 - \lambda_n)(\alpha, \beta) \frac{X}{n}
 \end{aligned}$$

where  $\lambda_n(\alpha, \beta) = (\alpha + \beta - 2) / (\alpha + \beta + n - 2)$ . Note that when  $\alpha = 1 = \beta$ , the Bayes estimator for this loss function becomes the familiar maximum likelihood estimator  $\hat{p} = X/n$ .

**Example 5.6.2.** Multinomial/Dirichlet. If  $(\underline{X} | \underline{\theta}) \sim \text{Multinomial}_k(n, \underline{\theta})$ ,  $\underline{\theta} \sim \text{Dirichlet}(\underline{\alpha})$ , then

$$E(\underline{\boldsymbol{\theta}}) = \frac{\underline{\alpha}}{\sum \alpha_i},$$

and

$$d_{\Lambda}(\underline{X}) = E(\underline{\boldsymbol{\theta}} | \underline{X}) = \frac{\underline{\alpha} + \underline{X}}{\sum \alpha_i + n}.$$

**Example 5.6.3.** Normal with normal prior. If  $(X | \theta = \theta) \sim N(\theta, \sigma^2)$ ,  $\theta \sim N(\mu, \tau^2)$ , then

$$(\theta | X) \sim N\left(\frac{1/\tau^2}{1/\tau^2 + 1/\sigma^2} \mu + \frac{1/\sigma^2}{1/\tau^2 + 1/\sigma^2} X, \frac{1}{1/\tau^2 + 1/\sigma^2}\right).$$

Consequently

$$E(\theta) = \mu,$$

and

$$d_{\Lambda}(X) = E(\theta | X) = \frac{1/\tau^2}{1/\tau^2 + 1/\sigma^2} \mu + \frac{1/\sigma^2}{1/\tau^2 + 1/\sigma^2} X.$$

This remains true if  $L(\theta, a) = \rho(\theta - a)$  where  $\rho$  is convex and even (see e.g. Lehmann, TPE, page 244).

## 6. FINDING MINIMAX RULES

**Definition 6.1.** A prior  $\Lambda_0$  for which  $\mathbb{R}(\Lambda, d_\Lambda)$  is maximized is called a *least favorable prior*:

$$\mathbb{R}(\Lambda_0, d_{\Lambda_0}) = \sup_{\Lambda} \mathbb{R}(\Lambda, d_\Lambda).$$

**Theorem 6.2.** Suppose that  $\Lambda$  is a prior distribution on  $\Theta$  such that

$$(1) \quad \mathbb{R}(\Lambda, d_\Lambda) = \int_{\Theta} R(\theta, d_\Lambda) d\Lambda(\theta) = \sup_{\theta} R(\theta, d_\Lambda).$$

Then:

- (i)  $d_\Lambda$  is minimax.
- (ii) If  $d_\Lambda$  is unique Bayes with respect to  $\Lambda$ ,  $d_\Lambda$  is unique minimax.
- (iii)  $\Lambda$  is least favorable.

**Proof.** (i) Let  $d$  be another rule. Then

$$\begin{aligned} (a) \quad \sup_{\theta} R(\theta, d) &\geq \int_{\Theta} R(\theta, d) d\Lambda(\theta) \\ (b) \quad &\geq \int_{\Theta} R(\theta, d_\Lambda) d\Lambda(\theta) \\ &\quad \text{since } d_\Lambda \text{ is Bayes wrt } \Lambda \\ &= \sup_{\theta} R(\theta, d_\Lambda) \quad \text{by (1)}. \end{aligned}$$

Hence  $d_\Lambda$  is minimax.

(ii). If  $d_\Lambda$  is unique Bayes, then  $>$  holds in (b), so  $d_\Lambda$  is unique minimax.

(iii). Let  $\Lambda^*$  be some other prior distribution. Then

$$\begin{aligned} r_{\Lambda^*} &\equiv \int_{\Theta} R(\theta, d_{\Lambda^*}) d\Lambda^*(\theta) \leq \int_{\Theta} R(\theta, d_\Lambda) d\Lambda^*(\theta) \\ &\quad \text{since } d_{\Lambda^*} \text{ is Bayes wrt } \Lambda^* \\ &\leq \sup_{\theta} R(\theta, d_\Lambda) \\ &= \mathbb{R}(\Lambda, d_\Lambda) \equiv r_\Lambda \quad \text{by (1)}. \end{aligned}$$

□

**Corollary 6.3.** If  $d_\Lambda$  is Bayes with respect to  $\Lambda$  and has constant risk,  $R(\theta, d_\Lambda) = \text{constant}$ , then  $d_\Lambda$  is minimax.

**Proof.** If  $d_\Lambda$  has constant risk, then (1) holds. □

**Corollary 6.4.** Let

$$\Theta_\Lambda \equiv \{\theta \in \Theta : R(\theta, d_\Lambda) = \sup_{\theta'} R(\theta', d_\Lambda)\}$$

be the set of  $\theta$ 's where the risk of  $d_\Lambda$  assumes its maximum. Then  $d_\Lambda$  is minimax if  $\Lambda(\Theta_\Lambda) = 1$ . Equivalently,  $d_\Lambda$  is minimax if there is a set  $\Theta_\Lambda$  with  $\Lambda(\Theta_\Lambda) = 1$  and  $R(\theta, d_\Lambda) = \sup_{\theta'} R(\theta', d_\Lambda)$  for all  $\theta \in \Theta_\Lambda$ .

**Example 6.5.** Binomial( $n, \theta$ ) with  $\theta \sim \text{Beta}(\alpha, \beta)$ . Thus  $(\theta | X = x) \sim B(\alpha + x, \beta + n - x)$ , and as we computed in section 5,

$$d_\Lambda(X) = \frac{\alpha + X}{\alpha + \beta + n} = \frac{\alpha + \beta}{\alpha + \beta + n} \frac{\alpha}{\alpha + \beta} + \frac{n}{\alpha + \beta + n} \frac{X}{n}.$$

Consequently,

$$\begin{aligned} R(\theta, d_\Lambda) &= \left(\frac{n}{\alpha + \beta + n}\right)^2 \frac{\theta(1 - \theta)}{n} + \left(\frac{\alpha + n\theta}{\alpha + \beta + n} - \theta\right)^2 \\ &= \frac{1}{(\alpha + \beta + n)^2} \{n\theta(1 - \theta) + (\alpha - \alpha\theta - \beta\theta)^2\} \\ &= \frac{1}{(\alpha + \beta + n)^2} \{\alpha^2 + (n - 2\alpha(\alpha + \beta))\theta + ((\alpha + \beta)^2 - n)\theta^2\} \\ &= \frac{1}{(\alpha + \beta + n)^2} \{\alpha^2 + (n - 2\alpha(\alpha + \beta))\theta + ((\alpha + \beta)^2 - n)\theta^2\} \\ &= \frac{1}{(\alpha + \beta + n)^2} \alpha^2 \end{aligned}$$

if  $2\alpha(\alpha + \beta) = n$  and  $(\alpha + \beta)^2 = n$ . But solving these two equations yields  $\alpha = \beta = \sqrt{n}/2$ . Thus for these choices of  $\alpha$  and  $\beta$ , the risk of the Bayes rule is constant in  $\theta$ , and hence

$$d_M(X) = \frac{1}{1 + \sqrt{n}} \frac{1}{2} + \frac{\sqrt{n}}{1 + \sqrt{n}} \frac{X}{n}$$

is Bayes with respect to  $\Lambda = \text{Beta}(\sqrt{n}/2, \sqrt{n}/2)$  and has risk

$$R(\theta, d_M) = \frac{1}{4(1 + \sqrt{n})^2}.$$

Hence  $\text{Beta}(\sqrt{n}/2, \sqrt{n}/2)$  is least favorable and  $d_M$  is minimax. Note that  $d_M$  is a consistent estimator of  $\theta$ :

$$d_M(X) \rightarrow_p \theta$$

as  $n \rightarrow \infty$ , but its bias is of the order  $n^{-1/2}$  (the bias equals  $(1/2 - \theta)/(1 + \sqrt{n})$ ). Furthermore  $d_M$  is a (locally) regular estimator of  $\theta$ , but the limit distribution is *not* centered at zero if  $\theta \neq 1/2$ : under  $P_{\theta_n}$

$$\sqrt{n}(d_M(X) - \theta_n) \rightarrow_d N(1/2 - \theta, \theta(1 - \theta)).$$



If a least - favorable prior does not exist, then we can still consider improper priors or limits of proper priors:

Let  $\{\Lambda_k\}$  be a sequence of prior distributions, let  $d_k$  denote the Bayes estimator corresponding to  $\Lambda_k$ , and set

$$r_k \equiv \int_{\Theta} R(\theta, d_k) d\Lambda_k(\theta) .$$

Suppose that

$$(2) \quad r_k \rightarrow r < \infty \quad \text{as } k \rightarrow \infty .$$

**Definition 6.6.** The sequence of prior distributions  $\{\Lambda_k\}$  with Bayes risks  $\{r_k\}$  is said to be *least favorable* if  $r_{\Lambda} \equiv \mathbb{R}(\Lambda, d_{\Lambda}) \leq r$  for any prior distribution  $\Lambda$ .

**Theorem 6.7.** Suppose that  $\{\Lambda_k\}$  is a sequence of prior distributions with Bayes risks satisfying

$$(i) \quad r_k \rightarrow r ,$$

and  $d$  is an estimator for which

$$(ii) \quad \sup_{\theta} R(\theta, d) = r .$$

Then:

- A.  $d$  is minimax, and
- B.  $\{\Lambda_k\}$  is least - favorable.

**Proof.** A. Suppose  $d^*$  is any other estimator. Then

$$\sup_{\theta} R(\theta, d^*) \geq \int R(\theta, d^*) d\Lambda_k(\theta) \geq r_k$$

for all  $k \geq 1$ . Hence

$$\sup_{\theta} R(\theta, d^*) \geq r = \sup_{\theta} R(\theta, d) \quad \text{by (ii) ,}$$

so  $d$  is minimax.

B. If  $\Lambda$  is any prior distribution, then

$$\begin{aligned} r_{\Lambda} &= \int R(\theta, d_{\Lambda}) d\Lambda(\theta) \leq \int R(\theta, d) d\Lambda(\theta) \\ &\leq \sup_{\theta} R(\theta, d) = r \end{aligned}$$

by (ii). □

**Example 6.8.** Let  $X_1, \dots, X_n$  be iid  $N(\theta, \sigma^2)$  given  $\theta = \theta$ , and suppose that  $\theta \sim N(\mu, \tau^2)$  and that  $\sigma^2$  is known. Then it follows from 5.6 that

$$d_{\Lambda}(\underline{X}) = E(\theta | \underline{X}) = \frac{1/\tau^2}{1/\tau^2 + n/\sigma^2} \mu + \frac{n/\sigma^2}{1/\tau^2 + n/\sigma^2} \bar{X}$$

and

$$\begin{aligned}
 \mathcal{R}(\Lambda, d_\Lambda) &\equiv r_\Lambda = E\{\boldsymbol{\theta} - d_\Lambda(\underline{X})\}^2 \\
 &= EE\{[\boldsymbol{\theta} - d_\Lambda(\underline{X})]^2 | \underline{X}\} \\
 &= E \text{Var}\{\boldsymbol{\theta} | \underline{X}\} \\
 &= \frac{1}{1/\tau^2 + n/\sigma^2} \\
 &\rightarrow \frac{\sigma^2}{n} \quad \text{as } \tau^2 \rightarrow \infty \\
 &= R(\boldsymbol{\theta}, \bar{X})
 \end{aligned}$$

for all  $\boldsymbol{\theta} \in \Theta$ . Hence, by theorem 6.7,  $\bar{X}$  is a minimax estimator of  $\boldsymbol{\theta}$ .  $\square$

The remainder of this section is aimed at extending minimaxity of estimators from smaller models to larger ones.

**Lemma 6.9.** Suppose that  $X \sim P \in \mathbf{M} \equiv \{\text{all } P\text{'s on } \mathbf{X}\}$ , and that  $\nu: \mathbf{P} \subset \mathbf{M} \rightarrow R$  is a functional (e.g.  $\nu(P) = E_P X = \int x dP(x)$ ,  $\nu(P) = \text{Var}_P(X)$ ,  $\nu(P) = Pf$  for some fixed function  $f$ ; or, in the case of  $\mathbf{X} = R$ ,  $\nu(P) = F_P^{-1}(t)$  for fixed  $t \in (0, 1)$  where  $F_P(x) = P(X \leq x)$ ). Suppose that  $d$  is a minimax estimator of  $\nu(P)$  for  $P \in \mathbf{P}_0 \subset \mathbf{P}_1$ . If

$$\sup_{P \in \mathbf{P}_0} R(P, d) = \sup_{P \in \mathbf{P}_1} R(P, d),$$

then  $d$  is also minimax for estimating  $\nu(P)$ ,  $P \in \mathbf{P}_1$ .

**Proof.** Suppose that  $d$  is not minimax. Then there exists a rule  $d^*$  with smaller maximum risk:

$$\sup_{P \in \mathbf{P}_0} R(P, d^*) \leq \sup_{P \in \mathbf{P}_1} R(P, d^*) < \sup_{P \in \mathbf{P}_1} R(P, d) \leq \sup_{P \in \mathbf{P}_0} R(P, d)$$

which contradict the hypothesis that  $d$  is minimax for  $\mathbf{P}_0$ .  $\square$

**Example 6.10.** This continues the normal mean example 6.8. Suppose that  $\sigma^2$  is unknown. Thus, with  $\Theta = \{(\boldsymbol{\theta}, \sigma^2) : \boldsymbol{\theta} \in R, 0 < \sigma^2 < \infty\}$ , So, to get something reasonable, we need to restrict  $\sigma^2$ . Let

$$\begin{aligned}
 \mathbf{P}_0 &= \{N(\boldsymbol{\theta}, M) : \boldsymbol{\theta} \in R\}, \\
 \mathbf{P}_1 &= \{N(\boldsymbol{\theta}, \sigma^2) : \boldsymbol{\theta} \in R, 0 \leq \sigma^2 \leq M\}.
 \end{aligned}$$

Then  $\bar{X}$  is minimax for  $\mathbf{P}_0$  by our earlier calculation,  $\mathbf{P}_0 \subset \mathbf{P}_1$ , and

$$\sup_{\mathbf{P}_1} R(P, \bar{X}) = \sup_{\mathbf{P}_0} R(P, \bar{X}) = \frac{M}{n}.$$

Thus  $\bar{X}$  is a minimax estimator of  $\theta$  for  $\mathbf{P}_1$ . □

**Example 6.11.** Let  $X_1, \dots, X_n$  be iid  $P \in \mathbf{P}_\mu$  where

$$\mathbf{P}_\mu = \{\text{all probability measures } P : E_P|X| < \infty\}$$

and consider estimation of  $\nu(P) = E_P(X)$  with squared error loss for the families

- (i)  $\mathbf{P}_{b\sigma^2} \equiv \{P \in \mathbf{P}_\mu : \text{Var}_P(X) \leq M < \infty\}$ .
- (ii)  $\mathbf{P}_{br} \equiv \{P \in \mathbf{P}_\mu : P(a \leq X \leq b) = 1\}$  for some fixed  $a, b \in R$ .

Then:

A.  $\bar{X}$  is minimax for  $\mathbf{F}_{b\sigma^2} = \mathbf{P}_1$  by 6.9 since it is minimax for  $\mathbf{P}_0 \equiv \{N(\theta, \sigma^2) : \theta \in R, 0 < \sigma^2 \leq M\}$  and

$$\sup_{P \in \mathbf{P}_0} R(P, \bar{X}) = \sup_{P \in \mathbf{P}_1} R(P, \bar{X}).$$

B. Without loss of generality suppose  $a = 0$  and  $b = 1$ . Let

$$\begin{aligned} \mathbf{P}_1 &= \{P : P([0, 1]) = 1\}, \\ \mathbf{P}_0 &= \{P \in \mathbf{P}_1 : P(X = 1) = p, P(X = 0) = 1 - p \text{ for some } 0 < p < 1\}. \end{aligned}$$

For  $\mathbf{P}_0$  we know that the minimax estimator is

$$d_M(\underline{X}) = \frac{1}{1 + \sqrt{n}} \frac{1}{2} + \frac{\sqrt{n}}{1 + \sqrt{n}} \bar{X}.$$

Now, with  $E_P X \equiv \theta$ ,

$$\begin{aligned} R(P, d_M) &= \frac{1}{(1 + \sqrt{n})^2} \{\text{Var}_P(X) + (\frac{1}{2} - \theta)^2\} \\ &\leq \frac{1}{(1 + \sqrt{n})^2} \{\theta - \theta^2 + \frac{1}{4} - \theta + \theta^2\} \\ &\quad \text{since } 0 \leq X \leq 1 \text{ implies } E_P X^2 \leq E_P X \equiv \theta \\ &= \frac{1/4}{(1 + \sqrt{n})^2}. \end{aligned}$$

Thus

$$\sup_{P \in \mathbf{P}_1} R(P, d_M) = \sup_{P \in \mathbf{P}_0} R(P, d_M),$$

and by 6.9  $d_M$  is minimax.

- Notes.**
1. The restriction  $\sigma^2 \leq M$  in 6.10 (i) is crucial; note that  $\nu(P) = E_P(X)$  is discontinuous at every  $P$  in the sense that we can easily have  $P_m \rightarrow_d P$ , but  $\nu(P_m)$  fails to converge to  $\nu(P)$ .
  2. If  $\Theta = [-M, M] \subset R$ , then  $\bar{X}$  is no longer minimax in the normal case;

see Bickel (1981), *Ann. Statist.* **9**, 1301 - 1309. (The approximately least favorable densities for  $M \rightarrow \infty$  are  $\lambda_M(\theta) = M^{-1} \cos^2(\frac{\pi}{2} \frac{\theta}{M}) 1_{[-M, M]}(\theta)$ .)

## 7. ADMISSIBILITY AND INADMISSIBILITY

Our goal in this section is to establish some basic and simple results about admissibility / inadmissibility of some classical estimators. In particular, we will show that  $\bar{X}$  is admissible for the Gaussian location model for  $k = 1$ , but inadmissible for the Gaussian location model for  $k \geq 3$ . The fundamental work in this area is that of Stein (1956) and James and Stein (1961). For an interesting expository paper, see Efron and Morris (1977); for current work motivated by admissibility issues see e.g. Perlman and Marden (1990) or Strawderman and Brandwein (1991).

**Theorem 7.1.** Any unique Bayes estimator is admissible.

**Proof.** Suppose that  $d_\Lambda$  is unique Bayes with respect to  $\Lambda$  and is inadmissible. Then there exists an estimator  $d$  such that  $R(\theta, d) \leq R(\theta, d_\Lambda)$  with strict inequality for some  $\theta$ , and hence

$$\int_{\Theta} R(\theta, d) \Lambda(\theta) \leq \int_{\Theta} R(\theta, d_\Lambda) d\Lambda(\theta)$$

which contradicts uniqueness of  $d_\Lambda$ . Thus  $d_\Lambda$  is admissible.  $\square$

**Lemma 7.2.** If the loss function  $L(\theta, a)$  is squared error (or is convex in  $a$ ) and, with  $Q$  defined by  $Q(A) = \int P_\theta(X \in A) d\Lambda(\theta)$ , a.e.  $Q$  implies a.e.  $\mathbf{P}$ , then a Bayes rule with finite Bayes risk is unique a.e.  $\mathbf{P}$ . [a.e.  $\mathbf{P}$  means  $P(N) = 0$  for all  $P \in \mathbf{P}$ .]

**Proof.** See Lehmann TPE Corollary 4.1.2, page 240.  $\square$

**Example 7.3.** Consider the Bayes estimator of a normal mean,  $d_\Lambda = p_n \mu + (1 - p_n) \bar{X}$ ,  $p_n = (1/\tau^2)/(1/\tau^2 + n/\sigma^2)$ . The Bayes risk is finite and a.e.  $Q$  implies a.e.  $\mathbf{P}$ . Hence  $d_\Lambda$  is unique Bayes and admissible.

**Theorem 7.4.** If  $X$  is a random variable with mean  $\theta$  and variance  $\sigma^2$ , then  $aX + b$  is *inadmissible* as an estimator of  $\theta$  for squared error loss if:

- (i)  $a > 1$
- (ii)  $a < 0$
- (iii)  $a = 1, b \neq 0$ .

**Proof.** For any  $a, b$  the risk of the rule  $aX + b$  is

$$R(\theta, aX + b) = a^2\sigma^2 + \{(a-1)\theta + b\}^2 \equiv \rho(a, b).$$

(i) If  $a > 1$

$$\rho(a, b) \geq a^2\sigma^2 > \sigma^2 = \rho(1, 0),$$

so  $aX + b$  is dominated by  $X$ .

(ii) If  $a < 0$ , then  $(a-1)^2 > 1$  and

$$\begin{aligned} \rho(a, b) &\geq \{(a-1)\theta + b\}^2 = (a-1)^2\left\{\theta + \frac{b}{a-1}\right\}^2 \\ &> \left\{\theta + \frac{b}{a-1}\right\}^2 = \rho\left(0, -\frac{b}{a-1}\right). \end{aligned}$$

(iii). If  $a = 1, b \neq 0$ ,

$$\rho(1, b) = \sigma^2 + b^2 > \sigma^2 = \rho(1, 0),$$

so  $X + b$  is dominated by  $X$ . □

**Example 7.5, continued.** Thus  $a\bar{X} + b$  is inadmissible for  $a < 0$  or  $a > 1$ . When  $a = 0, d = b$  is admissible since it is the only estimator with zero risk at  $\theta = b$ . When  $a = 1, b \neq 0, d$  is inadmissible. What is left is  $\bar{X}$ .

**Theorem 7.6.** If  $X_1, \dots, X_n$  are iid  $N(\theta, \sigma^2)$ ,  $\theta \in \Theta = \mathcal{R}$  with  $\sigma^2$  known, then  $\bar{X}$  is an admissible estimator of  $\theta$ .

**Proof A. (Limiting Bayes method).** Suppose that  $\bar{X}$  is inadmissible (and  $\sigma = 1$ ). Then there is an estimator  $d^*$  such that  $R(\theta, d^*) \leq \frac{1}{n} = R(\theta, \bar{X})$  for all  $\theta$  with risk  $<$  for some  $\theta$ . Now

$$R(\theta, d) = E(\theta - d(X))^2$$

is continuous in  $\theta$  for every  $d$ , and hence there exists  $\epsilon > 0$  and  $\theta_0 < \theta_1$  so that

$$R(\theta, d^*) < \frac{1}{n} - \epsilon \quad \text{for all } \theta_0 < \theta < \theta_1.$$

Let

$$r_\tau^* \equiv \int_{\Theta} R(\theta, d^*) d\Lambda(\theta)$$

where  $\Lambda = N(0, \tau^2)$ . Thus

$$r_\tau \equiv \int_{\Theta} R(\theta, d_\tau) d\Lambda(\theta) = \frac{1}{1/\tau^2 + n} = \frac{\tau^2}{1 + n\tau^2}$$

so  $r_\tau \leq r_\tau^*$ . Thus

$$\begin{aligned}
 \frac{1/n - r_\tau^*}{1/n - r_\tau} &= \frac{\frac{1}{\sqrt{2\pi\tau}} \int_{-\infty}^{\infty} \left\{ \frac{1}{n} - R(\theta, d^*) \right\} \exp(-\theta^2/2\tau^2) d\theta}{1/n - \tau^2/(1 + n\tau^2)} \\
 &\geq \frac{\frac{1}{\sqrt{2\pi\tau}} \epsilon \int_{\theta_0}^{\theta_1} \exp(-\theta^2/2\tau^2) d\theta}{1/n(1 + n\tau^2)} \\
 &= \frac{n(1 + n\tau^2)}{\sqrt{2\pi\tau}} \epsilon \int_{\theta_0}^{\theta_1} \exp(-\theta^2/2\tau^2) d\theta \\
 &\rightarrow \infty (\theta_1 - \theta_0) = \infty \quad \text{as } \tau \rightarrow \infty.
 \end{aligned}$$

Hence

$$\frac{1}{n} - r_\tau^* > \frac{1}{n} - r_\tau$$

for  $\tau > \text{some } \tau_0$ , or,

$$r_\tau^* < r_\tau$$

for  $\tau > \text{some } \tau_0$ , which contradicts  $d_\tau$  Bayes wrt  $\Lambda_\tau$  with Bayes risk  $r_\tau$ . Hence  $\bar{X}$  is admissible.  $\square$

**Proof B. (Information inequality method).** The risk of any rule  $d$  is

$$\begin{aligned}
 R(\theta, d) &= E_\theta(d - \theta)^2 = \text{Var}_\theta[d(X)] + b^2(\theta) \\
 &\geq \frac{[1 + b'(\theta)]^2}{n} + b^2(\theta)
 \end{aligned}$$

since  $I(\theta) = 1$ . Suppose that  $d$  is any estimator satisfying

$$R(\theta, d) \leq 1/n \quad \text{for all } \theta.$$

Then

$$(a) \quad \frac{1}{n} \geq b^2(\theta) + \frac{[1 + b'(\theta)]^2}{n}.$$

But (a) implies that:

- (i)  $|b(\theta)| \leq 1/\sqrt{n}$  for all  $\theta$ ; i.e.  $b$  is bounded.
- (ii)  $b'(\theta) \leq 0$  since  $1 + 2b'(\theta) + [b'(\theta)]^2 \leq 1$ .
- (iii) There exists  $\theta_i \rightarrow \infty$  such that  $b'(\theta_i) \rightarrow 0$ . [If  $b'(\theta) \leq -\epsilon$  for all  $\theta > \theta_0$ , then  $b(\theta)$  is not bounded.]
- (iv)

$$\begin{aligned}
 b^2(\theta) &\leq \frac{1}{n} - \frac{[1 + b'(\theta)]^2}{n} \\
 &= - \frac{2b'(\theta) + [b'(\theta)]^2}{n}
 \end{aligned}$$

Hence  $b(\theta) = 0$ , and  $R(\theta, d) \equiv 1/n$ . □

Now let  $X_1, \dots, X_n$  be iid  $N_k(\theta, I)$ , and consider estimation of  $\theta \in R^k$  with loss function  $L(\theta, a) = |\theta - a|^2$ . Let  $\bar{X} = \sum_{i=1}^n X_i/n$ . Note that

$$R(\theta, \bar{X}) = E_\theta |\bar{X} - \theta|^2 = \frac{k}{n} \quad \text{for all } \theta.$$

**Theorem 7.7.** If  $k \geq 3$ , then  $\bar{X}$  is inadmissible.

**Proof.** Let  $g : R^k \rightarrow R^k$  have  $E|\frac{\partial}{\partial x_i} g_i(X)| < \infty$ , and consider estimators of the form

$$\hat{\theta}_n = \bar{X} + n^{-1}g(\bar{X}).$$

Now

$$\begin{aligned} E_\theta |\bar{X} - \theta|^2 - E_\theta |\hat{\theta}_n - \theta|^2 &= E_\theta |\bar{X} - \theta|^2 - E_\theta |\bar{X} - \theta + n^{-1}g(\bar{X})|^2 \\ &= -2n^{-1}E_\theta \langle \bar{X} - \theta, g(\bar{X}) \rangle - n^{-2}E_\theta |g(\bar{X})|^2. \end{aligned} \tag{a}$$

To proceed further we need an identity due to Stein.

**Lemma 7.8.** If  $X \sim N(\theta, \sigma^2)$ , and  $g$  is a function with  $E|g'(X)| < \infty$ , then

$$(i) \quad \sigma^2 E g'(X) = E(X - \theta)g(X).$$

If  $X \sim N_k(\theta, \sigma^2 I)$ , and  $g : R^k \rightarrow R^k$  has

$$E|\frac{\partial}{\partial x_i} g_i(X)| < \infty,$$

then:

$$(ii) \quad \sigma^2 E \frac{\partial}{\partial x_i} g_i(X) = E(X_i - \theta_i)g_i(X), \quad i = 1, \dots, k.$$

**Proof.** Without loss of generality, suppose that  $\theta = 0$  and  $\sigma^2 = 1$ . Integration by parts gives

$$\begin{aligned} E g'(X) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} g'(x) \exp(-\frac{x^2}{2}) dx \\ &= \frac{-1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} g(x) \exp(-\frac{x^2}{2}) \{-\frac{2x}{2}\} dx \\ &= E X g(X). \end{aligned}$$

[In applying integration by parts here, we have ignored the term

$uv|_{-\infty}^{\infty} = g(x)\phi(x)|_{-\infty}^{\infty}$ . To prove that this does in fact vanish, it is easiest to apply Fubini's theorem (twice!) in a slightly devious way as follows. Let  $\phi(t) = (2\pi)^{-1/2} \exp(-t^2/2)$  be the standard normal density. Since  $\phi'(t) = -t\phi(t)$ , we have both

$$\phi(x) = - \int_x^{\infty} \phi'(t) dt = \int_x^{\infty} t\phi(t) dt$$

and

$$\phi(x) = \int_{-\infty}^x \phi'(t) dt = - \int_{-\infty}^x t\phi(t) dt .$$

Therefore we can write

$$\begin{aligned} E g'(X) &= \int_{-\infty}^{\infty} g'(x)\phi(x) dx \\ &= \int_0^{\infty} g'(x) \int_x^{\infty} t\phi(t) dt dx - \int_{-\infty}^0 g'(x) \int_{-\infty}^x t\phi(t) dt dx \\ &= \int_0^{\infty} t\phi(t) \left\{ \int_0^t g'(x) dx \right\} dt - \int_{-\infty}^0 t\phi(t) \left\{ \int_t^0 g'(x) dx \right\} dt \\ &= \left( \int_0^{\infty} + \int_{-\infty}^0 \right) \{ t\phi(t)[g(t) - g(0)] \} dt \\ &= \int_{-\infty}^{\infty} tg(t)\phi(t) dt = EXg(X) . \end{aligned}$$

Here the third equality is justified by the hypothesis  $E|g'(X)| < \infty$  and Fubini's theorem.]

To prove (ii), write  $X^{(i)} \equiv (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_k)$ . Then

$$\begin{aligned} \sigma^2 E \frac{\partial}{\partial x_i} g_i(X) &= \sigma^2 EE \left\{ \frac{\partial}{\partial x_i} g_i(X) \mid X^{(i)} \right\} \\ &= EE \{ (X_i - \theta_i) g_i(X) \mid X^{(i)} \} \quad \text{by (i)} \\ &= E(X_i - \theta_i) g_i(X) . \end{aligned}$$

□

Now we return to the proof of the theorem. Using (ii) in (a) shows that the right side of (a) equals

$$(b) \quad -2n^{-2} E_{\theta} \sum_{i=1}^k \frac{\partial g_i}{\partial x_i}(\bar{X}) - n^{-2} E_{\theta} |g(X)|^2 .$$

Now let  $\psi : R^k \rightarrow R$  be twice differentiable, and set

$$g(x) = \nabla \{ \log \psi(x) \} = \frac{1}{\psi(x)} \left( \frac{\partial \psi(x)}{\partial x_1}, \dots, \frac{\partial \psi(x)}{\partial x_k} \right) .$$



Thus

$$\begin{aligned} \frac{\partial g_i}{\partial x_i}(x) &= \frac{1}{\psi(x)} \frac{\partial^2}{\partial x_i^2} \psi(x) - \frac{1}{\psi^2(x)} \left( \frac{\partial \psi(x)}{\partial x_i} \right)^2 \\ &= \frac{1}{\psi(x)} \frac{\partial^2}{\partial x_i^2} \psi(x) - g_i(x)^2 \end{aligned}$$

and

$$\sum_{i=1}^k \frac{\partial g_i}{\partial x_i}(x) = \frac{1}{\psi(x)} \nabla^2 \psi(x) - |g(x)|^2 .$$

Hence the right side of (b) is

$$\begin{aligned} \text{(c)} \quad &= n^{-2} E_\theta |g(\bar{X})|^2 - 2 n^{-2} E_\theta \left\{ \frac{1}{\psi(\bar{X})} \nabla^2 \psi(\bar{X}) \right\} \\ &> 0 \end{aligned}$$

if  $\psi(x) \geq 0$  and  $\nabla^2 \psi \leq 0$ ,  $g \neq 0$  (i.e.  $\psi$  superharmonic). Here is one example of such a function:

A. Suppose that  $\psi(x) = |x|^{-(k-2)} = \{x_1^2 + \dots + x_k^2\}^{-(k-2)/2}$ . then

$$g(x) = \nabla \log \psi(x) = - \frac{k-2}{|x|^2} x$$

and

$$\nabla^2 \psi(x) = 0 ,$$

so  $\psi$  is *harmonic*. Thus

$$\hat{\theta}_n = \left( 1 - \frac{k-2}{n|\bar{X}|^2} \right) \bar{X}$$

and

$$\begin{aligned} E|\bar{X} - \theta|^2 - E|\hat{\theta} - \theta|^2 &= n^{-2} E_\theta |g(\bar{X})|^2 \\ &= \left( \frac{k-2}{n} \right)^2 E_\theta |\bar{X}|^{-2} \\ &= \left( \frac{k-2}{\sqrt{n}} \right)^2 E_\theta |\sqrt{n}(\bar{X} - \theta) + \sqrt{n}\theta|^{-2} \\ &= \left( \frac{k-2}{\sqrt{n}} \right)^2 E_0 |X + \sqrt{n}\theta|^{-2} \\ &= \begin{cases} \left( \frac{k-2}{n} \right)^2 E_0 |n^{-1/2} X + \theta|^{-2} = O(n^{-2}), & \theta \neq 0 \\ \left( \frac{k-2}{n} \right)^2 \frac{1}{k-2} = \frac{k-2}{n}, & \theta = 0 \end{cases} \end{aligned}$$

since  $|X|^2 \sim \chi_k^2$  with  $E(1/\chi_k^2) = 1/(k-2)$ . Hence

$$\frac{E_0|\hat{\theta} - \theta|^2}{E_0|\bar{X} - \theta|^2} = \frac{2}{k} < 1.$$

For general  $\theta$ ,

$$\begin{aligned} R(\theta, \hat{\theta}) &= \frac{k}{n} - \frac{(k-2)^2}{n} E_0\left(\frac{1}{|X + \sqrt{n}\theta|^2}\right) \\ &= \frac{k}{n} \left(1 - \frac{k-2}{k} E_0\left(\frac{k-2}{\chi_k^2(\delta^2)}\right)\right) \end{aligned}$$

since  $|X + \sqrt{n}\theta|^2 \sim \chi_k^2(\delta^2)$  with  $\delta^2 = n|\theta|^2/2$ . Thus

$$\begin{aligned} \frac{R(\theta, \hat{\theta})}{R(\theta, \bar{X})} &= \left(1 - \frac{k-2}{k} E_0\left(\frac{k-2}{\chi_k^2(\delta^2)}\right)\right) \\ &= 2/k \quad \text{when } \theta = 0, \\ &\rightarrow 1 \quad \text{as } n \rightarrow \infty \text{ for fixed } \theta \neq 0, \\ &\rightarrow 1 \quad \text{as } |\theta| \rightarrow \infty \text{ for fixed } n. \end{aligned}$$

□

**Remark 7.9.** Note that the James - Stein estimator

$$\hat{\theta}_n = \left(1 - \frac{k-2}{n|\bar{X}|^2}\right)\bar{X}$$

is not regular at  $\theta = 0$ : if  $\theta_n = tn^{-1/2}$ , then  $\sqrt{n}(\bar{X} - \theta_n) \stackrel{d}{=} Z \sim N_k(0, \sigma^2 I)$  under  $P_{\theta_n}$ , so that

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta) &= \sqrt{n}(\bar{X} - \theta) - \frac{k-2}{|\sqrt{n}(\bar{X} - \theta_n) + t|^2} \{\sqrt{n}(\bar{X} - \theta_n) + t\} \\ &= Z - \frac{k-2}{|Z + t|^2}(Z + t) \end{aligned}$$

which has a distribution dependent on  $t$ .

**Remark 7.10.** It turns out that the James - Stein estimator  $\hat{\theta}_n$  is itself inadmissible; see Lehmann TPE pages 295, 299, 302-307.

**Remark 7.11.** Another interesting function  $\psi$  is

$$\psi(x) = \begin{cases} |x|^{-(k-2)} & |x| \geq \sqrt{k-2} \\ (k-2)^{-(k-2)/2} \exp\left(\frac{1}{2}[(k-2) - |x|^2]\right) & |x| < \sqrt{k-2} \end{cases}.$$

For this  $\psi$ ,

$$g(x) = \nabla \log \psi(x) = \begin{cases} -(k-2)/|x|^2 & \text{if } |x| \leq \sqrt{k-2} \\ -x & |x| \geq \sqrt{k-2} \end{cases}$$

### 8. Asymptotic Theory of Bayes Estimators

We begin with a basic result for Bayesian estimation of a Bernoulli probability  $\theta$ .

Suppose that  $(Y_1, Y_2, \dots) | \theta = \theta$  are iid Bernoulli( $\theta$ ). Suppose that the prior distribution is an arbitrary distribution on  $(0, 1)$ . The Bayes estimator of  $\theta$  wrt squared error loss is the posterior mean  $d(\underline{Y}) = d(Y_1, \dots, Y_n) = E(\theta | \underline{Y})$ . Now this sequence is a (uniformly integrable) martingale; this type of martingale is sometimes known as *Doob's martingale*. Hence, by the martingale convergence theorem,

$$d(\underline{Y}_n) = E(\theta | \underline{Y}_n) \rightarrow_{a.s.} E(\theta | Y_1, Y_2, \dots).$$

To prove consistency, it remains only to show that  $E(\theta | Y_1, Y_2, \dots) = \theta$ .

To see this, we compute conditionally on  $\theta$ : now

$$E\left(\sum_{i=1}^n Y_i | \theta\right) = n\theta, \quad \text{Var}\left(\sum_{i=1}^n Y_i | \theta\right) = n\theta(1 - \theta).$$

Hence, with  $\hat{\theta}_n \equiv n^{-1} \sum_{i=1}^n Y_i$ ,

$$E(\hat{\theta}_n - \theta)^2 \leq n^{-1},$$

and by Chebychev's inequality,  $\hat{\theta}_n$  converges in probability to  $\theta$ . Therefore  $\hat{\theta}_{n_k} \rightarrow_p \theta$  almost surely for some subsequence  $n_k$ . Hence  $\theta = \tilde{\theta}$  a.s. where  $\tilde{\theta} \equiv \lim_k \hat{\theta}_{n_k}$  is measurable with respect to  $Y_1, Y_2, \dots$ . Thus we have

$$\begin{aligned} E(\theta | Y_1, Y_2, \dots) &= E(\tilde{\theta} | Y_1, Y_2, \dots) \\ &= \tilde{\theta} = \theta \quad \text{a.s. } P_\Lambda \end{aligned}$$

where

$$P_\Lambda(\underline{Y} \in A, \theta \in B) = \int_B P_\theta(\underline{Y} \in A) d\Lambda(\theta).$$

Therefore

$$P_\Lambda(d(\underline{Y}_n) \rightarrow \theta) = \int_{[0,1]} P_\theta(d(\underline{Y}_n) \rightarrow \theta) d\Lambda(\theta) = 1,$$

and this implies that

$$P_\theta(d(\underline{Y}_n) \rightarrow \theta) = 1 \quad \text{a.e. } \Lambda.$$

Thus the Bayes estimator  $d(\underline{Y}_n) \equiv E(\theta | \underline{Y}_n)$  is consistent for  $\Lambda$  a.e.  $\theta$ .

A general result of this type was established by Doob (1948), and has been generalized by Breiman, Le Cam, and Schwartz (1964), Freedman (1963), and Schwartz (1965). The latter authors prove that for smooth, finite - dimensional families  $(\theta, \Lambda)$  is consistent if and only if  $\theta$  is in the support of  $\Lambda$ . However the assumption of finite - dimensionality is important: Freedman (1963) gives an example of an inconsistent Bayes rule when the sample space is  $\mathbf{X} = \mathbb{Z}^+$  and  $\Theta$  is the collection of all probabilities on  $\mathbf{X}$ . The recent work of Diaconis and Freedman (1986) on the consistency of Bayes estimates is an outgrowth of the work of Freedman (1963).

**Asymptotic Efficiency of Bayes Estimators**

Now we consider a more general situation, but still with a real-valued parameter  $\theta$ . Let  $X_1, \dots, X_n$  be iid  $p_\theta(x)$  with respect to  $\mu$  where  $\theta \in \Theta \subset \mathbb{R}$ , and  $\theta_0 \in \Theta$  (an open interval) is true; write  $P$  for  $P_{\theta_0}$ .

**Assumptions:**

B1: (i) - (vi) of theorem 1.1 (TPE, page 406) hold. Thus it follows that

$$(1) \quad l(\theta) - l(\theta_0) = (\theta - \theta_0)\dot{l}_\theta(\underline{X}; \theta_0) - \frac{1}{2}(\theta - \theta_0)^2 \{nI(\theta_0) + R_n(\theta)\}$$

where

$$\frac{1}{n} R_n(\theta) \rightarrow_p 0;$$

recall problems 2.22 and 7.2, TPE. We strengthen this to:

B2: For any  $\epsilon > 0$  there exists a  $\delta > 0$  such that

$$P\left\{ \sup_{|\theta - \theta_0| \leq \delta} \left| \frac{1}{n} R_n(\theta) \right| \geq \epsilon \right\} \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

B3: For any  $\delta > 0$  there exists an  $\epsilon > 0$  such that

$$P\left\{ \sup_{|\theta - \theta_0| \geq \delta} \frac{1}{n} (l(\theta) - l(\theta_0)) \leq -\epsilon \right\} \rightarrow 1.$$

B4: The prior density  $\lambda$  of  $\theta$  is continuous and  $> 0$  for all  $\theta \in \Theta$ .

B5:  $E_\lambda|\theta| < \infty$ .

**Theorem 8.1.** Suppose that  $\lambda^*(t|\underline{X})$  is the posterior density of  $\sqrt{n}(\theta - T_n)$  where

$$T_n \equiv \theta_0 + \frac{1}{I(\theta_0)} \{n^{-1}\dot{l}_\theta(\underline{X}; \theta_0)\}.$$

(i) If B1 - B4 hold, then

$$(2) \quad d_{TV}(\Lambda^*(\cdot|\underline{X}), N(0, \frac{1}{I(\theta_0)})) = \int |\lambda^*(t|\underline{X}) - \sqrt{I(\theta_0)}\phi(t\sqrt{I(\theta_0)})| dt \rightarrow_p 0.$$

(ii) If B1 - B5 hold, then

$$(3) \quad \int (1 + |t|) |\lambda^*(t | \underline{X}) - \sqrt{I(\theta_0)} \phi(t\sqrt{I(\theta_0)})| dt \rightarrow_p 0 .$$

**Theorem 8.2.** If B1 - B5 hold, and  $\tilde{\theta}_n$  is the Bayes estimator wrt  $\lambda$  for squared error loss, then

$$(4) \quad \sqrt{n}(\tilde{\theta}_n - \theta_0) \rightarrow_d N(0, \frac{1}{I(\theta_0)}) ,$$

so that  $\tilde{\theta}_n$  is  $\sqrt{n}$ -consistent and asymptotically efficient.

**Example 8.3.** Suppose that  $X_1, \dots, X_n$  are iid  $N(\theta, \sigma^2)$  and that the prior on  $\theta$  is  $N(\mu, \tau^2)$ . Then the posterior distribution is given by  $(\theta | \bar{X}) \sim N(p_n \bar{X} + (1 - p_n)\mu, \sigma_n^2/n)$  where

$$p_n \equiv \frac{n/\sigma^2}{n/\sigma^2 + 1/\tau^2} , \quad \sigma_n^2 = \frac{n}{n/\sigma^2 + 1/\tau^2}$$

so that  $p_n \rightarrow 1$ ,  $\sqrt{n}(1 - p_n) \rightarrow 0$ , and  $\sigma_n^2 \rightarrow \sigma^2$ . Note that

$$T_n = \theta_0 + \frac{1}{1/\sigma^2} \frac{(\bar{X} - \theta_0)}{\sigma^2} = \bar{X} ,$$

and the Bayes estimator with squared error loss is

$$E(\theta | \bar{X}) = p_n \bar{X} + (1 - p_n)\mu .$$

Thus

$$\begin{aligned} (\sqrt{n}(\theta - \bar{X}) | \bar{X}) &\sim N(\sqrt{n}(1 - p_n)(\mu - \bar{X}), \sigma_n^2) \\ &\rightarrow_d N(0, \sigma^2) = N(0, \frac{1}{I(\theta_0)}) \quad \text{a.s.} \end{aligned}$$

in agreement with theorem 1, while

$$\begin{aligned} \sqrt{n}(E(\theta | \bar{X}) - \theta_0) &= \sqrt{n}(p_n \bar{X} + (1 - p_n)\mu - \theta_0) \\ &= p_n \sqrt{n}(\bar{X} - \theta_0) + \sqrt{n}(1 - p_n)(\mu - \theta_0) \\ &=_d p_n N(0, \sigma^2) + \sqrt{n}(1 - p_n)(\mu - \theta_0) \\ &\rightarrow_d N(0, \sigma^2) \end{aligned}$$

in agreement with theorem 2.  $\square$

**Proof of theorem 8.2.** Suppose that theorem 8.1 is proved, and write

$$(a) \quad \sqrt{n}(\tilde{\theta}_n - \theta_0) = \sqrt{n}(\tilde{\theta}_n - T_n) + \sqrt{n}(T_n - \theta_0) .$$

Since

$$(b) \quad \sqrt{n}(T_n - \theta_0) \rightarrow_d N(0, \frac{1}{I(\theta_0)}),$$

it suffices to show that

$$(c) \quad \sqrt{n}(\tilde{\theta}_n - T_n) \rightarrow_p 0.$$

To prove (c), note that by definition of  $\tilde{\theta}_n$

$$\tilde{\theta}_n = \int \theta \lambda(\theta | \underline{X}) d\theta = \int \left\{ \frac{t}{\sqrt{n}} + T_n \right\} \lambda^*(t | \underline{X}) dt,$$

and hence

$$\sqrt{n}(\tilde{\theta}_n - T_n) = \int t \lambda^*(t | \underline{X}) dt - \int t \sqrt{I(\theta_0)} \phi(t \sqrt{I(\theta_0)}) dt.$$

It follows that

$$\sqrt{n}|\tilde{\theta}_n - T_n| \leq \int |t| |\lambda^*(t | \underline{X}) - \sqrt{I(\theta_0)} \phi(t \sqrt{I(\theta_0)})| dt \rightarrow_p 0 \quad \text{as } n \rightarrow \infty$$

by theorem 1(ii). Thus (c) holds, and the proof is complete; it remains only to prove theorem 8.1.  $\square$

**Proof of theorem 8.1:** (i) By definition of  $T_n$

$$\begin{aligned} \lambda^*(t | \underline{X}) &= \frac{\lambda(T_n + n^{-1/2}t) \exp(l(T_n + n^{-1/2}t))}{\int \lambda(T_n + n^{-1/2}u) \exp(l(T_n + n^{-1/2}u)) du} \\ &\equiv e^{\omega(t)} \lambda(T_n + n^{-1/2}t) / C_n \end{aligned}$$

where

$$(a) \quad \omega(t) \equiv l(T_n + n^{-1/2}t) - l(\theta_0) - \frac{1}{2nI(\theta_0)} \{\dot{\mathbf{l}}_{\theta}(\underline{X}; \theta_0)\}^2$$

and

$$(b) \quad C_n \equiv \int e^{\omega(u)} \lambda(T_n + n^{-1/2}u) du.$$

Now if we can show that

$$J_{1n} \equiv \int |e^{\omega(t)} \lambda(T_n + n^{-1/2}t) - e^{-t^2 I(\theta_0)/2} \lambda(\theta_0)| dt \rightarrow_p 0,$$

then

$$(d) \quad C_n \rightarrow_p \int e^{-t^2 I(\theta_0)/2} \lambda(\theta_0) dt = \lambda(\theta_0) \sqrt{2\pi/I(\theta_0)},$$

and the left side of (2) is  $J_n/C_n$  where

$$(e) \quad J_n \equiv \int |e^{\omega(t)} \lambda(T_n + n^{-1/2}t) - C_n \sqrt{I(\theta_0)} \phi(t\sqrt{I(\theta_0)})| dt .$$

Further  $J_n \leq J_{1n} + J_{2n}$  where  $J_{1n}$  is defined in (c) and

$$(f) \quad J_{2n} \equiv \int |C_n \sqrt{I(\theta_0)} \phi(t\sqrt{I(\theta_0)}) - e^{-t^2 I(\theta_0)/2} \lambda(\theta_0)| dt \\ = | \frac{C_n \sqrt{I(\theta_0)}}{\sqrt{2\pi}} - \lambda(\theta_0) | \int e^{-t^2 I(\theta_0)/2} dt \rightarrow_p 0$$

by (d). Hence it remains only to prove that (c) holds.

(ii) The same proof works with  $J'_{1n}$  replacing  $J_{1n}$  where  $J'_{1n}$  has an additional factor of  $(1 + |t|)$ .

**Lemma 8.4.** The following identity holds:

$$(5) \quad \omega(t) \equiv l(T_n + n^{-1/2}t) - l(\theta_0) - \frac{1}{2nI(\theta_0)} \{\dot{\mathbf{i}}_{\theta}(\underline{X}; \theta_0)\}^2 \\ = -I(\theta_0) \frac{t^2}{2} - \frac{1}{2n} R_n(T_n + \frac{t}{\sqrt{n}}) \{t + \frac{1}{I(\theta_0)} \frac{\dot{\mathbf{i}}_{\theta}(\underline{X}; \theta_0)}{\sqrt{n}}\}^2 .$$

**Proof.** This follows immediately from the definition of  $T_n$ , (1), and algebra.

□

**Proof that  $J_{1n} \rightarrow_p 0$ :** Recall the definition of  $J_{1n}$  given in (c) of the proof of theorem 1. Divide the region of integration into:

- (i)  $|t| < M$ ;
- (ii)  $M \leq t \leq \delta \sqrt{n}$ ; and
- (iii)  $\delta \sqrt{n} < t < \infty$ .

For the region (i), the integral is bounded by  $2M$  times the supremum of the integrand over the set  $|t| \leq M$ , and by application of lemma 1 it is seen that this supremum converges to zero in probability if

$$(a) \quad \sup_{|t| \leq M} | \frac{1}{n} R_n(T_n + n^{-1/2}t) (t + \frac{1}{I(\theta_0)} \frac{\dot{\mathbf{i}}_{\theta}(\underline{X}; \theta_0)}{\sqrt{n}})^2 | \rightarrow_p 0$$

and

$$(b) \quad \sup_{|t| \leq M} |\lambda(T_n + n^{-1/2}t) - \lambda(\theta_0)| \rightarrow_p 0 .$$

Now  $\lambda$  is continuous by B4 and  $T_n \rightarrow_p \theta_0$  by B1, so  $\lambda$  is uniformly continuous in a neighborhood of  $\theta_0$  and

$$(c) \quad \theta_0 - \delta \leq T_n - \frac{M}{\sqrt{n}} \leq T_n - \frac{t}{\sqrt{n}} \leq T_n + \frac{M}{\sqrt{n}} \leq \theta_0 + \delta$$

for  $|t| \leq M$ , so (b) holds. Now  $n^{-1/2} \dot{\mathbf{i}}_{\theta}(\underline{X}; \theta_0)$  is bounded in probability (by B1 and the central limit theorem), so it follows from (c) and B2 that (a) holds.

(ii) For the region  $M \leq |t| \leq \delta\sqrt{n}$  it suffices to show that the integrand is bounded by an integrable function with probability close to one, since then the integral can be made small by choosing  $M$  sufficiently large. Since the second term of the integrand in  $J_{1n}$  is integrable, it suffices to find such an integrable bound for the first term. In particular, it can be shown that for every  $\epsilon > 0$  there exists a  $C = C(\epsilon, \delta)$  and an integer  $N = N(\epsilon, \delta)$  so that, for  $n \geq N$ ,

$$(d) \quad P(e^{\omega(t)} \lambda(T_n + \frac{t}{\sqrt{n}})) \leq C \exp(-t^2 I(\theta_0)/4) \text{ for all } |t| \leq \delta\sqrt{n} \geq 1 - \epsilon;$$

this follows from the assumption B2; see TPE, pages 463 - 464.

(iii) The integral over the region  $|t| \geq \delta\sqrt{n}$  converges to zero by argument similar to that for the region (ii), but now the assumption B3 comes into play; see TPE pages 464-465.

The proof for  $J'_{1n}$  requires only trivial changes, using the assumption B5.  $\square$

**Notes:**

1. This proof is from Lehmann, TPE, pages 454 - 465; this type of results apparently dates back to Laplace (1820), and was later rederived by Bernstein (1917), and von Mises (1931). More general versions have been established by Le Cam (1958), Bickel and Yahav (1969), and Ibragimov and Has'minskii (1972). See the discussion on page 482 of TPE. For a proof with weaker regularity conditions, see Hartigan (1983), chapter 11, section 3.
2. See Ibragimov and Has'minskii (1981) for more on the consistency and efficiency of Bayes estimators; also see Strasser (1981) as referenced by Lehmann.
3. Read the examples and discussion on pages 458 - 461 of TPE.
4. Analogues of these theorems for nonparametric and semiparametric problems are still in a relatively undeveloped state; Diaconis and Freedman (1986) is a good starting point. Cox (1993) treats Bayesian nonparametric regression.

Now we will continue the development for convex likelihoods started in section 4.7. The notation will be as in section 4.7, and we will make use of lemmas 4.7.1 and 4.7.2. Here, as in section 4.7, we will *not* assume that the distribution  $P$  governing the data is a member of the parametric model  $\mathbf{P} = \{P_\theta : \theta \in \Theta \subset R^d\}$ .

Here are the hypotheses which will be used.

Suppose that  $\lambda(\theta)$  is the prior density (with respect to Lebesgue measure  $\nu$  on  $R^d$ ). We will assume that:



- C1.  $\lambda(\theta) \leq C_1 \exp(C_2|\theta|)$  for all  $\theta$  for some positive constants  $C_1, C_2$ .
- C2.  $\lambda(\theta)$  is continuous at  $\theta_0$ .
- C3.  $\log p(x, \theta_0 + t) - \log p(x, \theta_0) = \psi(x)^T t + R(x; t)$  is concave in  $t$ .
- C4.  $E_P \psi(X_1) \psi(X_1)^T \equiv K$ ,  $E_P \psi(X_1) = 0$ .
- C5.  $E_P R(X_1; t) = -t^T J t / 2 + o(|t|^2)$  and  $Var_P(R(X_1; t)) = o(|t|^2)$  where  $J$  is symmetric, positive definite.
- C6.  $X_1, \dots, X_n$  are i.i.d.  $P$ .

**Theorem 8.5.** Suppose that C1 - C6 hold. Then the maximum likelihood estimator  $\hat{\theta}_n$  is  $\sqrt{n}$ -consistent for  $\theta_0 = \theta_0(P)$  and satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow_d N_d(0, J^{-1}K(J^{-1})^T).$$

Moreover the Bayes estimator  $\theta_n^* = E\{\theta | \underline{X}_n\}$  satisfies

$$\sqrt{n}(\theta_n^* - \hat{\theta}_n) \rightarrow_p 0,$$

and hence also

$$\sqrt{n}(\theta_n^* - \theta_0) \rightarrow_d N_d(0, J^{-1}K(J^{-1})^T).$$

**Proof.** Let  $L_n(\theta) \equiv \prod_{i=1}^n p(X_i, \theta)$  denote the likelihood function. Define random convex functions  $A_n$  by

$$\exp(-A_n(t)) = L_n(\hat{\theta}_n + t/\sqrt{n}) / L_n(\hat{\theta}_n).$$

By definition of the MLE,  $A_n$  achieves its minimum value of zero at  $t = 0$ . Then

$$\begin{aligned} \theta_n^* &= \frac{\int \theta L_n(\theta) \lambda(\theta) d\theta}{\int L_n(\theta) \lambda(\theta) d\theta} \\ &= \hat{\theta}_n + \frac{1}{\sqrt{n}} \frac{\int t \exp(-A_n(t)) \lambda(\hat{\theta}_n + t/\sqrt{n}) \exp(-C_2|\hat{\theta}_n|) dt}{\int \exp(-A_n(t)) \lambda(\hat{\theta}_n + t/\sqrt{n}) \exp(-C_2|\hat{\theta}_n|) dt} \end{aligned}$$

by the change of variable  $\theta = \hat{\theta}_n + t/\sqrt{n}$ .

**Claim:** The random functions  $A_n$  converge in probability uniformly on compact sets to  $t^T J t / 2$ .

**Proof.** Let

$$A_n^0(t) = n IP_n \{ \log p(x; \theta_0 + t/\sqrt{n}) - \log p(x; \theta_0) \}$$

From the proof of theorem 4.7.6 (with  $h(x; \theta) \equiv -\log p(x; \theta)$  and multiplication by  $-1$ )

$$(*) \quad A_n^0(t) = U_n^T t + \frac{1}{2} t^T J t + o_p(1),$$

and this holds uniformly in  $t$  in compact sets by lemma 4.7.1. Note that

$$A_n(t) = A_n^0(t + \sqrt{n}(\hat{\theta}_n - \theta_0)) - A_n^0(\sqrt{n}(\hat{\theta}_n - \theta_0)).$$

Since  $\sqrt{n}(\hat{\theta}_n - \theta_0) = O_p(1)$  so that for every  $\epsilon > 0$  there exists a compact set  $K \equiv K_\epsilon$  such that  $P(\sqrt{n}(\hat{\theta}_n - \theta_0) \in K) \geq 1 - \epsilon$ . Hence by the uniformity in  $t$  in the convergence in (\*),

$$\begin{aligned} A_n(t) &= U_n^T(t + \sqrt{n}(\hat{\theta}_n - \theta_0)) - U_n^T(\sqrt{n}(\hat{\theta}_n - \theta_0)) \\ &\quad + \frac{1}{2} \{(t + \sqrt{n}(\hat{\theta}_n - \theta_0))^T J (t + \sqrt{n}(\hat{\theta}_n - \theta_0)) \\ &\quad - \sqrt{n}(\hat{\theta}_n - \theta_0)^T J \sqrt{n}(\hat{\theta}_n - \theta_0)\} + o_p(1) \\ &= U_n^T t + \frac{1}{2} t^T J t + \sqrt{n}(\hat{\theta}_n - \theta_0)^T J t + o_p(1) \\ &= \frac{1}{2} t^T J t + o_p(1) \end{aligned}$$

where the last line follows because  $\sqrt{n}(\hat{\theta}_n - \theta_0) = -J^{-1}U_n + o_p(1)$ . Since  $A_n$  is convex and converges pointwise in probability to  $t^T J t / 2$ , it converges uniformly in probability on compact subsets by lemma 4.7.1.  $\square$

Now we return to the main proof. Define  $\gamma_n \equiv \inf_{|u|=1} A_n(u)$ . It converges in probability to  $\gamma_0 = \inf_{|u|=1} u^T J u / 2 > 0$ . By the same argument used in lemma 4.7.2 it follows that

$$A_n(t) \geq \gamma_n |t| \quad \text{for } |t| > 1.$$

Now the integrand in the numerator above converges in probability for each fixed  $t$  to

$$t \exp(-t^T J t / 2) \lambda(\theta_0) \exp(-C_2 |\theta_0|),$$

and similarly the integrand of the denominator converges for each fixed  $t$  to

$$\exp(-t^T J t / 2) \lambda(\theta_0) \exp(-C_2 |\theta_0|).$$

Furthermore the domination hypotheses of the following convergence lemma hold for both the numerator and denominator with dominating function

$$D(t) \equiv \begin{cases} 2C_1 & \text{if } |t| \leq 1 \\ C_1 |t| \exp(-\gamma_0 |t| / 2) & \text{if } |t| > 1 \end{cases}.$$

Hence the ratio of integrals converges in probability to

$$\frac{\int t \exp(-t^T J t / 2) \lambda(\theta_0) \exp(-C_2 |\theta_0|) dt}{\int \exp(-t^T J t / 2) \lambda(\theta_0) \exp(-C_2 |\theta_0|) dt} = 0.$$

□

**Lemma 8.6.** Suppose that  $X_n(t), Y_n(t)$  are jointly measure random functions,  $X_n(t, \omega), Y_n(t, \omega)$  for  $(t, \omega) \in K \times \Omega$  where  $K \subset \mathbb{R}^d$  is compact, that  $\lambda$  is a measure on  $\mathbb{R}^d$ , and:

- (i)  $Y_n(t) \rightarrow_p Y(t)$  and  $X_n(t) \rightarrow_p X(t)$  for  $\lambda$  almost all  $t \in K$ .
- (ii)  $\int_K Y_n(t) d\lambda(t) \rightarrow_p \int_K Y(t) d\lambda(t)$  with  $|\int_K Y(t) d\lambda(t)| < \infty$  almost surely.

Then  $\int_K X_n(t) d\lambda(t) \rightarrow_p \int_K X(t) d\lambda(t)$ .

**Proof.** It suffices to show almost sure convergence for some further subsequence of any given subsequence. By convergence in probability for fixed  $t$ , and then dominated convergence

$$\begin{aligned} H_n(\epsilon) &\equiv (P \otimes \lambda)(\{(\omega, t) : |X_n(t, \omega) - X(t, \omega)| > \epsilon\}) \\ &= \int_K P(\{\omega : |X_n(t, \omega) - X(t, \omega)| > \epsilon\}) dt \rightarrow 0 \end{aligned}$$

by the dominated convergence theorem, and similarly for  $\{Y_n\}$ . By replacing  $\epsilon$  by  $\epsilon_n \downarrow 0$  and extraction of subsequences we can find a subsequence  $\{n'\}$  so that the integrals  $\int_K X_n d\lambda$  converge and for some set  $N \subset K \times \Omega$  with  $P \otimes \lambda(N) = 0$  we get convergence for all  $(\omega, t) \in N^c$  of both  $X_{n'}$  and  $Y_{n'}$ . Then we have  $\lambda(\{t : (\omega, t) \in N\}) = 0$  for almost all  $\omega$ , Hence by Fatou's lemma applied to  $Y_{n'} \pm X_{n'}$  we deduce that

$$\int_K X_{n'}(t) d\lambda(t) \rightarrow_{a.s.} \int_K X(t) d\lambda(t).$$

□

**Corollary 8.7.** Suppose that:

- (i)  $G_n(t) \rightarrow_p G(t)$  for each fixed  $t \in \mathbb{R}^d$ .
- (ii)  $P(|G_n(t)| \leq D(t) \text{ for all } t \in \mathbb{R}^d) \rightarrow 1$ .
- (iii)  $\int D(t) dt < \infty$ .

Then  $\int G_n(t) dt \rightarrow_p \int G(t) dt$ .

**Proof.** Let  $\epsilon > 0$ ; choose a compact set  $K = K_\epsilon$  so that  $\int_K D(t) dt < \epsilon$ .

This is possible since  $\int D(t) dt < \infty$ . Then, using  $|G(t)| \leq D(t)$ ,

$$|\int G_n(t) dt - \int G(t) dt|$$

$$\begin{aligned}
 &\leq \left| \int_K G_n(t) dt - \int_K G(t) dt \right| \\
 &\quad + \int_{K^c} |G_n(t)| dt + \int_{K^c} |G(t)| dt \\
 &\leq \left| \int_K G_n(t) 1_{\{|G_n(t)| \leq D(t)\}} dt - \int_K G(t) dt \right| \\
 &\quad + \int_{K \cap \{|G_n(t)| > D(t)\}} |G_n(t)| dt + \int_{K^c \cap \{|G_n(t)| > D(t)\}} |G_n(t)| dt \\
 &\quad + 2 \int_{K^c} D(t) dt
 \end{aligned}$$

Now apply lemma 8.6 with  $Y_n(t) \equiv D(t)$  and

$$X_n(t) \equiv G_n(t) 1_{\{|G_n(t)| \leq D(t)\}}.$$

Then  $X_n(t) \rightarrow_p G(t)$  and the second hypothesis of the lemma holds easily. Hence  $\int X_n(t) dt \rightarrow_p \int G(t) dt$  so that the first term above  $\rightarrow_p 0$ . The second and third terms converge in probability to 0 because with the set over which the integral is taken is empty with arbitrarily high probability; and the third term is bounded by  $2\epsilon$  by choice of  $K$ .  $\square$

## 9. Nonparametric Bayes Estimation