Evaluating the Performance of Estimators (Section 7.3)

**Example:** Suppose we observe  $X_1, \ldots, X_n$  iid  $N(\theta, \sigma_0^2)$ , with  $\sigma_0^2$  known, and wish to estimate  $\theta$ .

Two possible estimators are:

 $\hat{\theta} = \bar{X} \equiv \text{sample mean}$  and  $\hat{\theta} = M \equiv \text{sample median}$ .

Which is better? How to measure performance?

Some possibilities:

- 1. Compare  $E|\bar{X} \theta|$  with  $E|M \theta|$ .
- 2. Compare  $E(\bar{X} \theta)^2$  with  $E(M \theta)^2$ .
- 3. Compare  $EL(\theta, \overline{X})$  with  $EL(\theta, M)$

where  $L(\cdot, \cdot)$  is an appropriate "loss function": the value  $L(\theta, a)$  is some measure of the loss incurred when the true value is  $\theta$  and our estimate is a.

absolute error loss:	$L(\theta, a) =  a - \theta $
squared error loss:	$L(\theta, a) = (a - \theta)^2$
game show loss:	$L(\theta, a) = I( a - \theta  > c)$
Stein's loss (for $\theta, a > 0$ ):	$L(\theta, a) = \frac{a}{\theta} - 1 - \log\left(\frac{a}{\theta}\right)$

Historically, estimators have been most frequently compared using Mean Squared Error:  $MSE(\theta) = E_{\theta}(\hat{\theta} - \theta)^2$ .

This is because the MSE can often be calculated or approximated (for large samples), and has nice mathematical properties.

## Admissible and Inadmissible Estimators

Let W = W(X) be an estimator of  $\tau = \tau(\theta)$ .

Define  $MSE_W(\theta) = E_{\theta}(W - \tau(\theta)^2)$ .

An estimator W is **inadmissible** (w.r.t. squared error loss) if there exists another estimator V = V(X) such that

 $MSE_V(\theta) \leq MSE_W(\theta) \quad \forall \theta \in \Theta$ 

with strict inequality for at least one value of  $\theta$ . (An estimator is inadmissible if there is another estimator that "beats" it.) An estimator which is **not inadmissible** is called **admissible**.

(Draw some pictures.)

Note: An admissible estimator may actually be very bad. An inadmissible estimator can sometimes be pretty good.

Note: If we are using a loss function  $L(\tau, a)$ , we also define inadmissible and admissible estimators in the same way, replacing the MSE by the more general notion of a **risk function**  $R(\theta, W) = E_{\theta}L(\tau(\theta), W(\mathbf{X})).$ 

**Examples:** Again, suppose we observe  $X_1, \ldots, X_n$  iid  $N(\theta, \sigma_0^2)$ , with  $\sigma_0^2$  known, and wish to estimate  $\theta$ .

Consider the estimator  $W \equiv 0$  that always estimates  $\theta$  by 0 regardless of the data X.

This is a very bad estimator, but it is admissible because it is great when  $\theta = 0$ . No non-degenerate estimator V can possibly beat W since it would have to satisfy

$$MSE_V(0) \leq MSE_W(0)$$
  

$$\Rightarrow E_0(V-0)^2 \leq E_0(W-0)^2$$
  

$$\Rightarrow E_0V^2 \leq 0 \Rightarrow P_0(V=0) = 1$$
  

$$\Rightarrow V \equiv 0.$$

Now consider the estimator  $M \equiv$  sample median.

We show later that the sample mean  $\overline{X}$  has a uniformly smaller MSE than M so that M is inadmissible. (The two MSE functions are constant, i.e., flat.)

However, M is not a bad estimate of  $\theta$ , and might be used if there were doubts about the normality assumption (perhaps the true distribution has thicker tails) or concern about outliers.

**Bias, Variance, and MSE** (for an estimator W of  $\tau(\theta)$ )

$$\begin{array}{rcl} \mathsf{Bias}_{W}(\theta) &=& E_{\theta}(W - \tau(\theta))\\ \mathsf{Var}_{W}(\theta) &=& E_{\theta}(W - E_{\theta}W)^{2} \equiv \mathsf{Var}_{\theta}(W) \end{array}$$
  
Fact:  $\mathsf{MSE}_{W}(\theta) &=& \mathsf{Bias}_{W}^{2}(\theta) + \mathsf{Var}_{W}(\theta)\\ \mathsf{MSE} &=& \mathsf{Bias}^{2} + \mathsf{Var} \quad (\mathsf{in brief}) \end{array}$ 

**Proof:** For any rv Y with finite second moment, we know  $EY^2 = (EY)^2 + Var(Y)$ .

Taking  $Y = W - \tau$  leads to

 $MSE = Bias^2 + Var$ 

since  $Var(W - \tau) = Var(W)$ .

**Terminology:** An estimator W with  $Bias_W(\theta) \equiv 0$ , that is,

$$E_{\theta}W = \tau(\theta) \quad \forall \, \theta \in \Theta$$

is said to be **unbiased**. If not, it is **biased**.

For an unbiased estimator, MSE = Var.

Example: Coin-tossing X1,...,Xn Fid Bernoulli(0)  $\hat{\Theta}_{MLE} = \frac{T}{n}, T = \sum X_i.$  $\hat{\Theta}_{Bayes} = (1-p)a + p(\frac{T}{n})$   $\frac{1}{2} \frac{1}{prior} \frac{1}{pr$  $MSE_{MLE}(\Theta) = Blas^{2} + Variance$  $= 0 + \Theta(1-\Theta)$  $MSE_{Bayes}(\Theta) = Bias^2 + VarianQ$  $[(1-p)a + p\theta - \theta]^{2} + p^{2} \Theta(H\theta) - (1-p)\theta$  $(1-p)^{2}(a-\theta)^{2} + p^{2} \Theta(1-\theta)$ 

Which is better? (according to MSE)

Answer: Neither dominates the other.

$$\mathsf{MSE}_{\mathsf{Bayes}}(a) = \frac{p^2\theta(1-\theta)}{n} < \frac{\theta(1-\theta)}{n} = \mathsf{MSE}_{\mathsf{MLE}}(a),$$

 $MSE_{Bayes}(0) = (1-p)^2 a^2 > 0 = MSE_{MLE}(0)$ ,

and similarly,  $MSE_{Bayes}(1) > MSE_{MLE}(1)$ .

Thus, the Bayes estimate is superior in the neighborhood of  $\theta = a$ , and the MLE is superior near  $\theta = 0$  and 1.

Note: both  $MSE_{Bayes}(\theta)$  and  $MSE_{MLE}(\theta)$  are parabolas (quadratic functions of  $\theta$ ).

**Note:** Regarding (in)admissibility, the above remarks prove nothing. But it can be shown that both the Bayes estimate and MLE are admissible here. Typically, Bayes estimates (with proper priors) are admissible.

Example: Estimating the Variance X1,---, Xn iid N(4,02)  $\gamma(\Theta) = \sigma^2$  $W = c \sum_{i=1}^{n} (X_i - \overline{x})^2 = c SS$ 

 $c = \frac{1}{n-1}$  usual unbiased estimator c= t is MLE

Which is better ?

 $\frac{SS}{\sigma^2} \sim \mathcal{R}_{n-1}^2 \xrightarrow{\text{mean} = n-1}_{\text{var} = 2(n-1)}$   $Thus \quad EcSS = c\sigma^2(n-1)$   $VarcSS = c\sigma^2(n-1)$   $WSE_{cSS}(\mu,\sigma^2) = E(cSS-\sigma^2)^2$   $= Bias^2 + Variance$   $= (c\sigma^2(n-1) - \sigma^2)^2 + c^2\sigma^4 \cdot 2(n-1)$   $= \sigma^4 ((c(n-1)-1)^2 + c^22(n-1))$ 

 $MSE_{SS}(\mu,\sigma^2) = \sigma^4 \mathcal{V}(c)$ with  $W(c) = (c(n-1)-1)^2 + 2(n-1)c^2$  $\Psi'(c) = 2(c(n-1)-1)(n-1) + 4(n-1)c$ = 2(n-1) [c(n+1)-1] $\psi^{pp}(c) = \alpha(n-1)(n+1) > 0$ for n≥2  $\Psi'(c) = 0$  when  $c = \frac{1}{n+1}$ . Both  $c = \frac{1}{n-1}$  and  $c = \frac{1}{n}$  give inadmissible estimators

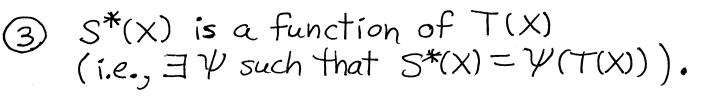
 $c = \frac{1}{n-1}$  is the best c with stein's loss function: (page 351)  $L(\sigma_{2}^{2}a) = \frac{\partial}{\partial z} - \log \frac{\partial}{\partial z} - 1$ Other Plausible 1055 functions  $L(\sigma^{2}, a) = \frac{\sigma^{2}}{a} - \log \frac{\sigma^{2}}{a} - 1$  $= \frac{0^{2}}{3} + \frac{1}{2} - 2$  $= (\log a - \log \sigma^2)^2$ What is the best c for these ? Comments: (state aurlier) Estimator with best c might not be admissible! MSE inappropriate (or dubious, anyway) for estimation of or2.

The Rao-Blackwell Theorem

If T=T(X) is a sufficient statistic for  $\Theta$ ,  $E_{A}S(X) = T(\theta)$  for all  $\theta$ , and  $E_{\Delta}(S(X) - T(\Theta))^2 < \infty$  for all  $\Theta_{\sigma}$ then  $S^*(X) = E(S(X)|T(X))$ satisfies  $E_{\Theta} S^{*}(X) = T(\Theta)$  for all  $\Theta$ , and  $E_{\Theta}(S^{*}(X) - \mathcal{T}(\Theta))^{2} \leq E_{\Theta}(S(X) - \mathcal{T}(\Theta))^{2}$ for all O. Notes:

① E(S(X)|T(X)) does not depend on  $\Theta$ because T(X) is sufficient so that  $\mathcal{L}(X|T)$  (and thus  $\mathcal{L}(S|T)$ ) does not depend on  $\Theta$ .

(2) Equality of MSE's for a particular  $\theta$ can occur iff  $P_{\theta}(S(x) = S^*(X)) = 1$ .



## **Proof:**

Recall: For any rv's X, Y with  $EX^2 < \infty$ , we have EX = E(E(X | Y))Var(X) = E(Var(X | Y)) + Var(E(X | Y))

Now apply these facts:

$$E_{\theta}[S^{*}(X)] = E_{\theta}[E_{\theta}(S(X) | T(X))]$$
  

$$= E_{\theta}S(X) = \tau(\theta)$$
  

$$E_{\theta}(S - \tau)^{2} = \operatorname{Var}_{\theta}(S)$$
  

$$= \underbrace{E[\operatorname{Var}(S | T)]}_{\geq 0} + \operatorname{Var}[\underbrace{E(S | T)}_{S^{*}}]$$
  

$$\geq \operatorname{Var}(S^{*}) = E_{\theta}(S^{*} - \tau(\theta))^{2}$$

Equality can occur only when  $E_{\theta} \operatorname{Var}(S | T) = 0$ . But

$$E_{\theta} \text{Var}(S \mid T) = E_{\theta} \{ E[(S - E(S \mid T))^2 \mid T] \}$$
  
=  $E_{\theta} \{ (S - E(S \mid T))^2 \} = E_{\theta} \{ (S - S^*)^2 \}$   
= 0 iff  $P_{\theta}(S = S^*) = 1$ .

Arguing more loosely,  $E_{\theta} \operatorname{Var}(S | T) = 0 \Rightarrow \operatorname{Var}(S | T) = 0 \Rightarrow$ S is a function of  $T \Rightarrow S^* \equiv E(S | T) = S$ .

Example: 
$$X_1, X_2, ..., X_n$$
 iid Bernaulli(p)  
 $T = \sum X_i$  is suff. stat. for p  
 $\mathcal{L}(X|T)$  puts equal prob. of  $\frac{1}{(T)}$  on  
all strings with T 1's and  
 $n-T$  0's.  
Generate from  $\mathcal{L}(X|T=t)$  by placing  
t 1's and  $n-t$  0's in an urn, and  
randomly drawing (without replacement)  
until the urn is empty.  
**(1)** Estimation of p.  
 $EX_i = p$  for all  $p$  so that  $S = X_i$  is  
an unbiased estimator of  $p$ .  
 $X_i$  is not a function of T, so it can be  
improved by conditioning (as in the  
Rao-Blackwell Theorem).  
 $S^* = E(S|T) = E(X_i|T) = P(X_i=1|T)$   
 $indicator rv$   
 $0-1 valued$   
 $S^* = T/n$  is the usual estimator (the sample  
proportion)!

Clearly  $ES^* = p \forall p$  $Var(s^{*}) = p(1-p)/n < Var S = p(1-p)$ verifying the conclusions of the R-B Thm.  $\forall P$ , (2) Estimation of  $p^2$ . Note that  $E_{X_1X_2} = P(X_1X_2=1)$ indicator  $= P(X_1 = X_2 = 1) = p^2.$ Thus  $S = X_1 X_2$  is an unbiased estimator of  $p^2$ . It is not a function of T, so it can be improved by conditioning.  $S^* = E(X_1X_2|T) = P(X_1X_2 = I|T)$  $= P(X_1 = X_2 = 1 | T) = T_1 \cdot T_{-1}$ By R-B Thm, S\* is an unbiased estimate of p<sup>2</sup> with smaller variance than S. This can be verified by straightforward calculations. For comparison, what is the MLE of p2?

The MLE of p is T/n, so the invariance principle for MLE's says the MLE of p2 is  $\left(\frac{T}{n}\right)^2$ . Clearly  $\frac{T}{n} \cdot \frac{T-1}{n-1}$  and  $\left(\frac{T}{n}\right)^2$  are very close when n is large. which is better? Neither dominates.  $\left(\frac{T}{n}\right)^2$  is biased, but the bias is negligible for large n. Estimation of p<sup>3</sup>, etc.  $S = X_1 X_2 X_3$  is unbiased for  $p^3$  $S^* = E(S|T) = P(X_1 = X_2 = X_3 = I|T)$ =  $T \cdot T - 1 \cdot T - 2$ is the Rao-Blackwell improvement on S. The pattern is now clear for p4, etc.

Suppose T = T(X) is a **complete** and sufficient statistic for  $\theta$ . Then ...

- (1) For any parameter  $\tau(\theta)$ , there is at most one unbiased estimator which is a function of T.
- (2) If S = S(X) is unbiased for  $\tau(\theta)$  (and  $Var(S) < \infty$  for all  $\theta$ ), then

$$S^*(X) = S^* = E(S | T)$$

is the UMVUE for  $\tau(\theta)$ .

Definition: S = S(X) is the UMVUE (uniformly minimum variance unbiased estimator) for  $\tau(\theta)$  if

$$E_{\theta}S = \tau(\theta)$$
 for all  $\theta$ ,

and

 $\operatorname{Var}_{\theta}(S) < \operatorname{Var}_{\theta}(S')$  for all  $\theta$ 

for any other unbiased estimator S'.

Note: For unbiased estimators, MSE = Variance.

Terminology: UMVUE = "best unbiased estimator"

(3) An unbiased estimator (with finite variance) which is a function of T is the UMVUE.

Proof of (D:

Suppose  $S_1(X) = \Psi_1(T(X)) [S_1 = \Psi_1(T)]$   $S_2(X) = \Psi_2(T(X)) [S_2 = \Psi_2(T)]$ and  $ES_1 = ES_2 = T(\Theta) \quad \forall \Theta$ . Then  $E(S_1 - S_2) = O \quad \forall \Theta$ or  $E_{\Theta} g(T) = O \quad \forall \Theta$  where  $g(t) = \Psi_1(t) - \Psi_2(t)$ . By completeness  $P(g(T) = O) = P_{\Theta}(S_1 = S_2) = 1 \quad \forall \Theta$ .

Thus  $S_1 = S_2$  (with prob 1 for all  $\theta$ ).

Proof of (2):  $S^*$  is unbiased and a fn. of T. Suppose W is <u>ony</u> unbiased estimator for  $T(\Theta)$ .

Then R-BThm says that  $W^* = E(W|T)$ is on unbiased estimator of  $T(\Theta)$  and

Var  $W^* \leq Var_{\Theta} W$  for all  $\Theta$ . (+) But  $W^*$  is unbiased and a function of T, so  $\square$  implies  $W^* = S^*$ . Thus (+) implies  $\operatorname{Var}_{\Theta} S^* \leq \operatorname{Var}_{\Theta} W$  for all  $\Theta$ , and S\* is the UMVUE. Proof of 3: Suppose  $E_{\Theta}S(X) = \mathcal{T}(\Theta) \ \forall \Theta$ and  $S(X) = \mathcal{V}(T(X))$ . Then  $S^* = E(S|T)$  is the UMVUE by(2).But S is a function of T so that E(S|T) = S.Thus S is the UMVUE.

**Example:** Observe  $X_1, \ldots, X_n$  iid Bernoulli(p).

• Find the UMVUE of *p*.

 $T = \sum_i X_i$  is a CSS. E(T/n) = p.

Since T/n is an unbiased estimator of p which is a function of the CSS T, it is the UMVUE.

• Find the best unbiased estimator of  $p^2$ .

 $E\left(\frac{T(T-1)}{n(n-1)}\right) = p^2 \left(\begin{array}{c} \text{Found indirectly using Rao-Blackwell.} \\ \text{For a direct argument, see below.} \end{array}\right)$ 

Since  $\left(\frac{T(T-1)}{n(n-1)}\right)$  is an unbiased estimator of  $p^2$  which is a function of the CSS T, it is the UMVUE.

Checking unbiasedness:

$$ET(T-1) = E(T^2) - ET = Var(T) + (ET)^2 - ET$$
  
=  $np(1-p) + (np)^2 - np = n(n-1)p^2$ 

Comment: "Estimate a parameter by its UMVUE" is another approach to estimation, but **not** a very good one. Often, no unbiased estimator exists, or the only one that exists is bad.

**Example:** Observe  $X_1, \ldots, X_n$  iid  $N(\mu, \sigma^2)$  with  $\theta = (\mu, \sigma^2)$  unknown.

Here  $T = (\bar{x}, s^2)$  is a CSS. (Recall the derivation: T is a 1-1 function of the natural SS for a 2pef.)

• Estimation of  $au(\mu,\sigma^2) = \mu$  :

 $\bar{x}$  is unbiased ( $E\bar{x} = \mu$ ) and a function of  $T \Rightarrow \bar{x}$  is UMVUE.

MLE of  $\theta$  is  $\hat{\theta} = (\bar{x}, n^{-1} \sum_{i} (X_i - \bar{x})^2)$ . So invariance principle says MLE of  $\mu$  is  $\tau(\hat{\theta}) = \bar{x}$ .

MOM estimate is also  $\bar{x}$  since  $E\bar{x} = \mu$ .

Note: For estimating  $\mu$ , the MLE, MOM, UMVUE all agree on  $\bar{X}$ . But Bayes estimate is different.

What about the sample median M?

M is an unbiased estimator of  $\mu$ . (Proof?) But it is **not** a function of the CSS T. Thus Rao-Blackwellizing M leads to the UMVUE (which we know is  $\bar{x}$ ) which has a strictly smaller variance than M.

Thus:  $E(M | T) = \overline{x}$  and  $Var(M) > Var(\overline{x}) = \sigma^2/n$ .

## • Estimating $au(\mu, \sigma^2) = \sigma^2$ :

Let  $SS = \sum_i (X_i - \bar{x})^2$ .

 $s^2 = SS/(n-1)$  is an unbiased estimator of  $\sigma^2$  and a function of the CSS T. Therefore  $s^2$  is the UMVUE.

By the invariance principle, the MLE of  $\sigma^2$  is SS/n. This is slightly biased.

• Estimation of  $au(\mu, \sigma^2) = \mu^2$  :

The MLE of  $\mu^2$  is  $(\bar{x})^2$  by invariance of MLE's.  $\bar{x}^2$  is biased for  $\mu^2$ :

$$E(\bar{x}^2) = \operatorname{Var}(\bar{x}) + (E\bar{x})^2 = \frac{\sigma^2}{n} + \mu^2 > \mu^2.$$

An unbiased estimate of  $\mu^2$  is  $W \equiv \bar{x}^2 - \frac{s^2}{n}$ :

$$E\left(\bar{x}^2 - \frac{s^2}{n}\right) = \left(\frac{\sigma^2}{n} + \mu^2\right) - \frac{\sigma^2}{n} = \mu^2.$$

Subtracting  $s^2/n$  removes (or corrects for) the bias in the MLE. W is the UMVUE since it is unbiased and a function of T. Which is better:  $\bar{x}^2$  or W?

For n > 3, W has slightly smaller MSE than  $\bar{x}^2$ . (Verify?) Thus  $\bar{x}^2$  is inadmissible for n > 3 (but is a perfectly reasonable estimator).

But W is also inadmissible because it sometimes takes on ''impossible'' values.

 $\mu^2 \geq 0$ , but W can be negative!

P(W < 0) is positive and will be sizeable when  $\mu$  is small ( $\approx 1/2$  when  $\mu = 0$ ).

A better estimate is clearly  $W_+ = \max(W, 0)$ .

Whenever  $W_+ \neq W$ , we know  $W_+$  is closer to the true value of  $\mu^2$ . More formally

$$E(W - \mu^2)^2 - E(W_+ - \mu^2)^2 = E\left[(W - \mu^2)^2 - (W_+ - \mu^2)^2\right]$$
  
=  $E\left[\underbrace{\{(W - \mu^2)^2 - (0 - \mu^2)^2\}I(W < 0)}_{\text{always} \ge 0 \text{ and sometimes } > 0}\right] > 0$ 

But  $W_+$  is biased! Oh, well.

No unbiased estimator of  $\mu^2$  exists which does **not** take on negative values.

**Fact:** There are situations where there are **no** unbiased estimators (and hence, no UMVUE exists).

**Example:** Observe  $X_1, \ldots, X_n$  iid Poisson( $\lambda$ ). There exists no unbiased estimator of  $1/\lambda$ .