

## Lecture 22: Point Estimation

Today we start Chapter 6 and with it the statistics part of the course. We saw in Lecture 20 (Random Samples) that it frequently occurs that we know a probability distribution except for the value of a parameter. In fact we had three examples

### 1. *The Election Example*

Bin (1, ?)

## 2. The Computer Failure Time Example

Exp (?)

## 3. The Random Number Example

U(0, ?)

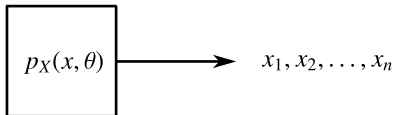
By convention the unknown parameter will be denoted  $\theta$ . So replace ? by  $\theta$  in the three examples. So  $\theta = p$  in example 1 and  $\theta = \lambda$  in Example 2 and  $\theta = B$  (so  $U(0, B)$ ) in Example 3.

If the population  $X$  is discrete we will write its pmf as  $p_X(x, \theta)$  to emphasize that it depends on the unknown parameter  $\theta$  and if  $X$  is continuous we will write its pdf as  $f_X(x, \theta)$  again to emphasize the dependence on  $\theta$ .

### *Important Remark*

$\theta$  is a fixed number, it is just that we don't know it. But we are allowed to make calculations with a number we don't know, that is the first thing we learn to do in high-school algebra, compute with "the unknown  $x$ ".

Now suppose we have an actual sample  $x_1, x_2, \dots, x_n$  from a population  $X$  whose probability distribution is known except for an unknown parameter  $\theta$ . For convenience we will assume  $X$  is discrete.



The idea of point estimation is to develop a theory of making a guess for  $\theta$  (“estimating  $\theta$ ”) in terms of  $x_1, x_2, \dots, x_n$ .  
So the big problem is

## The Main Problem (Vague Version)

What function  $h(x_1, x_2, \dots, x_n)$  of the items  $x_1, x_2, \dots, x_n$  in the sample should we pick to estimate  $\theta$  ?

### Definition

*Any function  $w = h(x_1, x_2, \dots, x_n)$  we choose to estimate  $\theta$  will be called an estimator for  $\theta$ .*

*As first one might ask -*

$$\left. \begin{array}{l} \text{find } h \text{ so that for every sample} \\ x_1, x_2, \dots, x_n \text{ we have} \\ h(x_1, x_2, \dots, x_n) = \theta. \end{array} \right\} (*)$$

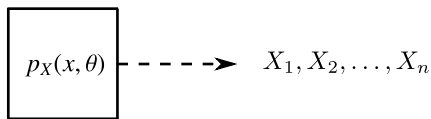
*This is hopelessly naive. Let's try something else*

### *The Main Problem (some what more precise)*

Give quantitative criteria to decide whether one estimator  $w_1 = h_1(x_1, x_2, \dots, x_n)$  for  $\theta$  is better than another estimator  $w_2 = h_2(x_1, x_2, \dots, x_n)$  for  $\theta$ .

The above version, though better, is not precise enough.

In order to pose the problem correctly we need to consider random samples from  $X$ , in other words go back before an actual sample is taken or “go random”.



Now our *function*  $h$  gives rise to a random variable (*statistic*)

$$W = h(X_1, X_2, \dots, X_n)$$

which I will call (for a while) an estimator *statistic*, to distinguish it from the estimator (*number*)  $w = h(x_1, x_2, \dots, x_n)$ . Once we have chosen  $h$  the corresponding estimator statistic will often be denoted  $\hat{\theta}$ .



## Main Problem (third version)

Find an estimator  $h(x_1, x_2, \dots, x_n)$  so that

$$P(h(X_1, X_2, \dots, X_n) = \theta) \quad (**)$$

is maximized

This is what we want but it is too hard to implement - after all we don't know  $\theta$ .

## Important Remark

We have made a huge gain by “going random”. The statement “maximize  $P(h(x_1, x_2, \dots, x_n) = \theta)$ ” does not make sense because  $h(x_1, x_2, \dots, x_n)$  is a fixed real number so either it is equal to  $\theta$  or it is not equal to  $\theta$ . But  $P(h(X_1, X_2, \dots, X_n) = \theta)$  does make sense because  $h(X_1, X_2, \dots, X_n)$  is a random variable.

Now we weaken  $(**)$  to something that can be achieved, in fact achieved surprisingly easily.

## Unbiased Estimators Main Problem (fourth version)

Find an estimator  $w = h(x_1, \dots, x_n)$  so that the expected value  $E(W)$  of the estimator statistic  $W = h(X_1, X_2, \dots, X_n)$  is equal to  $\theta$ .

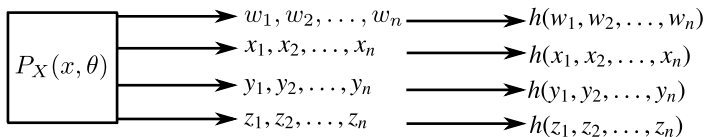
### Definition

*If an estimator  $W$  for an unknown parameter  $\theta$  satisfies  $E(W) = \theta$  then the estimator  $W$  is said to be unbiased.*

Intuitively, requiring  $E(W) = \theta$  is a good idea but we can make this more precise. Various theorems in probability e.g Chebyshev's inequality, tell us that if  $Y$  is a random variable and  $y_1, y_2, \dots, y_n$  are observed values of  $Y$  then the numbers  $y_1, y_2, \dots, y_n$  will tend to be near  $E(Y)$ .

Applying this to our statistic  $W$ - if we take many samples of size  $n$  and compute the value of our estimator  $h$  on each one to obtain many observed values of  $W$  then the resulting numbers will be near  $E(W)$ . But we want these to be near  $\theta$ . So we want

$$E(W) = \theta$$



I have run out of letters. In the above there are four samples of size  $n$  and four corresponding estimates  $h(w_1, \dots, w_n)$ ,  $h(x_1, \dots, x_n)$ ,  $h(y_1, \dots, y_n)$  and  $h(z_1, \dots, z_n)$  for  $\theta$ .

Imagine that instead of four we have one hundred estimates of size  $n$  and one hundred estimates. Then if  $E(W) = \theta$  most of these estimates will be close to  $\theta$ .

## Examples of Unbiased Estimators

Let's take another look at Problems 1 and 2 (pages 1 and 2)

For a Bernoulli random variable  $X \sim \text{Bin}(1, p)$  we have

$$E(X) = p.$$

Hence for the election example, **we are trying to estimate the mean in a Bernoulli distribution.**

For an exponential random variable  $X \sim \text{Exp}(\lambda)$  we have

$$E(X) = \frac{1}{\lambda}.$$

Hence for the Dell computer failure time example, **we are trying to estimate the reciprocal of the mean in an exponential distribution.** One approach is to choose an estimator for the mean, compute it then takes its reciprocal. If we use this approach then the problem again amount estimating the mean.

So in both cases we are trying to estimate the population mean  $E(X) = \mu$

However, in the second case we have to invert the estimate for  $\mu$  to get an estimate for  $\lambda$ .

In fact many other estimation problems amount to estimating the mean in some probability distribution. Accordingly we state this as a general problem.

### Problem

*Find an unbiased estimator for the population mean  $\mu$*

$$\boxed{\begin{array}{l} \theta = \mu \\ P_X(x, \theta) \end{array}} \dashrightarrow X_1, X_2, \dots, X_n$$

So we want  $h(x_1, x_2, \dots, x_n)$  so that

$$E(h(X_1, X_2, \dots, X_n)) = \mu$$

= the population mean.

Amazingly there is a very simple solution to this problem no matter what the underlying distribution is

### Theorem

*The sample mean  $\bar{X}$  is an unbiased estimator of the population mean  $\mu$ ; that is*

$$E(\bar{X}) = \mu$$

### Proof

*The proof is so simple, deceptively simple because the theorem is so important.*

$$\begin{aligned} E(\bar{X}) &= E\left(\frac{X_1 + \dots + X_n}{n}\right) \\ &= \frac{1}{n} (E(X_1) + \dots + E(X_n)) \end{aligned}$$

## Proof (Cont.)

But  $E(X_1) = E(X_2) = \dots = E(X_n) = \mu$  because all the  $X_i$ 's are samples from the population so they have the same distribution as the population so

$$\begin{aligned} E(\bar{X}) &= \frac{1}{n} (\underbrace{\mu + \mu + \dots + \mu}_{n \text{ times}}) \\ &= \frac{1}{n} (n\mu) \\ &= \mu \end{aligned}$$

□

There is lots of other unbiased estimators of  $\mu$  for any population. It is  $X_1$ , the first sample item (or any  $X_i$ ,  $1 \leq i \leq n$ ). This is because, as noted above,

$$E(X_1) = E(X_i) = E(X) = \mu, 1 \leq i \leq n.$$

For the problem of estimating  $p$  in  $\text{Bin}(1, p)$  we have

$$\bar{x} = \frac{\text{number of observed successes}}{n}$$

Since each of  $x_1, x_2, \dots, x_n$  is either 1 or 0 so

$$x_1 + x_2 + \dots + x_n = \# \text{ of } 1\text{'s.}$$

is the number of “successes” (voters who say “Trump” in 2020 (I am joking)) so

$$\bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n)$$

is the relative number of observed successes. This is the “common sense” estimator.

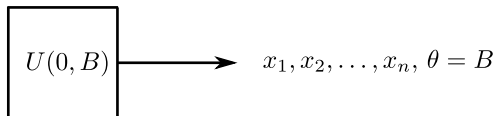


## An Example Where the “Common Sense” Estimator is Biased

Once we have a *mathematical* criterion for an estimator to be good we will often find to our surprise that “common sense” estimators do not meet this criterion.

We saw an example of this in the “Pandemonium jet fighter” Section 6.1, problem 14,(on page 263).

Another very similar problem occurs in Example 3 - estimate  $B$  from the uniform distribution  $U(0, B)$ .



The “common sense” estimator for  $B$  is  $w = \max(x_1, x_2, \dots, x_n)$ , the biggest number you observe. But it is intuitively clear that this estimate will be too small since it only gives the right answer if one of the  $x_i$ 's is equal to  $B$ . So the common sense estimator  $W = \max(x_1, x_2, \dots, x_n)$  is biased.

$$E(\text{Max}(X_1, \dots, X_n)) \underset{\neq}{<} B$$

Amazingly, if you do problem 32, page 274 you will see exactly by how much it *undershoots the mark*. We did this in class.

### Theorem

$$E(\text{Max}(X_1, X_2, \dots, X_n)) = \frac{n}{n+1}B$$

so  $\left(\frac{n+1}{n}\right) \text{Max}(X_1, X_2, \dots, X_n)$  is unbiased.

*Mathematics trumps common sense.*

## Minimum Variance Unbiased Estimators

We have seen that  $\bar{X}$  and  $X_1$  are both unbiased estimators of the population mean for any distribution. Common sense tells us that  $\bar{X}$  is better since it uses all the elements of the sample whereas  $X_1$  just uses one element of the sample (the first).

What mathematical criterion separates them. We have

$$V(X_1) = \sigma^2 = \text{the population variance}$$

$$V(\bar{X}) = \frac{\sigma^2}{n}$$

so if  $n$  is large then

$V(\bar{X})$  is a lot smaller than  $V(X_1)$ .

We will now see why small variance is good. First we state this as a general principle.

### The Principle of Minimum Variance Unbiased Estimation

Among all estimators of  $\theta$  that are unbiased, choose one that has minimum variance.

The resulting estimator is called a minimum variance unbiased estimator, MVUE.

## Theorem 1

$\bar{X}$  is a minimum variance unbiased estimator for the problems of

1. Estimating  $p$  in  $\text{Bin}(1, p)$
2. Estimating  $\mu$  in  $N(\mu, \sigma^2)$

Why is it good to minimize the variance?

We will now see why, assuming the estimator  $\hat{\theta}$  is unbiased.

Suppose  $\hat{\theta} = h(X_1, X_2, \dots, X_n)$  is an estimator statistic for an unknown parameter  $\theta$ .

### Definition

The mean squared error  $MSE(\hat{\theta})$  of the estimator  $\hat{\theta}$  is defined by

$$MSE(\hat{\theta}) = E((\hat{\theta} - \theta)^2)$$

so

$$MSE(\hat{\theta}) = \int \dots \int_{\mathbb{R}^n} (h(x_1, \dots, x_n) - \theta)^2 f_{X_1}(x_1) \dots f_{X_n}(x_n) dx_1 dx_2, \dots, dx_n.$$

or

$$= \sum_{\text{all } x_1, \dots, x_n} (h(x_1, \dots, x_n) - \theta)^2 P(X_1 = x_1) \dots P(X_n = x_n)$$

So  $MSE(\hat{\theta})$  is the square of the error  $h(x_1, x_2, \dots, x_n) - \theta$  of the estimate of  $\theta$  by  $\hat{\theta} = h(x_1, x_2, \dots, x_n)$  averaged over all  $x_1, x_2, \dots, x_n$ .

Obviously we want to minimize the mean squared error (after all it does measure an error). Here is the point - **if  $\hat{\theta}$  is unbiased this is the same as minimizing the variance  $V(\hat{\theta})$** . We now prove the last statement.

### Theorem

*If  $\hat{\theta}$  is unbiased then*

$$MSE(\hat{\theta}) = V(\hat{\theta})$$

*This is amazingly easy to prove.*

## Proof.

By definition

$$MSE(\hat{\theta}) = E((\hat{\theta} - \theta)^2).$$

But if  $\hat{\theta}$  is unbiased then  $E(\hat{\theta}) = \theta$  so

$$MSE(\hat{\theta}) = E((\hat{\theta} - E(\hat{\theta}))^2)$$

By definition the RHS is  $V(\hat{\theta})$ . □

Here is an important definition used a lot in the text. I essentially copied the definition that is in the text, on page 259.

## Definition (text page 259)

*The standard error of the estimator  $\hat{\theta}$ , denoted  $\sigma_{\hat{\theta}}$  is  $\sqrt{V(\hat{\theta})}$ . If the standard error itself contains unknown parameters whose values can be estimated, substitution of these estimates into  $\sigma_{\hat{\theta}}$  yields the **estimated standard error** denoted  $s_{\hat{\theta}}$*