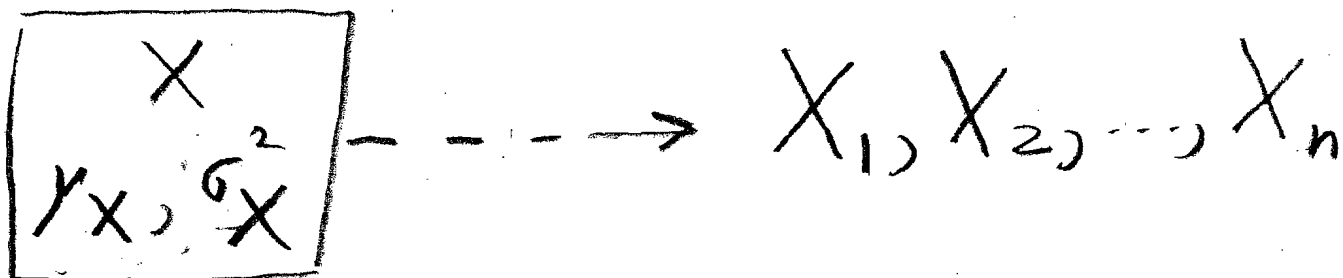


# Lecture 21

## The Sample Total and Mean and The Central Limit Theorem

### 1. Statistics and Sampling Distributions

Suppose we have a random sample from some population with mean  $\mu_X$  and variance  $\sigma_X^2$



and a function  $w = h(x_1, x_2, \dots, x_n)$  of  $n$  variables. Then (as we know) the combined random variable

$W = h(X_1, X_2, \dots, X_n)$   
is called a statistic.

If the population random variable  $X$  is discrete then  $X_1, X_2, \dots, X_n$  will all be discrete and since  $W$  is a combination of discrete random variables it too will be discrete.

The #64,000 question

How is  $W$  distributed?

More precisely, what is the pmf  $P_W(x)$  of  $W$ .

The distribution  $P_W(x)$  of  $W$  is called a "sampling distribution".

Similarly if the population random variable  $X$  is continuous we want to compute the pdf  $f_{W}(x)$  of  $W$ . (now it is continuous)

We will jump to § 5.5

3

The most common  $h(x_1, \dots, x_n)$  is a linear function

$$h(x_1, x_2, \dots, x_n) = a_1 x_1 + \dots + a_n x_n$$

where

$$W = a_1 X_1 + a_2 X_2 + \dots + a_n X_n$$

Proposition 1.1 (pg 219)

Suppose  $W = a_1 X_1 + \dots + a_n X_n$ .

Then

$$(i) \quad E(W) = E(a_1 X_1 + \dots + a_n X_n)$$

$$= a_1 E(X_1) + \dots + a_n E(X_n)$$

(ii) If  $X_1, X_2, \dots, X_n$  are independent then

$$V(a_1 X_1 + \dots + a_n X_n) = a_1^2 V(X_1) + \dots + a_n^2 V(X_n)$$

$$(so \quad V(cX) = c^2 V(X))$$

Now suppose  $X_1, X_2, \dots, X_n$  are a random sample from a population of mean  $\mu$  and variance  $\sigma^2$  so

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

$$V(X_i) = V(X) = \sigma^2, \quad 1 \leq i \leq n$$

and  $X_1, X_2, \dots, X_n$  are independent.

We recall

$$T_0 \equiv \text{the sample total} = X_1 + \dots + X_n$$

$$\bar{X} = \text{the sample mean} = \frac{X_1 + \dots + X_n}{n}$$

As an immediate consequence of the previous proposition we have

### Proposition B

Suppose  $X_1, X_2, \dots, X_n$  is a random sample from a population of mean  $\mu_X$  and variance  $\sigma_X^2$ . Then

$$(i) \quad E(T_0) = n \mu_X$$

$$(ii) \quad V(T_0) = n \sigma_X^2$$

$$(iii) \quad E(\bar{X}) = \mu_X$$

$$(iv) \quad V(\bar{X}) = \frac{\sigma_X^2}{n}$$

Proof (this is important)

6

$$i) E(T_0) = E(X_1 + \dots + X_n)$$

by the Prop.

$$= E(X_1) + \dots + E(X_n)$$

why

$$= \underbrace{\mu_X + \dots + \mu_X}_{n \text{ copies}}$$

$$= n \mu_X$$

$$ii) V(T_0) = V(X_1 + \dots + X_n)$$

by the Prop.

$$= V(X_1) + \dots + V(X_n)$$

$$= \sigma_X^2 + \dots + \sigma_X^2$$

$$= n \sigma_X^2$$

$$(iii) E(\bar{X}) = E\left(\frac{1}{n}(X_1 + \dots + X_n)\right)$$

$$= \frac{1}{n} E(X_1 + \dots + X_n)$$

by (i)

$$= \frac{1}{n} (n\mu_X)$$

$$= \mu_X$$

$$(iv) V(\bar{X}) = V\left(\frac{1}{n}(X_1 + \dots + X_n)\right)$$

by the Prop

$$= \frac{1}{n^2} V(X_1 + \dots + X_n)$$

by (ii)

$$= \frac{1}{n^2} (n\sigma_X^2)$$

$$= \frac{\sigma_X^2}{n}$$

□

# Remark

It is important to understand the symbols —  $\mu_X$  and  $\sigma_X^2$  are the mean and variance of the underlying population. In fact they are called the population mean and the population variance. Given a statistic  $W = h(X_1, \dots, X_n)$  we would like to compute  $E(W) = \mu_W$  and  $V(W) = \sigma_W^2$  in terms of the population mean  $\mu_X$  and



population variance  $\sigma^2_X$ .

9

So we solved this problem

for  $W = \bar{X}$  namely

$$\mu_{\bar{X}} = \mu_X$$

$$\text{and } \sigma^2_{\bar{X}} = \frac{1}{n} \sigma^2_X$$

Never confuse population  
quantities with sample  
quantities

# Corollary

$\sigma_{\bar{X}}$  = the standard deviation of  $\bar{X}$

$$= \frac{\sigma_X}{\sqrt{n}} = \frac{\text{population standard deviation}}{\sqrt{n}}$$

## Proof

$$\sigma_{\bar{X}} = \sqrt{V(\bar{X})}$$

$$= \sqrt{\frac{\sigma_X^2}{n}}$$

$$= \frac{\sqrt{\sigma_X^2}}{\sqrt{n}} = \frac{\sigma_X}{\sqrt{n}}$$

## 2. Sampling from a Normal 11

### Distribution

Theorem LCN (Linear combination of normal is normal)

Suppose  $X_1, X_2, \dots, X_n$

are independent and

$X_1 \sim N(\mu_1, \sigma_1^2), \dots, X_n \sim N(\mu_n, \sigma_n^2)$ .

Let  $W = a_1 X_1 + \dots + a_n X_n$ . Then

$W \sim N(a_1 \mu_1 + \dots + a_n \mu_n, a_1^2 \sigma_1^2 + \dots + a_n^2 \sigma_n^2)$

Proof At this stage

we can't prove  $W$  is normal

(we could if we have moment)

generating functions available). 12

But we can compute the mean and variance of  $W$  using Proposition L.

$$\begin{aligned} E(W) &= E(a_1 X_1 + \dots + a_n X_n) \\ &= a_1 E(X_1) + \dots + a_n E(X_n) \\ &= a_1 \mu_1 + \dots + a_n \mu_n \end{aligned}$$

and

$$\begin{aligned} V(W) &= V(a_1 X_1 + \dots + a_n X_n) \\ &= a_1^2 V(X_1) + \dots + a_n^2 V(X_n) \\ &= a_1^2 \sigma_1^2 + \dots + a_n^2 \sigma_n^2 \end{aligned}$$



Now we can state the theorem we need.

### Theorem N

Suppose  $X_1, X_2, \dots, X_n$  is a random sample from  $N(\mu, \sigma^2)$

$$\boxed{X \sim N(\mu, \sigma^2)} \text{ --- } \rightarrow X_1, X_2, \dots, X_n$$

Then  $T_0 \sim N(n\mu, n\sigma^2)$

and  $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

Proof The hard part is that  $T_0$  and  $\bar{X}$  are normal (this is Theorem LCN)

You show the mean of  $\bar{X}$  is  $\mu$  using either Proposition M or Theorem LCN and the same for showing the variance of  $\bar{X}$  is  $\frac{\sigma^2}{n}$ .  $\square$

Remark

It is very important for statistics that the sample variance

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

satisfies  $S^2 \sim \chi^2(n-1)$ .

This is one reason that the chi-squared distribution is so important.

### 3. The Central Limit Theorem (3.5.4) 15

In Theorem N we saw that if we sampled  $n$  times from a normal distribution with mean  $\mu$  and variance  $\sigma^2$  then

i)  $T_0 \sim N(n\mu, n\sigma^2)$

ii)  $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$

So both  $T_0$  and  $\bar{X}$  are still normal

The Central Limit Theorem says that if we sample  $n$  times with  $n$  large enough from any distribution with mean  $\mu$  and variance  $\sigma^2$  then

$T_0$  has approximately  $N(n\mu, n\sigma^2)$  distribution and  $\bar{X}$  has approximately  $N(\mu, \frac{\sigma^2}{n})$  distribution.

We now state the CLT.

16

## The Central Limit Theorem

$$\boxed{X, \mu, \sigma^2}$$

$$\text{---} \rightarrow X_1, X_2, \dots, X_n$$

$$\bar{X} \approx N(\mu, \sigma^2) \text{ provided } n > 30$$

### Remark

This result would not be satisfactory to professional mathematicians because there is no estimate of the error involved in the approximation



However an error estimate 17  
is known - you have to take a  
more advanced course.

The  $n > 30$  is a "rule of thumb".  
In this case the error will be  
negligible up to a large number of  
decimal places (but I don't know  
how many).

So the Central Limit Theorem  
says that for the purposes of  
sampling if  $n > 30$  then the  
sample mean behaves as if  
the sample were drawn from  
a NORMAL population with the  
same mean and variance of the  
actual population.

Example 5-27

A certain consumer organization reports the number of major defects for each new automobile that it tests.

Suppose that the number of such defects for a certain model is a random variable with mean 3.2 and standard deviation 2.4. Among 100 randomly selected cars of this model what is the probability that the average number of defects exceeds 4.

Solution

Let  $X_i = \#$  of defects for the  $i$ -th car

$$\boxed{\begin{array}{l} X \\ \mu=3.2, \sigma=2.4 \end{array}} \text{---} \rightarrow X_1, X_2, \dots, X_{100}$$

$n=100 > 30$  so we can use the CLT

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_{100}}{100}$$

so  $\bar{X}$  = average number of defects

So we want

$$P(\bar{X} > 4)$$

$$\text{Now } E(\bar{X}) = \mu = 3.2$$

20

$$V(\bar{X}) = \frac{\sigma^2}{n} = \frac{(2.4)^2}{100}$$

$$\text{Let } Y \sim N(3.2, \frac{(2.4)^2}{100})$$

By the CLT  $\bar{X} \approx Y$  so

$$\sigma_Y = \frac{2.4}{10} = .24$$

$$P(\bar{X} \geq 4) \approx P(Y \geq 4)$$

don't use  
correction  
for  
continuity

$$= P\left(\frac{Y - 3.2}{.24} \geq \frac{4 - 3.2}{.24}\right)$$

$$= P\left(Z \geq \frac{.8}{.24}\right) \quad 3.33$$

$$\approx 1 - \Phi(3.33) = 1 - .9996$$

$$= .0004$$

# How the Central Limit 21

## Theorem Gets Used More Often

The CLT is much more useful than one would expect. That is because many well-known distributions can be realized as sample totals of a sample drawn from another distribution.

I will state this as

### General Principle

Suppose a random variable  $W$  can be realized as a sample total  $W = T_n = X_1 + \dots + X_n$  from some  $X$  and  $n > 30$ .

Then  $W$  is approximately normal

## Examples (This isn't obvious)

1.  $W \sim \text{Bin}(n, p)$  with  $n$  large
2.  $W \sim \text{Gamma}(\alpha, \beta)$  with  $\alpha$  large
3.  $W \sim \text{Poisson}(\lambda)$  with  $\lambda$  large

We will do the example  
of  $W \sim \text{Bin}(n, p)$  and  
recover (more or less) the  
normal approximation to the  
binomial to

CLT  $\Rightarrow$  normal approx  
to binomial.

The point is

23

Theorem (sum of binomials is binomial)

Suppose  $X$  and  $Y$  are independent,  $X \sim \text{Bin}(m, p)$  and  $Y \sim \text{Bin}(n, p)$ . Then

$$W = X + Y \sim \text{Bin}(m+n, p)$$

"Proof"

For simplicity we will assume  $p = \frac{1}{2}$ .

Suppose Fred tosses a fair coin  $m$  times and then Jack tosses a fair coin  $n$  times.

Let  $X = \#$  of heads Fred observes

$Y = \#$  of heads Jack observes

So  $X \sim \text{Bin}(m, \frac{1}{2})$  and  $Y \sim \text{Bin}(n, \frac{1}{2})$

What is  $X+Y$ ?

Forget who was doing the tossing,  $X+Y$  is just the total number of heads in  $m+n$  tosses of a fair coin so

$X+Y \sim \text{Bin}(m+n, \frac{1}{2})$ .





Now suppose we have

$$\boxed{X \sim \text{Bin}(1, p)} \longrightarrow X_1, \dots, X_n$$

Then  $X_i \sim \text{Bin}(1, p)$   $\forall 1 \leq i \leq n$

so

$$T_0 = X_1 + X_2 + \dots + X_n \sim \text{Bin}(n, p)$$

Now if  $n > 30$  we know

$T_0$  is approximately normal so

if  $W \sim \text{Bin}(n, p)$  and  $n > 30$

then  $W \approx$  normal

$$E(W) = np \text{ and } V(W) = npq \text{ AND}$$

$$W \approx N(np, npq)$$

So we get the normal approximation to the binomial (with  $n > 30$  replacing  $np \geq 10$  and  $nq \geq 10$ )

### Remark

If  $p = \frac{1}{2}$  then the second

conditions give  $n > 20$

— so better than CLT but

if  $p = \frac{1}{5}$  then the second

conditions give  $n > 50$

— so worse than the CLT