

MAT 2377 3X (Spring 2011)

Sampling Distributions and Point Estimation of  
Parameters

Refer to Sections 7-1, 7-2, 7-3

**Random Sampling Terminology :** In practice we often repeat an experiment  $n$  times.

**Examples :**

- Randomly selecting  $n$  items from a production line and verifying the conformity of each of the selected items.
- Measuring the current in a thin copper wire  $n$  times.

We will model each of the observations as a random variable, i.e.  $X_i$  will represent the  $i$ th observation.

**Assumptions :**

1. We will assume **independent** trials.
2. Furthermore, we are repeating the same experiment  $n$  times, hence the random variables  $X_1, \dots, X_n$  are **identically distributed**. This common distribution is called the **population**.

**Random Sample :** We call the following collection of independent and identically distributed random variables  $X_1, \dots, X_n$  a **random sample**.

**Population Parameters :** Here are some population parameters that often need to estimate in practice :

- The mean of a population  $\mu$ , called the **population mean**.
- The variance of a population  $\sigma^2$ , called the **population variance**.
- The proportion  $p$  of items in a population that satisfy an attribute, called the **population proportion**.

**Estimation :** A function of a random sample is called a **statistic**. In practice we use a statistic

$$\hat{\Theta} = h(X_1, X_2, \dots, X_n)$$

to estimate an unknown population parameter  $\theta$ , then we say that  $\hat{\Theta}$  is a point estimator of  $\theta$ . Note that  $\hat{\Theta}$  is a random variable. The observed value of  $\hat{\Theta}$  which is

$$\hat{\theta} = h(x_1, x_2, \dots, x_n)$$

is called a point estimate of  $\theta$ .

Here are some common statistics.

1. The **sample mean**  $\bar{X}$  is defined as that average of  $X_1, \dots, X_n$ , that is

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}.$$

**Remarks :**

- We will use  $\bar{X}$  as a point estimator of the population mean  $\mu$ .
- The observed value of  $\bar{X}$  which is denoted  $\bar{x}$  is a point estimate for  $\mu$ .

2. The **sample variance**  $S^2$  is defined as follows :

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

**Remarks :**

- We will use  $S^2$  as a point estimator of the population variance  $\sigma^2$ .
- The observed value of  $S^2$  which is denoted  $s^2$  is a point estimate for  $\sigma^2$ .

3. The **sample standard deviation**  $S$  is defined as  $S = \sqrt{S^2}$ .

**Remarks :**

- We will use  $S$  as a point estimator of the population standard deviation  $\sigma$ .
- The observed value of  $S$  which is denoted  $s$  is a point estimate for  $\sigma$ .

4. The **sample proportion**  $\hat{P}$  is defined as follows :

$$\hat{P} = \frac{X}{n},$$

where  $X$  is the number of selected items that satisfy the attribute of interest among  $n$  items.

**Remarks :**

- We will use  $\hat{P}$  as a point estimator of the population proportion  $p$ .
- The observed value of  $\hat{P}$  which is denoted  $\hat{p} = x/n$  is a point estimate of  $p$ .

**Example 1 :** Consider the following data that represent the life of packaged magnetic disks exposed to corrosive gases (in hours) :

4, 86, 335, 746, 195

(a) Give point estimates for the mean and the standard deviation of the lifetime of a packaged magnetic disk exposed to corrosive gases.

(b) Let  $p$  denote the probability that a packaged magnetic disk exposed to corrosive gases will have a life of more than 100 hours. Find a point estimate for  $p$ .

### Expectations and Variances of a Linear Function of the Random Sample

We will give a result from Section 5-5. These results will allow us to study the properties of some estimators.

**Expectation :** Consider the following random variables  $X_1, \dots, X_n$  and the following constants  $c_0, c_1, \dots, c_n$ . If  $Y$  is defined as

$$Y = c_0 + c_1 X_1 + \dots + c_n X_n,$$

then it can be shown that the expected value of  $Y$  is

$$E[Y] = c_0 + c_1 E[X_1] + \dots + c_n E[X_n].$$

**Variance :** Consider the following random variables  $X_1, \dots, X_n$  and the following constants  $c_0, c_1, \dots, c_n$ . If  $Y$  is defined as

$$Y = c_0 + c_1 X_1 + \dots + c_n X_n,$$

and we assume that  $X_1, \dots, X_n$  are **independent**, then it can be shown that the variance of  $Y$  is

$$\sigma_Y^2 = c_1^2 \sigma_{X_1}^2 + \dots + c_n^2 \sigma_{X_n}^2.$$

**Corollary :** Let  $X$  be a random variable and  $a$  and  $b$  constants, then

$$E[aX + b] = a E[X] + b \quad \text{and} \quad V[aX + b] = a^2 \sigma_X^2.$$

**Example 2 :**

(a) Consider a random sample  $X_1, \dots, X_n$  from a population with mean  $\mu$  and variance  $\sigma^2$ . Show that the sample mean has the following mean and variance :

$$\mu_{\bar{X}} = \mu \quad \text{and} \quad \sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}.$$

(b) Let  $\hat{P} = X/n$  be the sample proportion from a population with probability of success  $p$ . Show that

$$E[\hat{P}] = p \quad \text{and} \quad \sigma_{\hat{P}}^2 = \frac{p(1-p)}{n}.$$

**Terminology :** The standard deviation of a point estimator is called the standard error of the estimate. So the standard error of the mean and of the sample proportion are respectively :

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} \quad \text{and} \quad \sigma_{\hat{P}} = \sqrt{\frac{p(1-p)}{n}}.$$

**Remark :** The standard error usually involves an unknown parameter, so in practice when giving data we replace this unknown parameter by its point estimate. So the estimated standard error of the mean and of the sample proportion become in this case :

$$\hat{\sigma}_{\bar{X}} = \frac{s}{\sqrt{n}} \quad \text{and} \quad \hat{\sigma}_{\hat{P}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

**Example 3 :** Refer to Example 1. Give the estimated standard errors for the mean from part (a) and the sample proportion from part (b).

## Sampling Distribution

The distribution of a statistic  $h(X_1, \dots, X_n)$  is known as a sampling distribution. The following two theorems concern the sampling distribution of the sample mean  $\bar{X}$ .

**Theorem :** Let  $X_1, \dots, X_n$  be **independent normal** random variables such that  $X_i$  follows a  $N(\mu_i, \sigma_i^2)$  distribution for  $i = 1, \dots, n$ . If we define  $Y$  as

$$Y = c_0 + c_1 X_1 + \dots + c_n X_n,$$

then it can be shown that  $Y$  follows a  $N(\mu_Y, \sigma_Y^2)$  distribution.

**Corollary :** Let  $X_1, \dots, X_n$  be a random sample from a  $N(\mu, \sigma^2)$  population, then the sample mean  $\bar{X}$  follows a  $N(\mu, \sigma^2/n)$  distribution.

**Remark :** So now we know that if the population is normally distributed, then the sample mean has a normal sampling distribution. Now in general a population is not necessarily normally distributed. The following theorem tells us that as long as we collect a large number of observations, then the sample mean is approximately normally distributed regardless of the underlying distribution.

**Central Limit Theorem :** Let  $X_1, \dots, X_n$  be a random sample from a population with mean  $\mu$  and variance  $\sigma^2$ . Define  $Z$  as

$$Z = \frac{\bar{X} - \mu_{\bar{X}}}{\sigma_{\bar{X}}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}.$$

Then, as  $n \rightarrow \infty$ , it can be shown that  $Z$  follows a **standard normal distribution**.

**Rule of thumb :** We will take  $n \geq 30$ , as large enough. That is, if  $n \geq 30$ , then  $\bar{X}$  follows approximately a  $N(\mu, \sigma^2/n)$  distribution.

**Example 4 :** The compressive strength of concrete has a mean of 2500 psi and a standard deviation of 50 psi. What is the probability that the mean of 40 specimens will be less than 2490 psi ?

**Comparing means from independent populations :** Let  $\bar{X}_i$  be the sample mean from a  $N(\mu_i, \sigma_i^2)$  population, for  $i = 1, 2$ . Let  $n_i$  be the size of the  $i$ th random sample.

We know that

$\bar{X}_i$  follows a  $N(\mu_i, \sigma_i^2/n_i)$  distribution.

Since  $\bar{X}_1 - \bar{X}_2$  is a linear function of independent normal random variables, then it is also normally distributed. Hence,

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}$$

is a standard normal random variable. **Remark :** By the central limit theorem, if  $n_1 \geq 30$  and  $n_2 \geq 30$ , then  $Z$  is approximately a standard normal random variable.

**Example 5 :** The effective life of a component in a jet-turbine aircraft is normally distributed with mean 5000 hours and standard deviation 40 hours. Suppose that when using an improved manufacturing process the effective is normally distributed with mean 5050 and standard deviation 30. Let  $\bar{X}_1$  be the average life of  $n_1 = 20$  components produced under the old manufacturing process and let  $\bar{X}_2$  be the average life of  $n_1 = 15$  components produced with the new manufacturing process. Find the following probability :

$$P(\bar{X}_1 - \bar{X}_2 > 10).$$

### Properties of Point Estimators - Refer to Section 7-3

**Definition :** Let  $\hat{\Theta}$  be a point estimator for the population parameter  $\theta$ . We say that it is an **unbiased** estimator if

$$E[\hat{\Theta}] = \theta.$$

If the estimator is not unbiased, then we define the difference

$$E[\hat{\Theta}] - \theta$$

is called the bias.

**Remark :** We have already shown that

$$E[\bar{X}] = \mu \quad \text{and} \quad E[\hat{P}] = p.$$

Hence the sample mean is an unbiased estimator for  $\mu$  and the sample proportion is an unbiased estimator for the population proportion  $p$ .

**Example 6 :** Show that the sample variance

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

is an unbiased estimator for the population variance  $\sigma^2$ , that is show that  $E[S^2] = \sigma^2$ .

### Comparing Unbiased Estimators

Suppose that  $\hat{\Theta}_1$  and  $\hat{\Theta}_2$  are unbiased estimators for the parameter  $\theta$ . We will say that  $\hat{\Theta}_1$  is best for estimating  $\theta$ , if

$$V[\hat{\Theta}_1] < V[\hat{\Theta}_2].$$

That is, we will prefer the estimator with the smaller variance.

**Example 7 :** Let  $X_1, X_2, X_3$  be a random sample from a population with mean  $\mu$  and variance  $\sigma^2$ .

$$\hat{\mu}_1 = \frac{X_1 + X_3}{2} \quad \text{and} \quad \hat{\mu}_2 = X_2.$$

- (a) Show that both estimators are unbiased estimators for  $\mu$ .
- (b) Which is best for estimating  $\mu$ ?

**Definition :** The **mean squared error** of an estimator  $\hat{\Theta}$  for the parameter  $\theta$  is defined as

$$\text{MSE}(\hat{\Theta}) = E[(\hat{\Theta} - \theta)^2] = V(\hat{\Theta}) + (\text{bias})^2.$$

**Remarks :**

- If the estimator is unbiased (i.e. bias=0) then the mean squared error of the estimator is equal to the variance of the estimator.
- We can use the mean squared error to compare estimators. The relative efficiency of  $\hat{\Theta}_2$  compared to  $\hat{\Theta}_1$  is defined as

$$\frac{\text{MSE}(\hat{\Theta}_1)}{\text{MSE}(\hat{\Theta}_2)}.$$

If the relative efficiency is less than 1, then we conclude that  $\text{MSE}(\hat{\Theta}_1)$  is a more efficient estimator of  $\theta$  than  $\text{MSE}(\hat{\Theta}_2)$ .

**Example 8 :** Let  $X_1, \dots, X_{20}$  be a random sample from a population with mean  $\mu = 10$  and variance  $\sigma^2 = 5$ . Consider the following two estimators of  $\mu$  :

$$\hat{\mu}_1 = \frac{X_1 + \dots + X_{19}}{20} \quad \text{and} \quad \hat{\mu}_2 = \frac{X_1 + X_{20}}{2}.$$

Which of the two estimators is more efficient estimator of  $\mu$ ?