# Section 5.2 - The Sampling Distribution of a Sample Mean

Statistics 104

Autumn 2004



Copyright ©2004 by Mark E. Irwin

## The Sampling Distribution of a Sample Mean

Example: Quality control check of light bulbs

Sample n light bulbs and look at the average failure time.

Take another sample, and another, and so on.

What is the mean of the sampling distribution? Variance? Standard Deviation?

What is 
$$P[\bar{X}_n \ge \mu_x]$$
 or  $P[\mu_x - c \le \bar{X}_n \le \mu_x + c]$ ?

We will consider the situation where the observations are independent (at least approximately). In the case of finite populations, we want  $N \gg n$ .

Under this assumption, we can use the same idea as we did to get the mean and variance of a binomial distribution.

$$\bar{X} = \frac{1}{n}(X_1 + X_2 + \dots + X_n)$$

$$\mu_{\bar{x}} = \frac{1}{n}\underbrace{(\mu_x + \mu_x + \dots + \mu_x)}_{n \text{ times}} = \mu_x$$

$$\sigma_{\bar{x}}^2 = \frac{1}{n^2}\underbrace{(\sigma_x^2 + \sigma_x^2 + \dots + \sigma_x^2)}_{n \text{ times}} = \frac{\sigma_x^2}{n}$$

$$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$$

So the sampling distribution of  $\overline{X}$  is centered at the same place as the observations, but less spread out as we should expect based on the law school example.

The smaller spread also agrees with the law of large numbers, which says  $\bar{X} \rightarrow \mu_x$  as increases.

Note the sampling distribution of  $\hat{p}$  is just a special case.  $\hat{p}$  is just an average of n 0's or 1's. The formulas have the same form

$$\mu_{\hat{p}} = p$$

$$\sigma_{\hat{p}}^{2} = \frac{p(1-p)}{n}$$

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

Lets suppose that the life times of the light bulbs have a gamma distribution with  $\mu_x = 2$  years and  $\sigma_x = 1$  year.

Sample n bulbs and calculate sample average

$$\mu_{\bar{x}} = \mu_x = 2$$
  
$$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}} = \frac{1}{\sqrt{n}}$$

Lets see how the sampling distribution changes as n increases along the sequence n=2,10,50,100

We will examine this two ways:

- The exact sampling distribution (which happens to be also a gamma distribution with the appropriate mean and standard deviation).
- Monte Carlo experiment with 10,000 samples of  $\bar{X}$  for each n. The blue line is the exact sampling distribution



Lifetime



Lifetime





Lifetime

4

6





As the sample size n increases, the sampling distribution of  $\overline{X}$  approaches a normal distribution.



## **Central Limit Theorem**

Assume that  $X_1, X_2, \ldots, X_n$  are independent and identically distributed with mean  $\mu_x$  and standard deviation  $\sigma_x$ . Then for large sample sizes n, the distribution  $\bar{X}$  is approximately  $N(\mu_x, \frac{\sigma_x}{\sqrt{n}})$ 

Note that the normal approximation for the distribution of  $\hat{p}$  is just a special case of the Central Limit Theorem (CLT).

#### What is a large sample size?

As we saw with the binomial distribution, how well the normal approximation for  $\hat{p}$  depended on how p, which influenced the skewness of the population distribution.

This same idea holds for the general for the CLT.

If the observations are normal,  $\overline{X}$  has precisely a normal distribution for any n, since, as we've discussed before, sums of normals are normal. However the farther the density looks like a normal, the bigger that n needs to be for the approximation to do a good job.



Section 5.2 - The Sampling Distribution of a Sample Mean



Section 5.2 - The Sampling Distribution of a Sample Mean

The sampling distributions based on observations from the Gamma(5,1) distribution ( $\mu_x = 5, \sigma_x = \sqrt{5}$ ) look more normal than the sampling distributions based on observations from the Exponential distribution ( $\mu_x = 1, \sigma_x = 1$ ). For every sample size, the distribution of  $\bar{X}$  is more normal for the Gamma distribution than the Exponential distribution.



Section 5.2 - The Sampling Distribution of a Sample Mean

The Central Limit Theorem allows us to make approximate probability statements about  $\bar{X}$  using a normal distribution, even though X is not normally distributed.

So for the light bulb example with n = 50, the  $P[\bar{X} \le 1.9]$  can be approximated by the normal distribution.

$$P[\bar{X} \le 1.9] = P\left[\frac{\bar{X} - 2}{\frac{1}{\sqrt{50}}} \le \frac{1.9 - 2}{\frac{1}{\sqrt{50}}}\right]$$
$$= P[Z \le -0.707]$$
$$\approx 0.2398$$

The true probability is 0.2433.



The CLT can also be used to make statements about sums of independently and identically distributed random variables. Let

$$S = X_1 + X_2 + \ldots + X_n = n\bar{X}$$

Then

$$\mu_S = n\mu_{\bar{x}} = n\mu_x$$
  

$$\sigma_S^2 = n^2 \sigma_{\bar{x}}^2 = n^2 \frac{\sigma_x^2}{n} = \sigma_x^2 n$$
  

$$\sigma_S = \sigma_x \sqrt{n}$$

So S is approximately  $N(n\mu_x, \sigma_x\sqrt{n})$  distributed.

### Relaxing assumptions for the CLT

The assumptions for the CLT can be relaxed to allow for some dependency and some differences between distributions. This is why much data is approximately normally distributed. The more general versions of the theorem say that when an effect is the sum of a large number of roughly equally weighted terms, the effect should be approximately normally distributed.

For example, peoples heights are influenced a (potentially) large number of genes and by various environmental effects. Histograms of adult men's and women's heights are both well described by normal densities.

Another consequence of this, is that  $\overline{X}$  based on a simple random samples, with fairly large sampling fractions are also approximately normally distributed.

## Sampling With and Without Replacement

Simple Random Sampling is sometimes referred to sampling without replacement. Once a member of the population is sampled, it can't be sampled again. As discussed before, the without replacement action of the sampling introduces dependency into the observations.

Another possible sampling scheme is sampling with replacement. In this case, when a member of a population is sampled, it is returned to the population and could be sampled again. This occurs if your sampling scheme is similar to repeated rolling of a dice. There is no dependency between observations in this case, as at each step, the members population that could be sampled are the same. This situation is also equivalent to drawing from an "infinite" population.

When SRS is used, the variance of the sampling distribution needs to be adjusted for dependency induced by the sampling.

The correction is based on the Finite Population Correction (FPC)

$$f = \frac{N-n}{N}$$

which is the fraction of the population which is not sampled.

Then the variance and standard deviation of  $\bar{X}$  are

$$\sigma_{\bar{x}}^2 = \frac{\sigma_x^2}{n} f$$
$$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}} \sqrt{f}$$

So when a bigger fraction of the population is sampled (so f is smaller), you get a smaller spread in the sampling distribution.

However when n is small relative to N, this correction has little effect. For sample, if 10% of the population is sampled, so f = 0.9, the standard deviation of the sampling distribution is about 95% ( $\sqrt{0.9}$ ) of the standard deviation for that of with replacement sampling. If a 1% sample is taken, the correction on the standard deviation is 0.995.

Except when fairly large sampling fractions occur, the FPC is usually not used.