

Approximate Confidence Intervals for Means of Non-Normal Probability Distributions

When we have independent observations X_1, X_2, \dots, X_n with common mean $E[X]$ and standard deviation $\sqrt{\text{Var}(X)}$, the sample mean \bar{X}_n has expected value $E[\bar{X}_n] = E[X]$ and standard deviation $\sqrt{\text{Var}(X)/n}$.

We know from the CLT that

$$\frac{\sqrt{n}(\bar{X}_n - E[X])}{\sqrt{\text{Var}[X]}} \stackrel{\text{approx}}{\sim} N(0,1)$$

for 'large' n .

If the X 's are normal, the standardized mean has a standard normal probability distribution. And if the X 's have some other probability distribution, the standardized mean still has an approximate standard normal distribution when the sample is large.

Confidence Intervals II

This enables us to determine an approximate 95% confidence interval for the expected value, $E[X]$, when the sample size is 'large' regardless of the underlying distribution:

$$\bar{X}_n \pm 1.96 \sqrt{\text{Var}(X)/n}$$

Why? Because

$$P\left(-1.96 < \frac{\sqrt{n}(\bar{X}_n - E[X])}{\sqrt{\text{Var}[X]}} < 1.96\right) \approx 0.95$$

or

$$P\left(\bar{X}_n - 1.96\sqrt{\frac{\text{Var}[X]}{n}} < E[X] < \bar{X}_n + 1.96\sqrt{\frac{\text{Var}[X]}{n}}\right) \approx 0.95$$

This is a powerful and amazing result, because all we need to assume is that we know

- 1) $\text{Var}[X]$ and
- 2) that n is 'large' enough

But number 2 is empirically verifiable, and number 1 is not necessary if we use a t-distribution, replacing $\text{Var}[X]$ with S^2 !!

Confidence Intervals II

Before we go on, let's examine the CLT.

Example:

To take a familiar extreme case, suppose the X 's are i.i.d. Bernoulli(θ) random variables. The Central Limit Theorem ensures that for large n

$$\frac{\sqrt{n}(\bar{X}_n - \theta)}{\sqrt{\theta(1-\theta)}} \stackrel{approx}{\sim} N(0,1)$$

So,

$$P\left(\frac{\sqrt{n}(\bar{X}_n - \theta)}{\sqrt{\theta(1-\theta)}} > 1.645\right) \approx 0.05$$

$$P\left(\frac{\sqrt{n}(\bar{X}_n - \theta)}{\sqrt{\theta(1-\theta)}} < -1.96\right) \approx 0.025$$

and so forth.

Confidence Intervals II

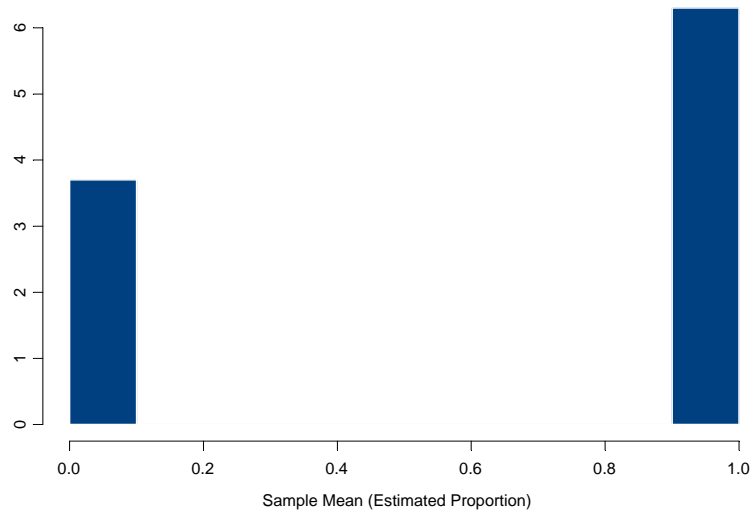
Now this is really quite remarkable: The Bernoulli probability distribution is not at all like the normal — instead of the normal's continuous, symmetric, bell-shaped density, the Bernoulli has all of its probability concentrated in two (usually unequal) blobs, on the points 0 and 1.

Yet, the standardized average of many independent Bernoulli random variables has approximately the same (standard normal) probability distribution as the standardized average of independent normal random variables.

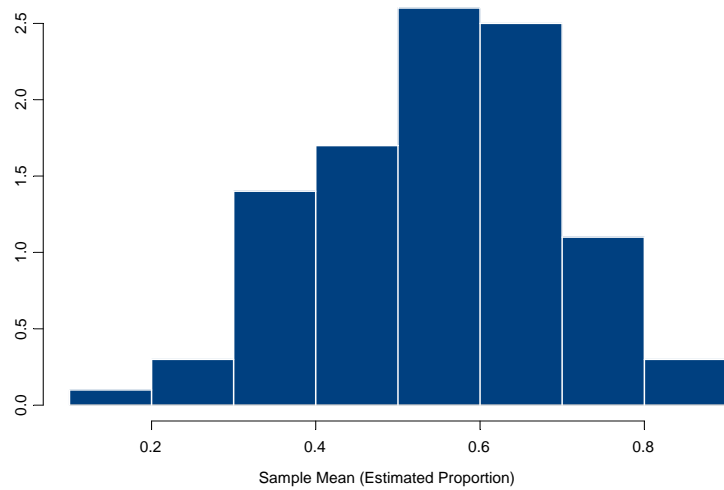
We've already seen this process in action. For the Bernoulli(0.6) distribution the histogram for \bar{X}_{25} already looks much like the normal density with the same mean and variance, and for $n = 100$, the similarity is even stronger.

Confidence Intervals II

Emperical Distribution of Sample Average of size $n=1$

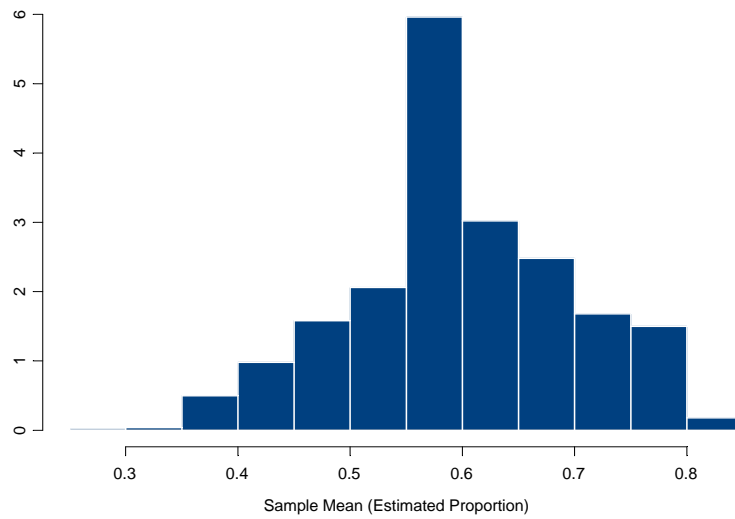


Emperical Distribution of Sample Average of size $n=10$

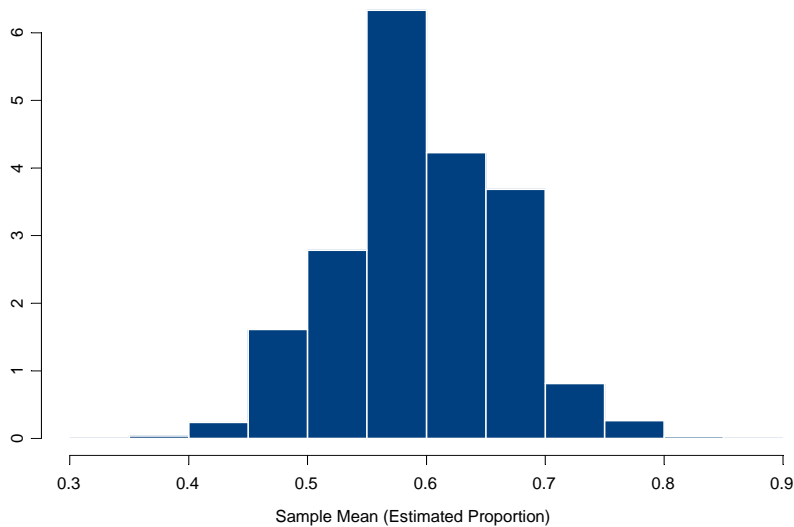


Confidence Intervals II

Emperical Distribution of Sample Average of size $n=25$

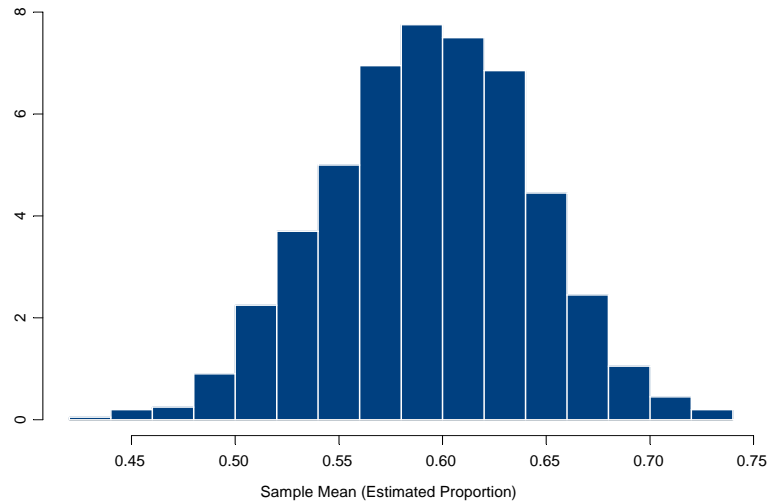


Emperical Distribution of Sample Average of size $n=50$

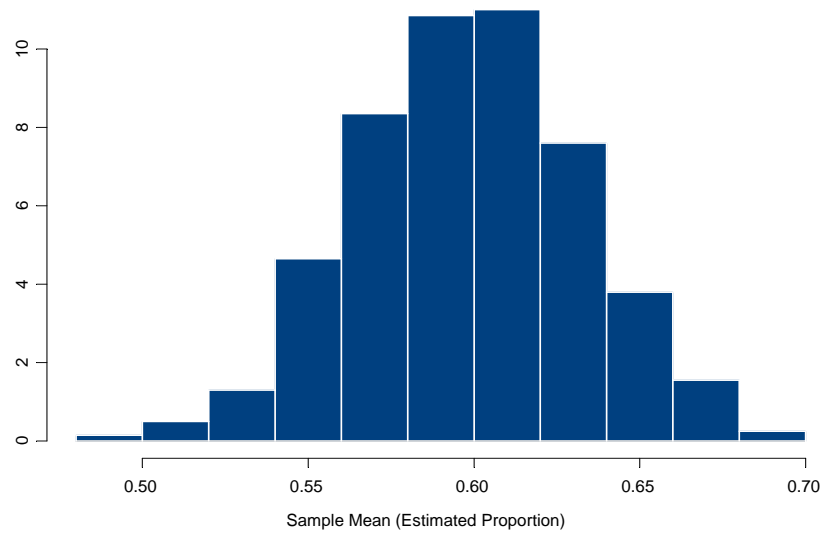


Confidence Intervals II

Empirical Distribution of Sample Average of size $n=100$



Empirical Distribution of Sample Average of size $n=200$



Confidence Intervals II

Let's compare some of the probabilities: For the Bernoulli(0.6),

$$\begin{aligned} 0.05 &\approx P\left(\frac{\sqrt{n}(\bar{X}_n - \theta)}{\sqrt{\theta(1-\theta)}} > 1.645\right) = P\left(\frac{\sqrt{n}(\bar{X}_n - 0.6)}{\sqrt{0.6(1-0.6)}} > 1.645\right) \\ &= P\left(\bar{X}_n > 0.6 + 1.645\sqrt{\frac{0.6 \times 0.4}{n}}\right) \\ &= P\left(\sum_{i=1}^n X_i > 0.6n + 1.645\sqrt{n}\sqrt{0.6 \times 0.4}\right) \end{aligned}$$

The last expression is the probability that a **Binomial(n, 0.6)** random variable (the sum of x's) will exceed the quantity on the right hand side.

For any value of n we can calculate this probability exactly and see how close it is to the probability that the standard normal will exceed 1.645 (which is 0.05). For n = 25, 100, and 1000, the exact probabilities are 0.029, 0.040, and 0.050.

Confidence Intervals II

The normal distribution plays a central role in statistics. This is not because many of the random variables that we observe have normal distributions. (Most do not.) It is because the normal can be used to approximate the distributions of averages (and of many other summary quantities as well) of samples from all type of distributions. The reason for this is the Central Limit Theorem.

What makes this really important is that it remains true when $\text{Var}(X)$ is replaced by the sample variance, S_n^2 :

$$\frac{\sqrt{n}(\bar{X}_n - E[X])}{S_n} \stackrel{\text{approx}}{\sim} N(0,1)$$

(for large n). This means that

$$\bar{X}_n \pm 1.96 SE(\bar{X}_n)$$

is an approximate 95% CI, not only when the individual X 's are normal and the standard error is given, but also when the X 's have some other probability distribution, and the standard error is estimated from the sample. (WOW!!)

Confidence Intervals II

Example:

In the Bernoulli(θ) case, the average, \bar{X}_n , is the proportion of successes, which is used to estimate the probability of success, θ . In this case, \bar{X}_n is often represented by the symbol $\hat{\theta}$ ("theta hat"), to indicate its role as an estimator of θ . The standardized mean is

$$\frac{\sqrt{n} (\bar{X}_n - \theta)}{\sqrt{\theta(1-\theta)}} = \frac{\sqrt{n} (\hat{\theta}_n - \theta)}{\sqrt{\theta(1-\theta)}}$$

If we estimate the variance, which is $\theta(1-\theta)$, by the sample variance S_n^2 the result

$$\frac{\sqrt{n}(\bar{X}_n - \theta)}{S_n} \stackrel{approx}{\sim} N(0,1)$$

Another natural estimate of the variance in this special case (Bernoulli) would be $\hat{\theta}(1-\hat{\theta})$ (i.e., $\bar{X}_n(1-\bar{X}_n)$).

$$\frac{\sqrt{n}(\bar{X}_n - \theta)}{\sqrt{\bar{X}_n(1-\bar{X}_n)}} = \frac{\sqrt{n}(\hat{\theta} - \theta)}{\sqrt{\hat{\theta}(1-\hat{\theta})}} \stackrel{approx}{\sim} N(0,1)$$

Confidence Intervals II

When n is large we can use either variance estimate -- they are essentially the same:

$$\begin{aligned} S_n^2 &= \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{n-1} = \frac{\sum_{i=1}^n X_i^2 - n(\bar{X}_n)^2}{n-1} \\ &= \frac{\sum_{i=1}^n X_i - n(\bar{X}_n)^2}{n-1} \\ &= \frac{n\bar{X}_n - n(\bar{X}_n)^2}{n-1} \\ &= \bar{X}_n \left(1 - \bar{X}_n\right) \left(\frac{n}{n-1}\right) \end{aligned}$$

Thus for a Bernoulli sample with large n , the sample variance S_n^2 is essentially just the same as

$$\bar{X}_n (1 - \bar{X}_n) = \hat{\theta}(1 - \hat{\theta})$$

because for large n the ratio $n/(n-1)$ is approximately 1.

Confidence Intervals II

What does all this tell us about the Student's t confidence interval?

Suppose we observe i.i.d. random variables X_1, X_2, \dots, X_n and want a 95% confidence interval for their expected value. If the distribution is normal, then the appropriate confidence interval is the Student's t interval:

$$\bar{X}_n \pm t_{n-1} S_n / \sqrt{n}$$

If the distribution of the X 's is normal, the coverage probability

$$P\left(\bar{X}_n - t_{n-1} \sqrt{\frac{S_n^2}{n}} < E[X] < \bar{X}_n + t_{n-1} \sqrt{\frac{S_n^2}{n}}\right) = 0.95$$

is exactly 0.95,

but even if the X 's have some other distribution the coverage probability is still approximately 0.95.

Confidence Intervals II

Why? The reasoning goes like this:

$$\begin{aligned} & P\left(-t_{n-1} < \frac{\sqrt{n}(\bar{X}_n - E[X])}{S_n} < +t_{n-1}\right) \\ & \approx P\left(-1.96 < \frac{\sqrt{n}(\bar{X}_n - E[X])}{S_n} < +1.96\right) \\ & \approx P\left(-1.96 < \frac{\sqrt{n}(\bar{X}_n - E[X])}{\sqrt{\text{Var}[X]}} < +1.96\right) \\ & = 0.95 \end{aligned}$$

This approximation gets better because

- 1) the t-distribution $\rightarrow N(0,1)$
- 2) $S_n^2 \rightarrow \text{Var}[X]$
- 3) $\frac{\sqrt{n}(\bar{X}_n - E[X])}{\sqrt{\text{Var}[X]}} \rightarrow N(0,1)$

as the sample size increases.

Confidence Intervals II

The Student's t confidence interval whose coverage probability is exactly 0.95 when the X 's are normal is still approximately 0.95 when they are not normal (as long as the sample is not too small).

This is terribly important in practice, because we never really know what distribution we're actually sampling from. Because they remain valid (covering the true mean with approximately the probability that they are supposed to) under quite general conditions, Student's t CI's are said to be "robust."

Statistics "works" in the real world because many of its procedures, like the Student's t CI, have this property of "robustness," or approximate validity under a wide range of probability models, i.e. under a wide range of assumptions about the probability distribution.