
The Central Limit Theorem

5.1 The Fundamental Theorem of Statistics

Statistics is the science that enables us to infer facts about large populations from limited samples of data. For instance, suppose we want to determine the average annual income of American households. One approach is to conduct a comprehensive survey of all American households. This approach would yield a very precise answer, but the cost would be huge. Another approach, less accurate, but more practical, is to conduct a limited survey of a relatively small “random sample” of households, say 500. Suppose, for the sake of argument, that our survey of 500 households yields an average annual income of \$45,000. What does this tell us about the average annual income of all American households? Based on this limited sample, to what extent is it valid to say that the average American household has an annual income of \$45,000? Statistics can answer this question precisely. Statistics can even determine the exact degree to which we can trust our findings.

The Central Limit Theorem is the principle that enables us to determine the degree of confidence to which we can trust the results of a limited survey. Roughly speaking, the Central Limit Theorem asserts that “sample means” are distributed “nearly normally” with a certain mean and standard deviation closely related to the mean and standard deviation of the whole population. The larger the samples, the closer is the distribution of their means to a normal distribution, and the larger an individual sample, the more likely it is that its mean accurately reflects the mean of the population. The Central Limit Theorem is important enough to be called the “Fundamental Theorem of Statistics”, although, reader be cautioned, this designation is not widely recognized among statisticians.

5.2 Elements of Statistics

We begin by recalling some basic definitions. By a *population*, we mean a set of items, cases, or trials, each possessing a property that can be observed, evaluated and recorded. Typically, the property to be evaluated and recorded is represented by a numerical variable, say x . For a finite population having N members, the *population mean* $\mu(x)$ and the *population standard deviation* $\sigma(x)$ are defined by the formulas:

$$\mu(x) = \frac{1}{N} \sum x \qquad \sigma^2(x) = \frac{1}{N-1} \sum (x - \mu)^2$$

These sums are understood to range over one measurement x for each member of the population.

By a *sample*, we mean a subset of a population. A sample is said to be of *size* n if it has n members. The *sample mean* \bar{x} and the *sample standard deviation* $s(x)$ are defined by the formulas:

$$\bar{x} = \frac{1}{n} \sum x \qquad s^2(x) = \frac{1}{n-1} \sum (x - \bar{x})^2$$

Ex: Suppose the population has 11 members, and the values of x are given by the data list:

$$x = 0, 1, 3, 3, 4, 6, 6, 6, 7, 9, 10$$

By direct calculation, we get $\sum x = 55$, so that $\mu(x) = 55/11 = 5$. To compute the standard deviation, it is convenient to set up a table like the one below:

x	$x - \mu$	$(x - \mu)^2$
0	-5	25
1	-4	16
3	-2	4
3	-2	4
4	-1	1
6	1	1
6	1	1
6	1	1
7	2	4
9	4	16
10	5	25

Therefore $\sum(x - \mu)^2 = 98$, so that $\sigma^2 = 98/10 = 9.8$ and $\sigma = \sqrt{9.8} \approx 3.1305$.

By a *statistical experiment* we mean a repeatable act of measurement whose outcome is, to some extent, unpredictable. For instance, when a six-faced die is rolled, the value (1 through 6) of the up face is unpredictable. Thus the act of rolling a six-faced die qualifies as a statistical experiment. The concept of “statistical experiment” applies to many different kinds of situations. For instance, the process of selecting a member at random from a large population and performing a measurement of a certain variable x qualifies as a statistical experiment. Similarly, a sample of size n can be thought of as a sequence of n repetitions of the experiment of selecting a random member from the population.

In general, the possible numerical outcomes of a statistical experiment are represented by a *random variable* x . A random variable is said to be *continuous* if it can assume any value in an interval of real numbers, with no interruptions. Analogously, a random variable is said to be *discrete* if it can assume only certain consecutive values separated by “forbidden intervals”. For instance, the random variable representing the weight of a snowflake is continuous, as two snowflakes can differ in weight by an arbitrarily small amount. On the other hand, the random variable representing the number of credit cards held by a random individual is discrete since its values are restricted to the discrete set $\{0, 1, 2, 3, \dots\}$.

The *frequency distribution* of a random variable x is found by counting the number of times each of its possible values occurs in a sequence of repetitions of the underlying experiment. Typically, some values occur more frequently than others. For instance, if a pair of dice is rolled and the sum x of the two up faces is recorded, then the possible values of x are $\{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$. However, not all of these values are equally likely. For the sake of illustration, suppose the double dice rolling experiment has been repeated 360 times. The resulting frequency distribution of x (defined as the sum of the up faces) might look something like the following:

x	2	3	4	5	6	7	8	9	10	11	12
Freq.	11	19	32	37	53	58	49	42	30	21	8

The *relative frequency distribution* of a random variable is obtained by dividing the observed frequencies by the total number of data points. For the double dice rolling experiment, repeated 360 times, we obtain the following relative frequency distribution:

x	2	3	4	5	6	7	8	9	10	11	12
Rel. Freq.	.031	.053	.089	.103	.147	.161	.136	.117	.083	.058	.022

The relative frequencies displayed above are calculated on the basis of repeated experimentation. They represent approximations to the true probabilities. The more often the experiment is repeated, the more accurately the relative frequencies reflect the true probabilities. In fact, by a simple theoretical analysis, it can be shown that the true probabilities $p(x)$ for the double dice rolling experiment are as follows:

x	2	3	4	5	6	7	8	9	10	11	12
$p(x)$	1/36	1/18	1/12	1/9	5/36	1/6	5/36	1/9	1/12	1/18	1/36

When statisticians talk about the *probability distribution* of a random variable, it is tacitly assumed that the underlying experiment can be repeated as often as desired, yielding, in the limit, relative frequencies that are equal to the true probabilities.

For a discrete random variable x , the mean $\mu(x)$ (also called the *expected value*) and the standard deviation $\sigma(x)$ are defined by the formulas:

$$\mu(x) = \sum p(x)x$$

$$\sigma^2(x) = \sum p(x)(x - \mu)^2 = \sum p(x)x^2 - \mu^2,$$

where these sums are understood to range over all possible values of x .

If the probabilities $p(x)$ are approximations based on a limited number of repetitions of the underlying experiment, then the symbols \bar{x} and $s(x)$ are used to denote the mean and standard deviation, respectively. The actual formulas for \bar{x} and $s(x)$ are the same as the formulas for $\mu(x)$ and $\sigma(x)$, with \bar{x} instead of μ in the formula for $s(x)$.

Ex: For the double dice rolling experiment, using the exact probabilities, we get

$$\mu(x) = \sum p(x)x = 7.000$$

$$\sigma^2(x) = \sum p(x)(x - \mu)^2 = 5.833$$

$$\sigma(x) = 2.415$$

Using the experimental relative frequencies based on 360 repetitions, we get

$$\bar{x} = \sum p(x)x = 6.974$$

$$s^2(x) = \sum p(x)(x - \bar{x})^2 = 5.8113$$

$$s(x) = 2.411$$

5.3 Probability Density Functions

For a discrete random variable x , the relative frequency or probability of each of its possible values can be calculated by repeated experimentation, or, occasionally, by theoretical reasoning. As we saw in the previous section, the distribution of x is completely specified by listing the probability of each of its possible values. However, for a continuous random variable, because the possible values for x are so numerous and so closely packed, the relative frequency or probability of any specific value of x is indistinguishable from 0. For instance, if x represents the height of a typical American adult male, then the likelihood of finding a subject whose height is exactly 69.500000000 inches is exactly 0. However, if someone asks the question “what is the probability of finding a subject whose height in inches lies somewhere in the interval $69.4 \leq x \leq 69.6$?”, then the answer can be found precisely by calculating the percent of subjects whose heights lie in this range. The answer is a specific positive real number.

In general, the probability distribution of a continuous random variable x is represented by a *probability density function* $f(x)$ having the following properties:

- 1) $f(x) \geq 0$ for all $x \in \mathbb{R}$.
- 2) The area below the graph of $y = f(x)$ and above the x -axis is exactly equal to 1.
- 3) The probability of x lying in the interval $a \leq x \leq b$ is equal to the area of the region under the graph of $y = f(x)$, above the x -axis, and between $x = a$ and $x = b$.

Properties 2) and 3) can best be expressed by the formulas

$$\text{a) } \int_{-\infty}^{\infty} f(x) dx = 1$$

$$\text{b) } p(a \leq x \leq b) = \int_a^b f(x) dx$$

In terms of the probability density function $f(x)$, the mean μ and standard deviation σ of the continuous random variable x are defined by the formulas

$$\mu = \int_{-\infty}^{\infty} xf(x) dx \qquad \sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

Ex: Consider the probability density function defined by:

$$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } 0 \leq x \leq 1 \\ 2 - x & \text{if } 1 \leq x \leq 2 \\ 0 & \text{if } x > 2 \end{cases}$$

The reader can easily verify the following calculations at his/her leisure:

$$\mu = \int_0^2 xf(x) dx = \int_0^1 x^2 dx + \int_1^2 x(2 - x) dx = 1$$

$$\sigma^2 = \int_0^2 (x - \mu)^2 f(x) dx = \int_0^1 (x - 1)^2 x dx + \int_1^2 (x - 1)^2 (2 - x) dx = \frac{1}{6}$$

$$p(.5 \leq x \leq 1.75) = \int_{.5}^{1.75} f(x) dx = \int_{.5}^1 x dx + \int_1^{1.75} (2 - x) dx = .84375$$

5.4 Normal Distributions

The most natural and important distributions in statistics are the *normal distributions*. A normal distribution is completely determined by its mean and standard deviation. If a normal distribution has mean μ and standard deviation σ , then its probability density function is

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

The graph of this function has a characteristic shape often referred to as a “bell curve”.

μ

Examples of normal distributions are found everywhere in nature. For instance, the heights of adult American males are normally distributed, as are their weights. Most mental test scores such as SAT scores and IQ scores are normally distributed. In fact, just about any random variable whose value is determined by multiple random influences is normally distributed.

As the reader may recall from elementary courses in statistics, a normal distribution of mean μ and standard deviation σ always satisfies the so-called “empirical rule”:

- a) $p(\mu - \sigma \leq x \leq \mu + \sigma) \approx .68$
- b) $p(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx .95$
- c) $p(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx .997$

In the interest of notational brevity, if a random variable x is known to be normally distributed with mean μ and standard deviation σ , then we may want to express this fact symbolically as

$$\text{dist}(x) = \mathcal{N}(\mu, \sigma)$$

Ex: Suppose $\text{dist}(x) = \mathcal{N}(20, 3)$. Using the empirical rule and the symmetry of the bell curve, we can derive the following selected facts, which the reader is advised to verify

$$p(17 \leq x \leq 23) \approx .68$$

$$p(14 \leq x \leq 26) \approx .95$$

$$p(x \geq 26) \approx .025$$

$$p(x \geq 29) \approx .0015$$

$$p(17 \leq x \leq 26) \approx .815$$

Among the class of normal distributions, there is a very special one whose values can be found tabulated in every statistics textbook, namely the *standard normal distribution* $\mathcal{N}(0, 1)$. It turns out that the probability values associated with any normal distribution $\mathcal{N}(\mu, \sigma)$ can be computed in terms of the probability values of $\mathcal{N}(0, 1)$. The principle that enables us to do so is the so-called *Standardization Principle*.

Standardization Principle: If $\text{dist}(x) = \mathcal{N}(\mu, \sigma)$ and if $z = \frac{x - \mu}{\sigma}$, then $\text{dist}(z) = \mathcal{N}(0, 1)$.

In particular, $p(a \leq x \leq b) = p\left(\frac{a - \mu}{\sigma} \leq z \leq \frac{b - \mu}{\sigma}\right)$.

Ex: Suppose $\text{dist}(x) = \mathcal{N}(20, 3)$ and suppose we want to evaluate $p(17.9 \leq x \leq 23.6)$. We begin by calculating the standardized z -scores of the endpoints of our interval, namely $(17.9 - 20)/3 = -.7$ and $(23.6 - 20)/3 = 1.2$. So $p(17.9 \leq x \leq 23.6) = p(-.7 \leq z \leq 1.2)$. Now we turn to the standard normal table. Typically such tables display only the probabilities for intervals of the form $0 \leq z \leq c$. But we can exploit the symmetry of the bell curve to calculate the probability of any interval. In this case, we can see proceed as follows:

$$\begin{aligned} p(-.7 \leq z \leq 1.2) &= p(-.7 \leq z \leq 0) + p(0 \leq z \leq 1.2) \\ &= p(0 \leq z \leq .7) + p(0 \leq z \leq 1.2) \\ &= .2580 + .3849 = .6429 \end{aligned}$$

Therefore $p(17.9 \leq x \leq 23.6) = .6429$.

5.5 The Central Limit Theorem and Applications

Let x denote a random variable associated with a large population. Assume that the probability distribution of x , not necessarily normal, has mean μ and standard deviation σ . Consider the set of all samples of size n selected from the population, and let \bar{x} denote the mean of a typical such sample. The Central Limit Theorem gives information about the distribution of the sample means \bar{x} .

The Central Limit Theorem: Under the above assumptions, let $\mu(\bar{x})$ denote the mean of all the means of samples of size n , and let $\sigma(\bar{x})$ denote the standard deviation of the sample means. Then we have $\mu(\bar{x}) = \mu$, $\sigma(\bar{x}) = \sigma/\sqrt{n}$, and

$$\text{dist}(\bar{x}) \approx \mathcal{N}(\mu, \sigma/\sqrt{n})$$

The approximation to a normal distribution becomes ever more accurate as $n \rightarrow \infty$. ■

Ex: The National Fraternity of $\text{MY}\Delta$ has 20,000 members whose college GPA's are distributed with mean 2.8 and standard deviation 0.75. The fraternity is organized into 400 local chapters with 50 members in each chapter. What percent of the chapters have average GPA's of at least 3.00?

Solution: Let x denote GPA. By the CLT, the average GPA's \bar{x} of the chapters are distributed like the normal distribution $\mathcal{N}(2.8, .75/\sqrt{50}) = \mathcal{N}(2.8, 0.106)$. By the Standardization Principle, $p(\bar{x} \geq 3.00) = p(z \geq (3.00 - 2.80)/0.106) = p(z \geq 1.887) = 3\%$.

The most powerful applications of the Central Limit Theorem are to questions having to do with determining confidence intervals for estimates of population parameters when these estimates are based on limited samples. For instance, suppose a random survey of 400 households is conducted in the hope of determining, to a reasonable degree of accuracy, the average household income of all American households. Intuitively it seems clear that the larger the sample, the more accurate the result is likely to be. But who's to say that a sample of size 400 is large enough to yield even a crude estimate? Imagine, for the sake of illustration, that our survey produces a sample mean of \$45,000. How likely is it that the true mean income μ of all American households lies in $\$40,000 \leq \mu \leq \$50,000$? The Central Limit Theorem can provide answers to these and similar questions.

Let x denote a random variable associated with a large population. Our aim is to estimate the population mean μ by measuring the mean \bar{x} of a sample of size n . If we assume that the sample standard deviation s is a reasonably accurate approximation to the population standard deviation σ , then, by the Central Limit Theorem, we have

$$\text{dist}(\bar{x}) \approx \mathcal{N}(\mu, s/\sqrt{n}).$$

Therefore, by the Standardization Principle, we get

$$\text{dist}\left(\frac{\bar{x} - \mu}{s/\sqrt{n}}\right) \approx \mathcal{N}(0, 1).$$

Now given a real number α , $0 < \alpha < 1$, let $z_{\alpha/2}$ denote the value of z in the standard normal distribution satisfying the relation

$$p(0 \leq z \leq z_{\alpha/2}) = 0.5 - \alpha/2.$$

In other words, if $\text{dist}(z) = \mathcal{N}(0, 1)$, then

$$p(-z_{\alpha/2} \leq z \leq z_{\alpha/2}) = 1 - \alpha.$$

Thus we have

$$p\left(-z_{\alpha/2} \leq \frac{\bar{x} - \mu}{s/\sqrt{n}} \leq z_{\alpha/2}\right) = 1 - \alpha,$$

which is equivalent to the relation

$$p\left(\bar{x} - z_{\alpha/2} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{x} + z_{\alpha/2} \frac{s}{\sqrt{n}}\right) = 1 - \alpha.$$

In the above relation, the number $(1 - \alpha)100\%$ is called the *degree of confidence*. The number $E = z_{\alpha/2} \cdot s/\sqrt{n}$ is called the *margin of error*. Thus, according to the above equation, we can be confident at the $(1 - \alpha)100\%$ level that the true value of the population mean lies in the interval $\bar{x} - E \leq \mu \leq \bar{x} + E$. Conversely, if we insist on a $(1 - \alpha)100\%$ level of confidence for μ to lie in the interval $\bar{x} - E \leq \mu \leq \bar{x} + E$, then we must take $n \geq (z_{\alpha/2} \cdot s/E)^2$.

Ex: A random variable x is distributed with unknown mean μ and unknown standard deviation. However, a sample of size 100 yields $\bar{x} = 500$ and $s = 80$.

- With what degree of confidence can we assert that $490 \leq \mu \leq 510$.
- With what degree of confidence can we assert that $480 \leq \mu \leq 520$.
- Determine a 95% confidence interval for μ .
- How large a sample would we need if we wanted to estimate μ with 99% confidence and a margin of error not exceeding 5 units?

Solution:

- Here $n = 100$, $s = 80$, and $E = 10$. Thus $z_{\alpha/2} = E\sqrt{n}/s = 1.25$. From the tables, we get $\alpha/2 = p(z \geq 1.25) = 0.1056$. Thus $\alpha \approx 0.21$, so that $1 - \alpha \approx .79$, which means that the degree of confidence is approximately 79%.
- Here $n = 100$, $s = 80$, and $E = 20$. Thus $z_{\alpha/2} = E\sqrt{n}/s = 2.50$. From the tables, we get $\alpha/2 = p(z \geq 2.50) = 0.0062$. Thus $\alpha \approx 0.0124$, so that $1 - \alpha \approx .9876$, which means that the degree of confidence is better than 98%.
- Here $n = 100$, $s = 80$, and $\alpha = .05$. Thus $\alpha/2 = .025$. Using the tables, with $.025 = p(z \geq z_{\alpha/2})$, we get $z_{\alpha/2} = 1.96$. Therefore $E = z_{\alpha/2} \cdot s/\sqrt{n} = 15.68$. So the corresponding 95% confidence interval is $[484.32, 515.68]$.
- Here $s = 80$, $E = 5$, $\alpha = .01$, and $\alpha/2 = .005$. From the tables, we get $z_{\alpha/2} = 2.575$. Therefore $n \geq (z_{\alpha/2} \cdot s/E)^2 \approx 1700$. A sample of size at least 1700 would be needed.

5.5 Exercises

1. A discrete random variable x with values in the range $\{6, 7, 8, 9, 10\}$ is distributed with the following (incomplete) probability profile:

$$p(x = 6) = 0.10 \quad p(x = 7) = 0.25 \quad p(x = 8) = 0.15 \quad p(x = 9) = 0.20$$

- Calculate $p(x = 10)$
- Calculate $p(x \geq 8)$
- Calculate $p(6 \leq x \leq 8)$
- Calculate the mean $\mu(x)$
- Calculate the standard deviation $\sigma(x)$

2. A continuous random variable x is distributed with probability density function

$$f(x) = \frac{1}{\pi} \frac{1}{x^2 + 1}$$

- Sketch the graph of $y = f(x)$.
- Evaluate approximately $p(0 \leq x \leq 1)$.
- Evaluate approximately $p(-2 \leq x \leq 2)$
- Evaluate approximately $p(x \geq .5)$

3. A continuous random variable x is distributed with probability density function

$$f(x) = \begin{cases} 0 & \text{if } x < -1 \\ 0.5 + 0.5x & \text{if } -1 \leq x \leq 0 \\ 0.5 & \text{if } 0 \leq x \leq 1 \\ 1.0 - 0.5x & \text{if } 1 \leq x \leq 2 \\ 0 & \text{if } x > 2 \end{cases}$$

- (a) Sketch the graph of $y = f(x)$.
 - (b) Find $\mu(x)$.
 - (c) Find $\sigma(x)$.
 - (d) Evaluate exactly $p(0 \leq x \leq 1.5)$.
 - (e) Evaluate exactly $p(-1 \leq x \leq 1)$.
 - (f) Evaluate exactly $p(x \geq .5)$.
4. Suppose the random variable x possesses a normal distribution with mean $\mu = 100$ and standard deviation $\sigma = 15$. That is, $\text{dist}(x) = \mathcal{N}(100, 15)$.
- (a) Find $p(x \geq 120)$.
 - (b) Find $p(90 \leq x \leq 110)$.
 - (c) Find $p(x \leq 70)$.
 - (d) Find $p(x \geq 150)$.
 - (e) Find $p(95 \leq x \leq 105)$.
5. Medium sized watermelons weigh an average (mean) of 5 lb each with a standard deviation of 1.2 lb. The watermelons are packed for shipment in crates of 36. Thus, on average, the typical crate weighs 180 lb. What percent of the crates weigh between 174 and 186 lb?
6. The tail-lengths of the inhabitants of the distant planet of Mephiston are normally distributed with unknown mean μ and unknown standard deviation. However, a random sample of 100 Mephistonians shows an average tail-length of 150 cm with standard deviation of 25 cm.
- (a) Based on this sample, determine a 90% confidence interval for μ .
 - (b) How large of a sample would we need if we wanted to estimate μ with 95% confidence and a margin of error not exceeding 2.5 cm?
7. The heights of freshly cut Douglas fir trees of the Grade A variety, supplied by Farmer Jim's Nursery, are distributed with mean $\mu = 7.5$ ft and standard deviation $\sigma = .5$ ft. The trees are shipped to vendors in lots of size 36. Let \bar{x} denote the average height of the trees in a typical lot.
- (a) Calculate $p(7.4 \leq \bar{x} \leq 7.6)$.
 - (b) Determine a 99% confidence interval for \bar{x} .

-
-
8. The heights of medium-sized Frazier Fir Christmas trees are normally distributed with unknown mean μ and unknown standard deviation σ . However, a random sample of 100 such trees shows an average height of 84 inches with standard deviation of 6 inches.
- (a) Based on this sample, determine a 90% confidence interval for μ .
 - (b) How large of a sample would we need if we wanted to estimate μ with 95% confidence and a margin of error not exceeding 2 inches?
9. The weights of Supreme Deluxe Garden Fresh tomatoes are distributed with mean $\mu = 18$ oz and standard deviation 4.2 oz. The tomatoes are shipped to supermarkets in packages of three dozen. Let \bar{x} denote the average weight of the tomatoes in a typical package. Calculate $p(17 \leq \bar{x} \leq 19)$.
10. The wing-spans of sugar plum fairies are distributed with unknown mean μ . However, a specific random sample of 64 sugar plum fairies shows an average wing-span of 3.0 feet with standard deviation 0.50 feet.
- (a) Based on this sample, determine a 90% confidence interval for μ .
 - (b) Based on this sample, with what level of confidence could we infer that $2.80 \leq \mu \leq 3.20$?