# Sampling Distributions
## Module 7

Statistics 251: Statistical Methods

Updated 2021

## Three Types of Distributions

**data distribution**
the distribution of a variable in a *sample*

**population distribution**
the *probability distribution* of a *single observation* of a variable

**sampling distribution**
the *probability distribution* of a *statistic*

## Terms I

**sampling distribution**: a probability distribution of a statistic; it is a distribution of *all possible samples* (random samples) from a population and how often each outcome occurs in repeated sampling (of the same size $n$). Given simple random samples of size $n$ from a given population with a measured characteristic such as mean $\overline{X}$, proportion $(\hat{\pi})$[1], or standard deviation ($s$) for each sample, the probability distribution of all the measured characteristics is called a sampling distribution. It is the distribution of all possible samples (outcomes) of that statistic.
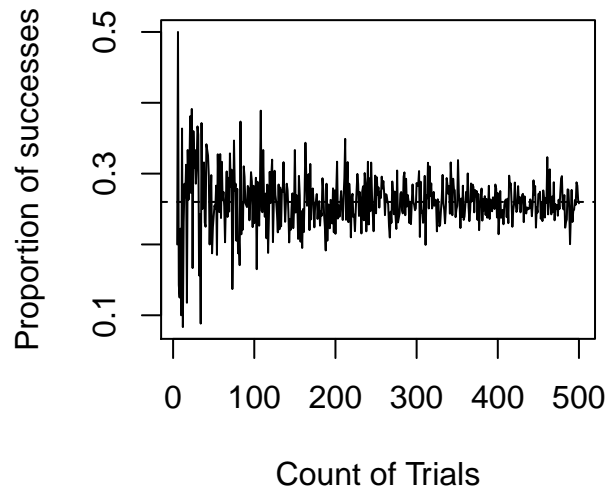Use of a statistic to estimate the parameter is the main function of inferential statistics as it provides the properties of the statistic.

## Terms II

**law of large numbers** states that as the number of repetitions of an experiment is increased, the relative frequency obtained in the experiment tends to become ever closer to the theoretical probability. Even though the outcomes do not happen according to any set pattern or order (overall), the long-term observed relative frequency will approach the theoretical probability

---

[1] $\pi$ and $\hat{\pi}$ are NOT $3.14159\ldots$, they are being used like $\mu$ and the other Greek letters we are using for notation.

**Simulation of LLN**



Count of Trials

## Central Limit Theorem (CLT)

*Definition*

The sampling distribution of the sample mean is approximately normal with mean $\mu_X$ and standard deviation (of the sampling distribution of the sample mean) $se = \frac{\sigma_X}{\sqrt{n}}$, provided $n$ is sufficiently large.

## Sampling distribution of the Sample Mean

If we take $n$ observations of a quantitative variable and then compute the mean $(\bar{x})$ of those observations in the sample, then $\bar{x}$ is the sample mean statistic.

Assumptions: Each observation $x$ has the same probability distribution with mean $\mu$ and standard deviation $\sigma$, and the observations are *independent*.

## Properties of the Sampling Distribution of $\bar{x}$

(1) The *mean* of the sampling distribution is $\mu$

(2) The *standard deviation* of the sampling distribution is $se = \frac{\sigma}{\sqrt{n}}$

(3) The *shape* of the sampling distribution becomes more like a normal distribution as $n$ increases

## Sampling distribution of the Sample Mean

$$\overline{X} \sim N\left(\mu, se_{mean}\right)$$

$$\text{Standard error of the mean: } \sigma_{\overline{X}} = se_{mean} = \frac{\sigma}{\sqrt{n}}$$

$$z = \frac{\overline{X} - \mu}{se_{mean}}$$

Sample sizes should be $n \geq 30$ for the sample mean If a distribution is already inherently normal, the sample size stipulation can be ignored.

## Sampling distribution of the Sample Proportion ($\hat{\pi}$)

If we make $n$ observations, and count the number of observations on which an outcome happens (call this $x$), then $\hat{\pi} = \frac{x}{n}$ is the *sample proportion* statistic.

Assumptions: $x$ has a binomial distribution where $n$ is the number of trials and the probability of the outcome on each trial is $\pi$.

## Properties of the Sampling Distribution of $\hat{\pi}$

(1) The *mean* of the sampling distribution is $\pi$.

(2) The *standard deviation* of the sampling distribution is $\sqrt{\pi(1-\pi)/n}$.

(3) The *shape* of the sampling distribution becomes more like a normal distribution as $n$ increases.

## Sampling distribution of $\hat{\pi}$

$$\hat{\pi} \sim N\left(\pi, se_{\hat{\pi}}\right)$$

$$\text{Standard error of the proportion: } \sigma_{\hat{\pi}} = se_{\hat{\pi}} = \sqrt{\frac{\pi(1-\pi)}{n}}$$

$$z = \frac{\hat{\pi} - \pi}{se_{\hat{\pi}}}$$

Sample sizes should be $n \geq 60$ for the sample proportion

## Properties of the Mean Total

The mean total is the average value of a distribution multiplied by the total number of trials. If we take $n$ observations of a quantitative variable and then compute the mean total (sum) ($\hat{\tau} = n\bar{x}$) of those observations in the sample, then $\hat{\tau}$ is the sample total statistic.

Assumptions: Each observation $x$ has the same probability distribution with mean $\tau = n\mu$ and standard deviation $\sqrt{n}\sigma$ (or maybe easier to see it as $\sigma(\sqrt{n})$), and the observations are *independent*.

## Properties of the Sampling Distribution of $\hat{\tau}$

(1) The *mean total* of the sampling distribution is $\tau = n\mu$

(2) The *standard deviation (of the total)* of the sampling distribution is $se = \sqrt{n}\sigma$

(3) The *shape* of the sampling distribution becomes more like a normal distribution as $n$ increases

## Sampling distribution of the Sample Total (Sum)

$$\hat{\tau} = n\overline{X} \quad \tau = n\mu \quad se_{sum} = \sqrt{n}\sigma$$

$$\hat{\tau} \sim N(\tau, se_{sum}) \; with \; se_{sum} = \sqrt{n}\sigma$$

$$z = \frac{n\overline{X} - n\mu}{se_{sum}} = \frac{\hat{\tau} - \tau}{se_{sum}}$$

## Simulation example

The linked file shows how taking multiple random samples of the same size from the same population will produce a normal distribution of the sample means. The examples show a normal distribution, exponential distribution, and a binomial distribution.

CLT simulation

## CLT for sample mean ($\overline{X}$) and sample sum/total ($\hat{\tau}$)

for sample mean ($\overline{X}$) and total ($\hat{\tau}$)

The level of a particular pollutant, nitrogen dioxide ($NO_2$), in the exhaust of a hypothetical model of car, that when driven in city traffic, has a mean level of 2.1 grams per mile ($g/m$) and a standard deviation of 0.3 $g/m$. Suppose a company has a fleet of 35 of these cars.

(a) What is the mean and standard deviation of the sampling distribution of the sample mean?

mean: $\mu_X = \mu = 2.1$ and $se_{mean} = \frac{\sigma}{\sqrt{n}} = \frac{0.3}{\sqrt{35}} = 0.0507$

$\overline{X} \sim N(\mu, se_{mean}) = \overline{X} \sim N(2.1, 0.0507)$

## CLT for $\overline{X}$ and $\hat{\tau}$ solutions

(b) find the probability that the mean $NO_2$ level is less than 2.03 $g/m$

$$P(\overline{X} < 2.03) = P\left(Z < \frac{2.03 - 2.1}{0.0507}\right) = P(Z < -1.38) = 0.083793$$

(c) Mandates by the EPA state that the average of the fleet of these cars cannot exceed 2.2 $g/m$, find the probability that the fleet $NO_2$ levels from their fleet exceed the EPA mandate

$$P(\overline{X} > 2.2) = 1 - P\left(Z < \frac{2.2 - 2.1}{0.0507}\right)$$

$$= 1 - P(Z < 1.97) = 1 - 0.975581 = 0.024419$$

## CLT for $\overline{X}$ and $\hat{\tau}$ solutions

(d) At most, 25% of these cars exceed what *mean $NO_2$* value?

Find the $z$ score that represents the top 25%, which is the same as the bottom 75% (is also $Q3$) and what is needed to find $z_{0.75} = 0.67449$. Next use $z = \frac{\overline{X} - \mu}{se_{mean}}$ and solve for $\overline{X}$: $\overline{X} = z(se_{mean}) + \mu$

$$\overline{X} = (0.67449)(0.0507) + 2.1 = 2.134197$$

## CLT for $\overline{X}$ and $\hat{\tau}$ solutions

(e) what is the mean and standard deviation of the total amount (sum), in $g/m$, of $NO_2$ in the exhaust for the fleet?

$$\tau = n\mu = 35(2.1) = 73.5$$

$$se_{sum} = \sqrt{n}\sigma = \sqrt{35}(0.3) = 1.774824$$

$$\hat{\tau} \sim N(\tau, se_{sum}) = \hat{\tau} \sim N(73.5, 1.7748)$$

## CLT for $\overline{X}$ and $\hat{\tau}$ solutions

(f) find the probability that the total amount of $NO_2$ for the fleet is between 70 and 75 $g/m$

$$P(70 < \hat{\tau} < 75) = P\left(\frac{70 - 73.5}{1.7448} < Z < \frac{75 - 73.5}{1.7748}\right)$$

$$= P(-2.01 < Z < 0.86) = P(Z < 0.86) - P(Z < -2.01)$$

$$= 0.805105 - 0.022216 = 0.78289$$

## CLT for proportion ($\hat{\pi}$)

Mars company claims that 10% of the M&M's it produces are green. Suppose that candies are packaged at random in bags that contain 60 candies.
(a) Describe the sampling distribution of the sample proportion (what should the distribution look like?); calculate the mean proportion and standard deviation of the sampling distribution of the sample proportion of green M&M's in bags that contain 60 candies (calculate $\pi$ and $se$).
(b) What is the probability that a bag of 60 candies will have more than 13% green M&M's?

## CLT for $\hat{\pi}$ solutions

(a) Describe the sampling distribution of the sample proportion; calculate the mean proportion and standard deviation of the sampling distribution of the sample proportion of green M&M's in bags that contain 60 candies.

The distribution of the sample proportion will be approximately normal since $n \geq 60$. The mean proportion $\pi = 0.1$ and the standard error is $\sqrt{\frac{\pi(1-\pi)}{n}} = \sqrt{\frac{(0.1)(1-0.1)}{60}} = 0.0387$ (the standard deviation of the sampling distribution of the sample proportion). Thus

$$\hat{\pi} \sim N(0.1, 0.0387)$$

## CLT for $\hat{\pi}$ solutions

(b) What is the probability that a bag of 60 candies will have more than 13% green M&M's?

$$P(\hat{\pi} > 0.13) = P\left(Z > \frac{0.13 - 0.1}{0.0387}\right)$$

$$= P(Z > 0.78) = 1 - P(Z < 0.78) = 1 - 0.782305$$

$$= 0.2177$$