## Confidence Intervals for a Population Mean

Let $X_1, \ldots, X_n$ be a simple random sample from a population with mean $\mu$ and variance $\sigma^2$. Let $\bar{X}$ be the sample mean, and $S_n$ be the sum of sample observation. If $n$ **is sufficiently large**,

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

and

$$S_n \sim N\left(n\mu, n\sigma^2\right)$$

Let $X_1, \ldots, X_n$ be a **large** ($n > 30$) random sample from a population with mean $\mu$ and standard deviation $\sigma$, so that $\bar{X}$ is approximately normal. Then a level $100(1-\alpha)\%$ confidence interval for $\mu$ is

$$\bar{X} \pm z_{\frac{\alpha}{2}} \sigma_{\bar{X}}$$

where $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$. When the value of $\sigma$ is unknown, it can be replaced with the sample standard deviation $s$.

## Small Sample Confidence Intervals for a Population Mean

Let $X_1, \ldots, X_n$ be a small ($n < 30$) sample from a *normal* population with mean $\mu$. Then the quantity

$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$$

has a Student's $t$ distribution with $n - 1$ degrees of freedom, denoted $t_{n-1}$.

When $n$ is large, the distribution of quantity $\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}}$ is very close to normal, so the normal curve can be used, rather than the Student's $t$.

# Hypothesis Testing

## Large-Sample Tests for a Population Mean

Let $X_1, \ldots, X_n$ be a **large** ($n > 30$) sample form a population with mean $\mu$ and standard deviation $\sigma$.
To test a null hypothesis of the form $H_0 : u \le u_0$, $H_0 : u \ge \mu_0$, $H_0 := \mu_0$:

- Compute the $z$-score:

$$z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

If $\sigma$ is unknown it may be approximated with $s$.
- Computer the $P$-value. The $P$-value is an area under the normal curve, which depends on the alternate hypothesis as follows:

| Alternate Hypothesis | $P$-**Value** |
|---|---|
| $H_1 : \mu > \mu_0$ | Area to the right of $z$ |
| $H_1 : \mu < \mu_0$ | Area to the left of $z$ |
| $H_1 : \mu = \mu_0$ | Sum of the areas cut off by $z$ and $-z$ |

## Drawing Conclusions from the Results of Hypothesis Tests

Let $\alpha$ be any value between 0 and 1. Then, if $P \le \alpha$,

- The result of the test is said to be statistically significant at the $100\alpha\%$ level.
- The null hypothesis is rejected at the $100\alpha\%$ level.
- When reporting the result of the hypothesis test, report the $P$-value, rather than just comparing it to the 5% or 1%.

## Small-Sample Tests for a Population Mean

Let $X_1, \ldots, X_n$ be a **small** ($n \le 30$) random sample from a *normal* population with mean $\mu$ (unknown) and a standard deviation $\sigma$.
To test a null hypothesis of the form $H_0 : \mu \le \mu_0$, $H_0 : \mu \ge \mu_0$, or $H_0 : \mu = \mu_0$:

- Compute the test statistic

$$t^* = \frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}}$$

- Compute the $P$-value. The $P$-value is an area under the Student's $t$ curve with $n - 1$ degrees of freedom, which depends on the alternate hypothesis as follows

| Alternate Hypothesis | $P$-**Value** |
|---|---|
| $H_1 : \mu > \mu_0$ | Area to the right of $z$ |
| $H_1 : \mu < \mu_0$ | Area to the left of $z$ |
| $H_1 : \mu = \mu_0$ | Sum of the areas cut off by $z$ and $-z$ |

- If $\sigma$ is known, the test statistic is $z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$

## Large-Sample Tests for the Difference Between Two Means

Let $X_1, \ldots, X_n$ and $Y_1, \ldots, Y_n$ be **large** ($n_x > 30$ and $n_y > 30$) independent random samples from populations with mean $u_x$ and $u_y$ and standard deviation $\sigma_x$ and $\sigma_y$, respectively.
The test statistic is as follows:

$$z^* = \frac{(\bar{X} - \bar{Y}) - \delta_0}{\sqrt{\frac{\sigma_X^2}{n_x} + \frac{\sigma_Y^2}{n_y}}}$$

If $\sigma_X$ and $\sigma_Y$ are unknown they may be replaced by $s_X$ and $s_Y$, respectively

| Null Hypothesis | Alternative Hypothesis | $p$-**value** |
|---|---|---|
| $H_0 : \mu_x - \mu_y \le \delta_0$ | $H_1 : \mu_x - \mu_y > \delta_0$ | $P(Z \ge z^*)$ |
| $H_0 : \mu_x - \mu_y \ge \delta_0$ | $H_1 : \mu_x - \mu_y < \delta_0$ | $P(Z \le z^*)$ |
| $H_0 : \mu_x - \mu_y = \delta_0$ | $H_1 : \mu_x - \mu_y \ne \delta_0$ | $2 \times P(Z \ge |z^*|)$ |

## Small-Sample Tests for the Difference Between Two Means

### Population Variances Are Not Equal

Let $X_1, \ldots, X_{n,x}$ and $Y_1, \ldots, Y_{n,y}$ be samples from *normal* populations with means $\mu_X$ and $\mu_Y$ and standard deviations $\sigma_X$ and $\sigma_Y$, respectively. Assume the samples are drawn independently of each other.
**If $\sigma_x$ and $\sigma_Y$ are not known to be equal**, then, to test a null hypothesis of the form $H_0 : \mu_X - \mu_Y \le \Delta_0$, $H_0 : \mu_X - \mu_Y \ge \Delta_0$, or $H_0 : \mu_X - \mu_Y = \Delta_0$.

- Rounding down to the nearest integer, calculate

$$\nu = \frac{\left[\frac{s_X^2}{n_X} + \frac{s_Y^2}{n_Y}\right]^2}{\frac{\left(\frac{s_X^2}{n_X}\right)^2}{n_X - 1} + \frac{\left(\frac{s_Y^2}{n_Y}\right)^2}{n_Y - 1}}$$

- Compute the test statistic

$$t = \frac{(\bar{X} - \bar{Y}) - \Delta_0}{\sqrt{S_X^2/n_x + S_Y^2/n_Y}}$$

- Compute the $P$-value. The $P$-value is an area under the Student's $t$ curve with $v$ degrees of freedom, which depends on the alternate hypothesis as follows:

| Alternate Hypothesis | $P$-**value** |
|---|---|
| $H_1 : \mu_X - \mu_Y > \Delta_0$ | Area to the right of $t$ |
| $H_1 : \mu_X - \mu_Y < \Delta_0$ | Area to the left of $t$ |
| $H_1 : \mu_X - \mu_Y \ne \Delta_0$ | Sum of the areas in the tails cut off |

## Population Variances Are Equal

Let $X_1, \ldots, X_{n,x}$ and $Y_1, \ldots, Y_{n,y}$ be samples from *normal* populations with means $\mu_X$ and $\mu_Y$ and standard deviations $\sigma_X$ and $\sigma_Y$, respectively. Assume the samples are drawn independently of each other.
If $\sigma_X$ and $\sigma_Y$ are known to be qual, then, to test a null hypothesis of the fortm $H_0 : \mu_X - \mu_Y \le \Delta_0$, $H_0 : \mu_X - \mu_Y \ge \Delta_0$, or $H_0 : \mu_x - \mu_y = \Delta_0$:

- Compute

$$s_p = \sqrt{\frac{(n_X - 1)s_X^2 + (n_Y - 1)s_Y^2}{n_X + n_Y - 2}}$$

- Compute the test statistic

$$t = \frac{(\bar{X} - \bar{Y}) - \Delta_0}{s_p \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}}}$$

- Compue the $P$-value. The $P$-value is an area under the Student's $t$ curve with $n_X + n_Y - 2$ degrees of freedom, which depends on the alternate hypothesis as follows.

| Alternate Hypothesis | $P$-**value** |
|---|---|
| $H_1 : \mu_X - \mu_Y > \Delta_0$ | Area to the right of $t$ |
| $H_1 : \mu_X - \mu_Y < \Delta_0$ | Area to the left of $t$ |
| $H_1 : \mu_X - \mu_Y \ne \Delta_0$ | Sum of the areas in the tails cut off |

# Correlation vs. Causation

## Correlation

A correlation coefficient (denoted $r$) deasures the strength and direction of a linear relationship between two variables. Let $(x_1, y_1), \ldots, (x_n, y_n)$ represent bivariate data, then the correlation coefficient is

$$r = \frac{1}{n-1} \sum_{i=1}^{n} \left(\frac{x_i - \bar{x}}{s_x}\right)\left(\frac{x_i - \bar{x}}{s_x}\right)$$

$$= \frac{\sum_{i=1}^{n} x_i y_i - n\bar{x}\bar{y}}{\sqrt{\sum_{i=1}^{n} x_i^2 - n\bar{x}^2}\sqrt{\sum_{i=1}^{n} y_i^2 - n\bar{y}^2}} \qquad = \frac{SSR}{SST}$$

## The Least-Squares Line

For an equation of the form

$$y_1 = \beta_0 + \beta_1 x_i + \epsilon_i$$

$y_i$ is called the **dependent variable**, $x_i$ is called the **indepedent variable**, $\beta$ is called the **regression coefficients** (the least squares coefficients), and $\epsilon_i$ is called the **error**.
Also, $r^2$ is the **proportion of variance in $y$ explained by regression**.

$$e_1 = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{B}_1 x_i$$

$$\hat{\beta}_1 = \hat{B}_1 = \frac{\sum_i^n y_i - n\bar{x}\bar{y}}{\sum_i^n x_i^2 - n\bar{x}^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

## Uncertainties in the Least-Squares Coefficients

Using some assumptions,

- The quantity $\hat{\beta}$ is *normally distributed* random variables.
- The means of $\hat{\beta}$ is the true values of $\hat{\beta}$.
- The *standard deviations* of $\beta$ is estimated with

$$s_{\beta_0} = s\sqrt{\frac{1}{n} + \frac{\bar{x}}{2}\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

$$s_{\beta_1} = \frac{s}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}}$$

where

$$s = \sqrt{\frac{(1 - r^2) \sum_{i=1}^{n} (y_i - \bar{y})^2}{n - 2}}$$

is an estimate of the error standard deviation $\sigma$.

### Confidence Intervals for Coefficients

Under assumptions, the quantities $\frac{\widehat{\beta_1} - \beta_1}{s_{\beta_1}}$ and $\frac{\widehat{\beta_1} - \beta_1}{s_{\beta_1}}$ have Student's $t$ distributions with $n - 2$ degrees of freedom.
Level $100(1 - \alpha)\%$ confidence intervals for $\beta_0$ and $\beta_1$ are given by

$$\bar{\beta}_0 \pm t_{n-2} \times s_{\bar{\beta}_0} \qquad \bar{\beta}_1 \pm t_{n-2} \times s_{\bar{\beta}_1}$$

Level $100(1 - \alpha)\%$ confidence intervals for the quantity $\beta_0 + \beta_1 x$ is given by

$$\widehat{\beta_0} + \widehat{\beta_1} x \pm t_{n-2, \alpha/2} \times s_{\widehat{y}}$$

where

$$s_{\widehat{y}=s} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^{n} (x_i - \bar{x})^2}}$$

### Checking Assumptions

If the plot of residuals versus fitted values

- Shjows no substantial trend or curve, and
- Is **homoscedastic**, that is, the vertical spread does not varu too much along the horizontal length of the plot, except perhaps near the edges,

then it is *likely*, but not certain, that the assumptions of the linear model hold.

However, if the residual plot *does* show a substantial trend or curve, or is **heteroscedastic**, it is certain that the assumptions of the linear plot *do not* hold.

## Miscellaneous Notes

- A **test statistic** is a function of the sample data whose value is used to test a hypothesis
- A **p-value** is a measure of the disagreement between a sample and $H_0$.
- The smaller the $P$-value, the more certain we can be that $H_0$ is false and vice versa.
- For **large samples**, we approximate the population standad deviation $\sigma$ using the sample standard deviation $s$.
- The correlation coefficient is called the **sample correlation ($r$)**, and the it is an estimate of the population correleation ($\rho$).
- Some properties of the correlation coeffcient ($r$):

1. $-1 \leq r \leq 1$, $r$ is unitless.
2. If the points lie exactly on a horizontal or vertical line, the correlation coefficient is undefined, because one of the standard deviations is equal to zero.
3. Whenever $r \neq 0$, $x$ and $y$ are said to be correlated. If $r = 0$, $x$ and $y$ are said to be uncorrelated.
4. Correlation coefficient is unaffected by the units in whicht he measurements are made.

- For **small samples**, $s$ may be far from $\sigma$, which invalidates this large-sample method. However, when the population is approximately normal, the Student's $t$ distribution can be used.

- The **pooled** sample variance is

$$s_p^2 = \frac{(n_X - 1)s_X^2 + (n_Y - 1)S_Y^2}{n_X + n_Y - 2}$$

- The correleation coefficient remains unchanged under each of the following operations
  - Multiplying each value of a variable by a positive constant.
  - Adding a constant to each of a variable.
  - Interchanging the values of $x$ and $y$.

- A **goodness-of-fit statistic** measures how well a model explains a given set of data.
- $SST$ = total sum of squares = $\sum_{i=1}^{n} (y_i - \bar{y})^2$
- $SSE$ = error sum of squares = $\sum_{i=1}^{n} (y_i - \widehat{y})^2$
- $SSR$ = regression sum of squares = $SST - SSE$
- The following assumptions are satisfied.

  - The errors $\epsilon_1, \ldots, \epsilon_0$ are random and independent. In particular, the magnitude of any error $\epsilon_i$ does not influence the value of the next error $\epsilon_{i+1}$.
  - The errors $\epsilon_1, \ldots, \epsilon_0$ all have mean 0.
  - The errors $\epsilon_1, \ldots, \epsilon_0$ all have the same variance, which we denote by $\sigma^2$.
  - The errors $\epsilon_1, \ldots, \epsilon_0$ are normally distributed.
  - Margin of error = $t_{\frac{\alpha}{2}, \sqrt{n}} \times \frac{\sigma}{\sqrt{n}}$

- The sample variance is calculated

$$\frac{1}{N - 1} \sum_{i=0}^{n} (x - \bar{x})^2$$