

DATA1220-55 Cheat Sheet

Sarah E. Grabinski

2024-11-17

Formulas

Sample Proportion

$$\hat{p} = \frac{\text{count (something)}}{\text{count (everything)}} = \frac{\text{count (something)}}{n}$$

Sample Mean

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

Population Variance

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

Sample Variance

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

Population Standard Deviation

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}}$$

Sample Standard Deviation

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

Table 1: Standard Errors of Sample Statistics

Measure	SE	Calculation
Mean	$SE_{\bar{x}}$	$\frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}}$
Paired Difference in Means	$SE_{\bar{x}_{\text{difference}}}$	$\frac{\sigma_{\text{difference}}}{\sqrt{n}} \approx \frac{s_{\text{difference}}}{\sqrt{n}}$
Difference in Means	$SE_{\bar{x}_1 - \bar{x}_2}$	$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \approx \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$
Proportion	$SE_{\hat{p}}$	$\sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$
Difference in Proportions ($H_0: p_1 - p_2 = \mu$)	$SE_{\hat{p}_1 - \hat{p}_2}$	$\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \approx \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$
Difference in Proportions ($H_0: p_1 - p_2 = 0$)	$SE_{\hat{p}_1 - \hat{p}_2}$	$\sqrt{\frac{p_{\text{pool}}(1-p_{\text{pool}})}{n_1} + \frac{p_{\text{pool}}(1-p_{\text{pool}})}{n_2}} \approx \sqrt{\frac{\hat{p}_{\text{pool}}(1-\hat{p}_{\text{pool}})}{n_1} + \frac{\hat{p}_{\text{pool}}(1-\hat{p}_{\text{pool}})}{n_2}}$

Terms

Population The entire group being researched (e.g. sample, study, target)

Sample A subset of the population, ideally random and large enough to be representative

Sample size The total number of subjects or observations in the sample, represented by n .

Reliability The consistency of the observed measurements from a sample. Data from a sample is considered reliable estimate of the sample statistic when there is very little bias or measurement error.

Validity The degree to which the sample statistic approximates the population parameter. A sample statistic is considered a valid estimate of the population parameter when the sample is large and/or representative of the study population.

Median The middle value in the data separating the top 50% from the bottom 50%. Found by arranging all values from lowest to highest and taking the middle value (or mean of the 2 middle values)

Quartile Each of the 4 equal groups into a which a population can be divided. The divisions between the quartiles are $Q1 = 0.25$ (25th percentile), $Q2 = 0.50$ (50th percentile, median), and $Q3 = 0.75$ (75th percentile).

Interquartile Range (IQR) The difference between the 3rd quartile ($Q3 = 0.75$) and the 1st quartile ($Q1 = 0.25$). The middle 50% of the data.

Mean Also called the average. The sum of all values in the sample divided by number of values in the sample. μ (mu) represents the mean of a population, and \bar{x} represents the mean of a sample.

Variance Dispersion (spread) around the mean, determined by averaging the squared differences of all values from the mean. σ^2 (sigma squared) represents the variance of a population, and s^2 represents the variance of a sample.

Standard Deviation The square root of the variance. Also measures dispersion (spread) around the mean, but in the same units as the variable. σ (sigma) represents the standard deviation of a population, and s represents the standard deviation of a sample.

Central Limit Theorem The distribution of a sample statistic approximates the normal distribution N (population parameter, standard error) as $n \rightarrow \infty$.

Sampling Distribution the distribution of theoretically possible sample statistics from all samples of size n that can be taken from a population

Standard Error The standard deviation of a sampling distribution. Reflects how variable a sample statistic is expected to be from sample to sample.

Confidence Interval A range of values within which you expect the “true” population parameter to fall if you repeated the study an infinite number of times. The confidence level is the percentage of samples whose confidence interval would capture the “true” population parameter. A confidence intervals upper and lower bounds are found by calculating point estimate \pm critical value \times standard error.

Critical Value The number which defines the upper and lower bounds of a confidence interval from a given distribution. Its value corresponds to the probabilities $\alpha/2$ and $1 - \alpha/2$.

Null Hypothesis There is no meaningful relationship in the data. Represented as H_0 , gives the null distribution under which the hypothesis is tested.

Alternate Hypothesis There is something meaningful in the data. Represented as H_A , indicates whether the hypothesis test is one-sided (left- or right-tailed) or two-sided (both tails).

Test Statistic The standardized value of the observed sample statistic under the null hypothesis H_0 , used to find the p-value of a hypothesis test.

Type I Error The probability of rejecting the null hypothesis H_0 when H_0 is actually “true.” Represented by α .

Type II Error The probability of failing to reject the null hypothesis H_0 when H_0 is not actually “true.”

Inference

Means

Table 2: Sample Statistics for Inference of Population Means

Measure	Sample Statistic	Population Parameter
Mean	\bar{x}	μ
Paired Difference in Means	$\bar{x}_{\text{difference}}$	$\mu_{\text{difference}}$
Difference in Means	$\bar{x}_1 - \bar{x}_2$	$\mu_1 - \mu_2$
Standard Deviation	s	σ

Assumptions

- **Independence:** sample observations are independent (i.e. random sample).
- **Sample size:** the sample size should be greater than 30 ($n \geq 30$) with no extreme outliers.
- **Normality:** when the sample size n is small, observations come from a normally distribution population. This condition relaxes as $n \rightarrow \infty$.
- **Validity:** sample statistics approximate the population parameters ($\bar{x} \approx \mu$, $s \approx \sigma$)

Single Mean (\bar{x}) - One-Sample t -test

- A **1-sample** t -test tests if the mean (μ) of a population is different from a null value (μ_0).
- Sample statistics \bar{x} (mean) and s (standard deviation) and the sample size n are used to infer the sampling distribution of the mean $\bar{x} \sim N(\mu, SE_{\bar{x}})$.
- To account for using $s \approx \sigma$ in the standard error, confidence intervals and hypothesis tests are based on the T distribution (Student’s t) with the parameter degrees of freedom $(df) = n - 1$.

Confidence Interval

The confidence interval for a single population mean μ is calculated using the point estimate \bar{x} .

$$\bar{x} \pm T_{df}^* \times SE_{\bar{x}}$$

The standard error of \bar{x} is estimated from the observed standard deviation s and the sample size n .

$$SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}}$$

The critical value from the t distribution with degrees of freedom $df = n - 1$ is...

$$T_{df}^* = T_{df, \alpha/2} = T_{df, 1-\alpha/2}$$

A critical value is calculated from the t distribution in R using the function `qt()`. This function takes a probability `p` ($\alpha/2$ or $1 - \alpha/2$) and degrees of freedom `df` ($n - 1$).

```
qt(alpha/2, df=n-1)
qt(1-alpha/2, df=n-1)
```

Hypothesis Test

The null hypothesis for a 1-sample t -test states that the population mean μ is equal to some null value μ_0 .

$$H_0: \mu = \mu_0$$

The null distribution of the sample mean \bar{x}_0 , given the null hypothesis that $\mu = \mu_0$, is $\bar{x}_0 \sim N(\mu_0, SE_{\bar{x}})$. The standard error of the null sample mean \bar{x}_0 is the same as the standard error of the observed sample mean \bar{x} .

$$SE_{\bar{x}_0} = SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}}$$

A confidence interval for the null population mean μ under the null hypothesis $H_0: \mu = \mu_0$ is calculated using μ_0 as the point estimate.

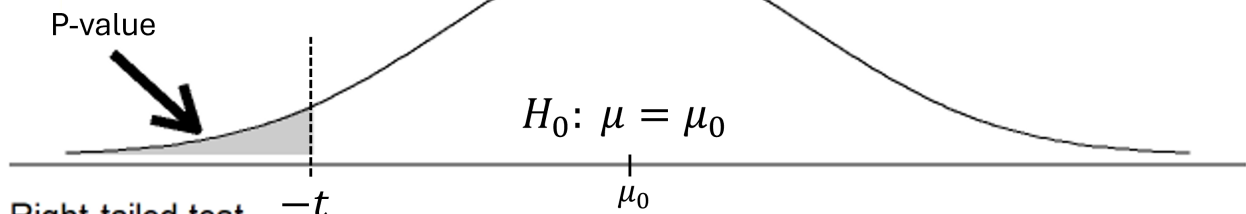
$$\mu_0 \pm T_{df}^* \times SE_{\bar{x}}$$

The alternate hypothesis of a 1-sample t -test states that the population mean μ is greater than, less than, or not equal to some null value μ_0 .

- $H_A: \mu < \mu_0$, left-tailed test (one-sided)
- $H_A: \mu > \mu_0$, right-tailed test (one-sided)
- $H_A: \mu \neq \mu_0$, two-tailed test (two-sided)

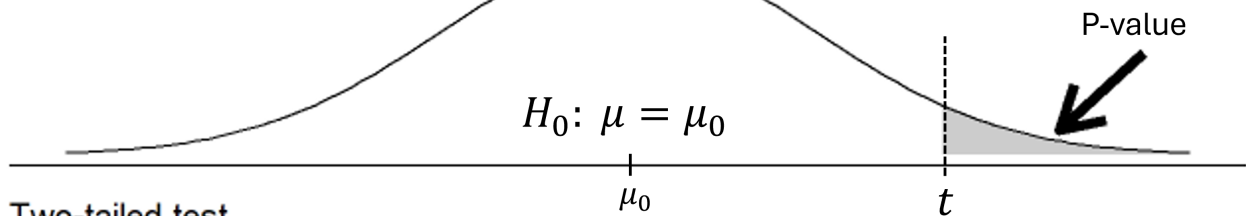
Left-tailed test

$$H_A: \mu < \mu_0$$



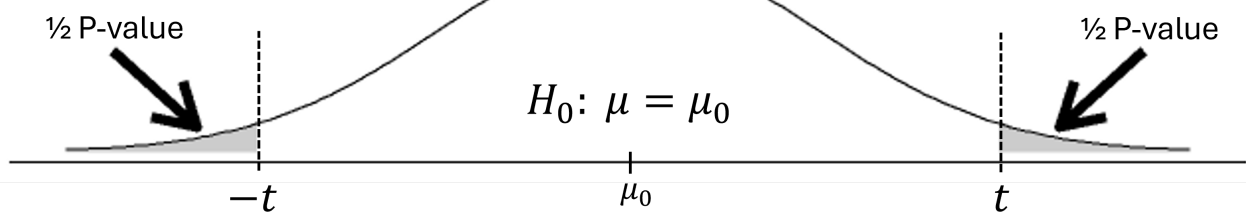
Right-tailed test

$$H_A: \mu > \mu_0$$



Two-tailed test

$$H_A: \mu \neq \mu_0$$



The test statistic t (the t -statistic) calculates how many standard errors ($SE_{\bar{x}}$) away from the null hypothesis μ_0 that the observed statistic \bar{x} is.

$$t = \frac{\bar{x} - \mu_0}{SE_{\bar{x}}} = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

The t -statistic is used to find the probability of your sample having the mean \bar{x} if the null distribution $N(\mu_0, SE_{\bar{x}})$ were the “true” distribution in your population.

The probability of the sample statistic \bar{x} under the null hypothesis $\mu = \mu_0$ and $\bar{x} \sim N(\mu_0, SE_{\bar{x}})$ is calculated from the t distribution in R using the function `pt()`. This function takes the test statistic t (called `q` or quantile in R) and the degrees of freedom `df` ($n - 1$) as parameters. This probability is known as the p-value.

```
# left-tailed hypothesis test
pt(-t, df=n-1)

# right-tailed hypothesis test
pt(t, df=n-1, lower.tail=F)

# two-tailed hypothesis test
pt(-t, df=n-1)*2
```

```
pt(t, df=n-1, lower.tail=F)*2
pt(-t, df=n-1)+pt(t, df=n-1, lower.tail=F)
```

A small p-value indicates that the probability of taking a sample of size n with your observed sample mean \bar{x} from the null sampling distribution $\bar{x}_0 \sim N(\mu_0, SE_{\bar{x}_0})$ is very low.

If the p-value for your t -statistic is less than your significance threshold α (i.e. the Type I Error Rate), then this is sufficient evidence that the null hypothesis $H_0: \mu = \mu_0$ may not be true. Reject the null hypothesis H_0 and accept the alternate hypothesis H_A .

If the p-value for your t -statistic is greater than α , this is not sufficient evidence against the null hypothesis $H_0: \mu = \mu_0$. You fail to reject the null hypothesis H_0 .

Paired Means

- A **paired means** t -test is a special instance of a 1-sample t -test. It tests if the mean of the difference between 2 paired measures ($\mu_{\text{difference}}$) in a population is different from a null value (μ_0). Typically, $\mu_0 = 0$.
- Sample statistics mean $\bar{x}_{\text{difference}}$ and standard deviation $s_{\text{difference}}$ are calculated using the difference between the 2 paired measures $x_{\text{difference}} = x_1 - x_2$, not the individual measures x_1 or x_2 .
- Sample statistics $\bar{x}_{\text{difference}}$ (mean) and $s_{\text{difference}}$ (standard deviation) and the sample size n are used to infer the sampling distribution of the mean $\bar{x}_{\text{difference}} \sim N(\mu_{\text{difference}}, SE_{\bar{x}_{\text{difference}}})$.
- To account for using $s \approx \sigma$ in the standard error, confidence intervals and hypothesis tests are based on the T distribution (Student's t) with the parameter degrees of freedom (df) = $n - 1$.

Confidence Interval

The confidence interval for the population difference in paired means $\mu_{\text{difference}}$ is calculated using the point estimate $\bar{x}_{\text{difference}}$.

$$\bar{x}_{\text{difference}} \pm T_{\text{df}}^* \times SE_{\bar{x}_{\text{difference}}}$$

The standard error of $\bar{x}_{\text{difference}}$ is estimated from the observed standard deviation $s_{\text{difference}}$ and the sample size n .

$$SE_{\bar{x}_{\text{difference}}} = \frac{\sigma_{\text{difference}}}{\sqrt{n}} \approx \frac{s_{\text{difference}}}{\sqrt{n}}$$

The critical value from the t distribution with degrees of freedom $\text{df} = n - 1$ is...

$$T_{\text{df}}^* = T_{\text{df}, \alpha/2} = T_{\text{df}, 1-\alpha/2}$$

A critical value is calculated from the t distribution in R using the function `qt()`. This function takes a probability p ($\alpha/2$ or $1 - \alpha/2$) and degrees of freedom `df` ($n - 1$).

```
qt(alpha/2, df=n-1)
qt(1-alpha/2, df=n-1)
```

Hypothesis Test

The null hypothesis for a paired means t -test states that the population mean $\mu_{\text{difference}}$ is equal to some null value μ_0 . Typically, $\mu_0 = 0$.

$$H_0 : \mu_{\text{difference}} = \mu_0$$

The null distribution of the sample mean $\bar{x}_{\text{difference}0}$, given the null hypothesis that $\mu_{\text{difference}} = \mu_0$, is $\bar{x}_{\text{difference}0} \sim N(\mu_0, SE_{\bar{x}_{\text{difference}}})$. The standard error of the null sample mean $\bar{x}_{\text{difference}0}$ is the same as the standard error of the observed sample mean $\bar{x}_{\text{difference}}$.

$$SE_{\bar{x}_{\text{difference}0}} = SE_{\bar{x}_{\text{difference}}} = \frac{\sigma_{\text{difference}}}{\sqrt{n}} \approx \frac{s_{\text{difference}}}{\sqrt{n}}$$

A confidence interval for the null population mean $\mu_{\text{difference}}$ under the null hypothesis $H_0 : \mu_{\text{difference}} = \mu_0$ is calculated using μ_0 as the point estimate.

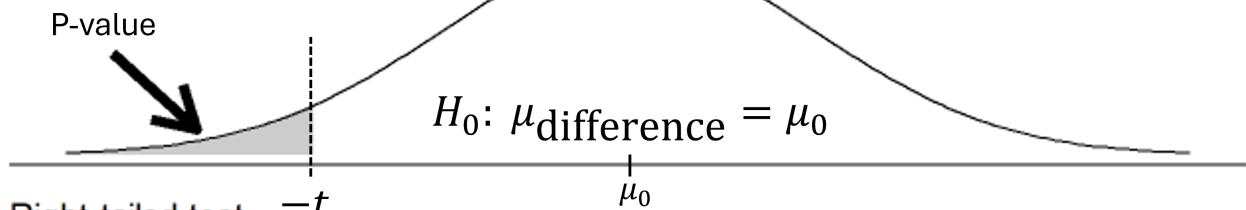
$$\mu_0 \pm T_{\text{df}}^* \times SE_{\bar{x}_{\text{difference}}}$$

The alternate hypothesis of a 1-sample t -test states that the population mean $\mu_{\text{difference}}$ is greater than, less than, or not equal to some null value μ_0 .

- $H_A : \mu_{\text{difference}} < \mu_0$, left-tailed test (one-sided)
- $H_A : \mu_{\text{difference}} > \mu_0$, right-tailed test (one-sided)
- $H_A : \mu_{\text{difference}} \neq \mu_0$, two-tailed test (two-sided)

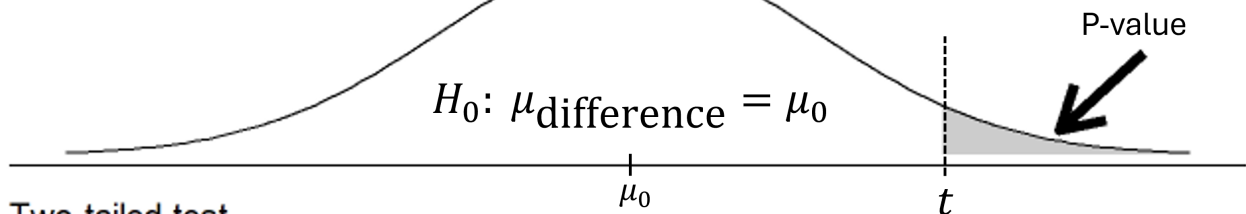
Left-tailed test

$$H_A: \mu_{\text{difference}} < \mu_0$$



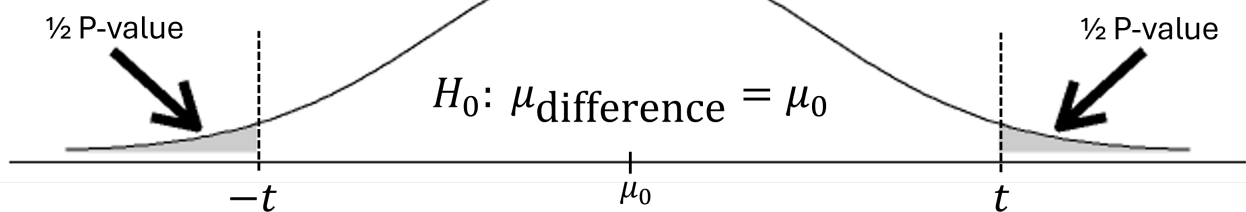
Right-tailed test

$$H_A: \mu_{\text{difference}} > \mu_0$$



Two-tailed test

$$H_A: \mu_{\text{difference}} \neq \mu_0$$



The test statistic t (the t -statistic) calculates how many standard errors ($SE_{\bar{x}_{\text{difference}}}$) away from the null hypothesis μ_0 that the observed statistic $\bar{x}_{\text{difference}}$ is.

$$t = \frac{\bar{x}_{\text{difference}} - \mu_0}{SE_{\bar{x}_{\text{difference}}}} = \frac{\bar{x}_{\text{difference}} - \mu_0}{\frac{s_{\text{difference}}}{\sqrt{n}}}$$

The t -statistic is used to find the probability of your sample having the mean $\bar{x}_{\text{difference}}$ if the null distribution $N(\mu_0, SE_{\bar{x}_{\text{difference}}})$ were the “true” distribution in your population for $\bar{x}_{\text{difference}}$.

The probability of the sample statistic $\bar{x}_{\text{difference}}$ under the null hypothesis $\mu_{\text{difference}} = \mu_0$ and $\bar{x}_{\text{difference}} \sim N(\mu_0, SE_{\bar{x}_{\text{difference}}})$ is calculated from the t distribution in R using the function `pt()`. This function takes the test statistic t (called `q` or quantile in R) and the degrees of freedom `df` ($n - 1$) as parameters. This probability is known as the p-value.

```
# left-tailed hypothesis test
pt(-t, df=n-1)

# right-tailed hypothesis test
pt(t, df=n-1, lower.tail=F)

# two-tailed hypothesis test
pt(-t, df=n-1)*2
```



```
pt(t, df=n-1, lower.tail=F)*2
pt(-t, df=n-1)+pt(t, df=n-1, lower.tail=F)
```

A small p-value indicates that the probability of taking a sample of size n with your observed sample mean $\bar{x}_{\text{difference}}$ from the null sampling distribution $\bar{x}_{\text{difference}0} \sim N(\mu_0, SE_{\bar{x}_{\text{difference}}})$ is very low.

If the p-value for your t -statistic is less than your significance threshold α (i.e. the Type I Error Rate), then this is sufficient evidence that the null hypothesis $H_0: \mu_{\text{difference}} = \mu_0$ may not be true. Reject the null hypothesis H_0 and accept the alternate hypothesis H_A .

If the p-value for your t -statistic is greater than α , this is not sufficient evidence against the null hypothesis $H_0: \mu_{\text{difference}} = \mu_0$. You fail to reject the null hypothesis H_0 .

Two Means

- A **2-sample** t -test tests if the difference between 2 means ($\mu_1 - \mu_2$) of a population is different from a null value (μ_0). Often, $\mu_0 = 0$.
- Sample statistics $\bar{x}_1 - \bar{x}_2$ (difference in means), standard deviations (s_1 and s_2), and the sample sizes n_1 and n_2 are used to infer the sampling distribution of the difference in means $\bar{x}_1 - \bar{x}_2 \sim N(\mu_1 - \mu_2, SE_{\bar{x}_1 - \bar{x}_2})$.
- To account for using $s \approx \sigma$ in the standard error, confidence intervals and hypothesis tests are based on the T distribution (Student's t) with the parameter degrees of freedom $(df) = \min(n_1, n_2) - 1$.

Confidence Interval

The confidence interval for a difference in means $\mu_1 - \mu_2$ is calculated using the point estimate $\bar{x}_1 - \bar{x}_2$.

$$\bar{x}_1 - \bar{x}_2 \pm T_{df}^* \times SE_{\bar{x}_1 - \bar{x}_2}$$

The standard error of $\bar{x}_1 - \bar{x}_2$ is estimated from the 2 observed standard deviations s_1 and s_2 and the 2 sample sizes n_1 and n_2 .

$$SE_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \approx \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

The critical value from the t distribution with degrees of freedom $df = \min(n_1, n_2) - 1$ is...

$$T_{df}^* = T_{df, \alpha/2} = T_{df, 1-\alpha/2}$$

A critical value is calculated from the t distribution in R using the function `qt()`. This function takes a probability `p` ($\alpha/2$ or $1 - \alpha/2$) and degrees of freedom `df` ($\min(n_1, n_2) - 1$).

```
qt(alpha/2, df=min(n1, n2)-1)
qt(1-alpha/2, df=min(n1, n2)-1)
```

Hypothesis Test

The null hypothesis for a 2-sample t -test states that the population difference in means $\mu_1 - \mu_2$ is equal to some null value μ_0 . Often, $\mu_0 = 0$.

$$H_0: \mu = \mu_0$$

The null distribution of the difference in sample means $\bar{x}_1 - \bar{x}_2$, given the null hypothesis that $\mu_1 - \mu_2 = \mu_0$, is $\bar{x}_1 - \bar{x}_2 \sim N(\mu_0, SE_{\bar{x}_1 - \bar{x}_2})$. The standard error of the null difference in sample means $\bar{x}_1 - \bar{x}_2$ is the same as the standard error of the observed difference in sample means.

$$SE_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \approx \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

A confidence interval for the null difference in population means $\mu_1 - \mu_2$ under the null hypothesis $H_0: \mu_1 - \mu_2 = \mu_0$ is calculated using μ_0 as the point estimate.

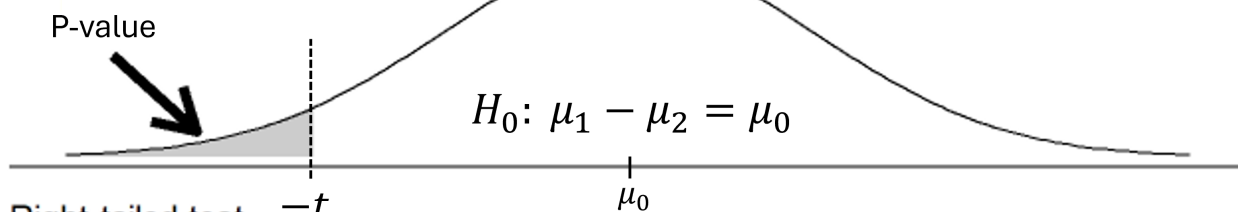
$$\mu_0 \pm T_{df}^* \times SE_{\bar{x}_1 - \bar{x}_2}$$

The alternate hypothesis of a 2-sample t -test states that the population difference in means $\mu_1 - \mu_2$ is greater than, less than, or not equal to some null value μ_0 .

- $H_A: \mu_1 - \mu_2 < \mu_0$, left-tailed test (one-sided)
- $H_A: \mu_1 - \mu_2 > \mu_0$, right-tailed test (one-sided)
- $H_A: \mu_1 - \mu_2 \neq \mu_0$, two-tailed test (two-sided)

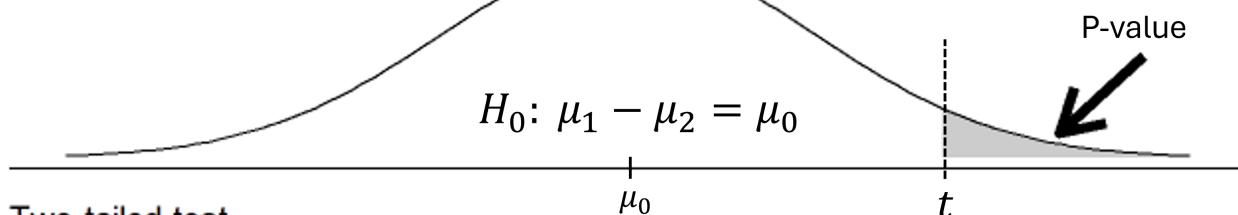
Left-tailed test

$$H_A: \mu_1 - \mu_2 < \mu_0$$



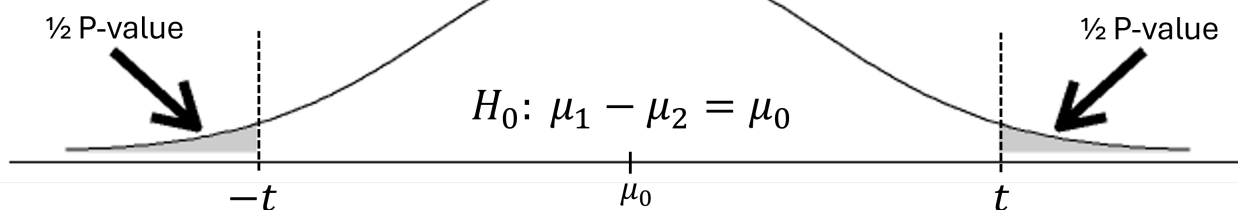
Right-tailed test

$$H_A: \mu_1 - \mu_2 > \mu_0$$



Two-tailed test

$$H_A: \mu_1 - \mu_2 \neq \mu_0$$



The test statistic t (the t -statistic) calculates how many standard errors ($SE_{\bar{x}_1 - \bar{x}_2}$) away from the null hypothesis μ_0 that the observed statistic $\bar{x}_1 - \bar{x}_2$ is.

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - \mu_0}{SE_{\bar{x}_1 - \bar{x}_2}} = \frac{(\bar{x}_1 - \bar{x}_2) - \mu_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

When the null hypothesis is that the population difference in means is 0 ($H_0: \mu_1 - \mu_2 = 0$), the calculation for the test statistic simplifies to the below.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{SE_{\bar{x}_1 - \bar{x}_2}} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

The t -statistic is used to find the probability of your sample having the difference in means $\bar{x}_1 - \bar{x}_2$ if the null distribution $N(\mu_0, SE_{\bar{x}_1 - \bar{x}_2})$ were the “true” distribution in your population.

The probability of the sample statistic $\bar{x}_1 - \bar{x}_2$ under the null hypothesis $\mu_1 - \mu_2 = \mu_0$ and $\bar{x}_1 - \bar{x}_2 \sim N(\mu_0, SE_{\bar{x}_1 - \bar{x}_2})$ is calculated from the t distribution in R using the function `pt()`. This function takes the test statistic t (called `q` or quantile in R) and the degrees of freedom `df` (`minimum(n1, n2) - 1`) as parameters. This probability is known as the p-value.

```
# left-tailed hypothesis test
pt(-t, df=min(n1, n2)-1)

# right-tailed hypothesis test
pt(t, df=min(n1, n2)-1, lower.tail=F)

# two-tailed hypothesis test
pt(-t, df=min(n1, n2)-1)*2
pt(t, df=min(n1, n2)-1, lower.tail=F)*2
pt(-t, df=min(n1, n2)-1)+pt(t, df=min(n1, n2)-1, lower.tail=F)
```

A small p-value indicates that the probability of taking samples of size n_1 and n_2 with your observed difference in sample means $\bar{x}_1 - \bar{x}_2$ from the null sampling distribution $\bar{x}_1 - \bar{x}_2 \sim N(\mu_0, SE_{\bar{x}_1 - \bar{x}_2})$ is very low.

If the p-value for your t -statistic is less than your significance threshold α (i.e. the Type I Error Rate), then this is sufficient evidence that the null hypothesis $H_0: \mu_1 - \mu_2 = \mu_0$ may not be true. Reject the null hypothesis H_0 and accept the alternate hypothesis H_A .

If the p-value for your t -statistic is greater than α , this is not sufficient evidence against the null hypothesis $H_0: \mu_1 - \mu_2 = \mu_0$. You fail to reject the null hypothesis H_0 .

Proportions

Table 3: Sample Statistics for Inference of Population Proportions

Measure	Sample Statistic	Population Parameter
Proportion	\hat{p}	p
Difference in Proportions	$\hat{p}_1 - \hat{p}_2$	$p_1 - p_2$

Assumptions

- **Independence:** sample observations are independent (i.e. random sample).
- **Sample size:** the sample size should be greater than 20 ($n \geq 20$) with at least 10 successes ($np \geq 10$) and 10 failures ($n(1 - p) \geq 10$).
- **Validity:** sample statistics approximate the population parameters ($\bar{x} \approx \mu$, $s \approx \sigma$)

Single Proportion (\hat{p}) - One-Sample Z -test

- A **1-sample Z -test** tests if the mean proportion (p) for a population is different from a null value (p_0).
- Sample statistic \hat{p} (proportion) and sample size n are used to infer the sampling distribution of the mean proportion $\hat{p} \sim N(p, SE_{\hat{p}})$.
- Confidence intervals and hypothesis tests for proportions are based on the Z distribution (standard normal).

Confidence Interval

The confidence interval for a population proportion p is estimated from the sample proportion \hat{p} .

$$\hat{p} \pm Z^* \times SE_{\hat{p}}$$

The standard error of \hat{p} is estimated from the sample proportion \hat{p} and the sample size n .

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

The critical value from the Z distribution is...

$$Z^* = Z_{\alpha/2} = Z_{1-\alpha/2}$$

A critical value is calculated from the Z distribution in R using the function `qnorm()`. This function takes a probability `p` ($\alpha/2$ or $1 - \alpha/2$), a `mean` (default = 0), and a standard deviation (`sd`, default = 1). However, you do not need to include the `mean` and `sd` parameters when you are using the default standard normal (Z) distribution.

```
qnorm(alpha/2)
```

```
qnorm(1-alpha/2)
```

Hypothesis Test

The null hypothesis of a 1-sample proportion- or Z -test states that the population proportion p is equal to some null value p_0 .

$$H_0: p = p_0$$

The null distribution of the sample proportion \hat{p}_0 , given the null hypothesis $H_0: p = p_0$, is $\hat{p} \sim N(p_0, SE_{\hat{p}_0})$. The standard error of the null sample proportion \hat{p}_0 under the null hypothesis $H_0: p = p_0$ is estimated using the null hypothesis value p_0 and the sample size n .

$$SE_{\hat{p}_0} = \sqrt{\frac{p_0(1-p_0)}{n}}$$

A confidence interval for the mean proportion under the null hypothesis is calculated using p_0 as the point estimate and $SE_{\hat{p}_0}$ under the null hypothesis.

$$p_0 \pm Z^* \times SE_{\hat{p}_0}$$

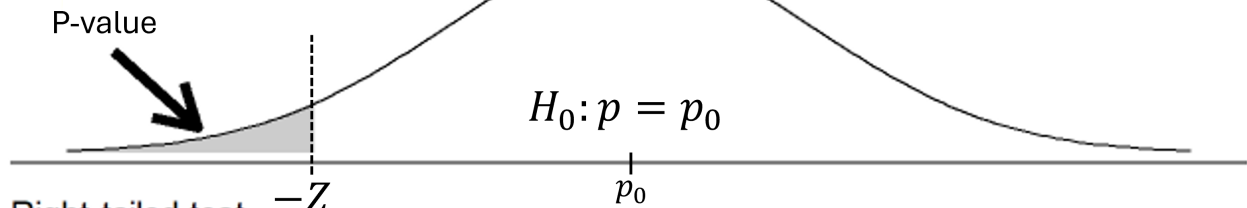
The alternate hypothesis of a 1-sample Z -test states that the population proportion p is greater than, less than, or not equal to some null value p_0 .

- $H_A: p < p_0$, left-tailed test (one-sided)
- $H_A: p > p_0$, right-tailed test (one-sided)

- $H_A: p \neq p_0$, two-tailed test (two-sided)

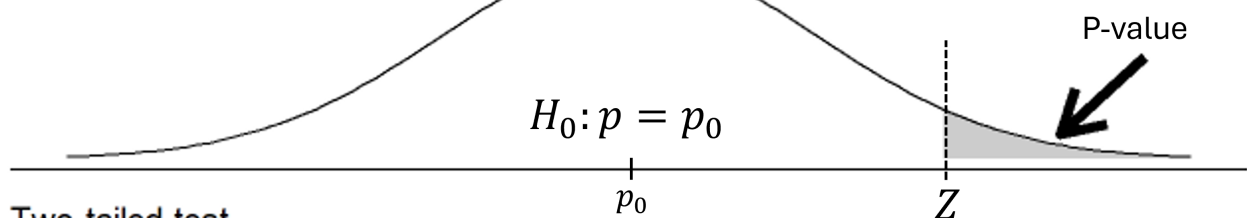
Left-tailed test

$$H_A: p < p_0$$



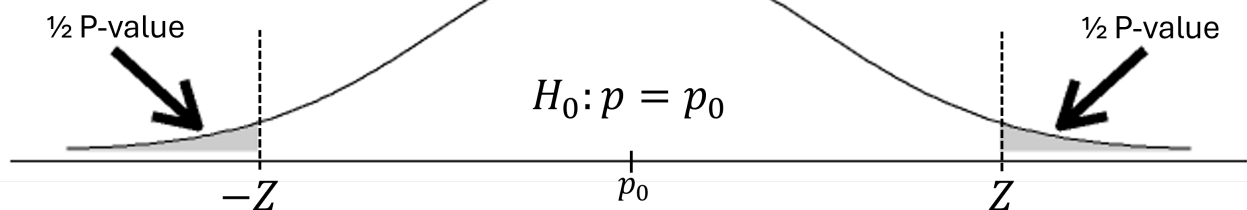
Right-tailed test

$$H_A: p > p_0$$



Two-tailed test

$$H_A: p \neq p_0$$



The test statistic Z (the Z -statistic) calculates how many standard errors ($SE_{\hat{p}_0}$) away from the null hypothesis p_0 that the observed statistic \hat{p} is.

$$Z = \frac{\hat{p} - p_0}{SE_{\hat{p}_0}} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

The Z -statistic is used to find the probability of your sample having the mean proportion \hat{p} if the null distribution $N(p_0, SE_{\hat{p}_0})$ were the “true” distribution in your population.

The probability of the sample statistic \hat{p} under the null hypothesis $H_0: p = p_0$ is calculated from the Z distribution in R using the function `pnorm()`. This function takes the test statistic Z (called `q` or quantile in R), the `mean` (default = 0), and the standard deviation (`sd`, default = 1) as parameters. This probability is known as the p-value. You do not need to include the `mean` and `sd` parameters when you are using the default standard normal (Z) distribution.

```
# left-tailed hypothesis test
pnorm(-Z)

# right-tailed hypothesis test
pnorm(Z, lower.tail=F)
```

```
# two-tailed hypothesis test
pnorm(-Z)*2
pnorm(Z, lower.tail=F)*2
pnorm(-Z)+pt(Z, lower.tail=F)
```

A small p-value indicates that the probability of taking a sample of size n with your observed sample proportion \hat{p} from the null sampling distribution $\hat{p}_0 \sim N(p_0, SE_{\hat{p}_0})$ is very low.

If the p-value for your Z -statistic is less than your significance threshold α (i.e. the Type I Error Rate), this is considered sufficient evidence that the null hypothesis $H_0: p = p_0$ may not be true. Reject the null hypothesis H_0 and accept the alternate hypothesis H_A .

If the p-value for your Z -statistic is greater than α , this is not sufficient evidence against the null hypothesis $H_0: p = p_0$. You fail to reject the null hypothesis H_0 .

Difference in Proportions ($\hat{p}_1 - \hat{p}_2$) - Two-Sample Z -test

- A **2-sample** Z -test for a difference in proportions tests if the difference ($p_1 - p_2$) between the mean proportion in population 1 (p_1) and the mean proportion in population 2 (p_2) is different from a null value μ . Typically, $\mu = 0$.
- Sample statistics \hat{p}_1 and \hat{p}_2 (proportions) and sample size n are used to infer the sampling distribution of the difference $\hat{p}_1 - \hat{p}_2 \sim N(p_1 - p_2, SE_{\hat{p}_1 - \hat{p}_2})$.
- Confidence intervals and hypothesis tests for the difference between proportions are based on the Z distribution (standard normal).

Confidence Interval

The confidence interval for a difference in population proportions $p_1 - p_2$ is estimated from the difference in sample proportions $\hat{p}_1 - \hat{p}_2$.

$$\hat{p}_1 - \hat{p}_2 \pm Z^* \times SE_{\hat{p}_1 - \hat{p}_2}$$

The standard error of $\hat{p}_1 - \hat{p}_2$ is estimated from the observed sample proportions, \hat{p}_1 and \hat{p}_2 , and the sample sizes, n_1 and n_2 .

$$\begin{aligned} SE_{\hat{p}_1 - \hat{p}_2} &= \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \\ &\approx \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \end{aligned}$$

The critical value from the Z distribution is...

$$Z^* = Z_{\alpha/2} = Z_{1-\alpha/2}$$

A critical value is calculated from the Z distribution in R using the function `qnorm()`. This function takes a probability `p` ($\alpha/2$ or $1 - \alpha/2$), a `mean` (default = 0), and a standard deviation (`sd`, default = 1). You

do not need to include the `mean` and `sd` parameters when you are using the default standard normal (Z) distribution.

```
qnorm(alpha/2)

qnorm(1-alpha/2)
```

Hypothesis Test

The null hypothesis of a 2-sample proportion- or Z -test states that the difference in population proportions $p_1 - p_2$ is equal to some null value p_0 .

$$H_0: p_1 - p_2 = p_0$$

The null distribution of $\hat{p}_1 - \hat{p}_2$ given the null hypothesis that $p_1 - p_2 = p_0$ is $\hat{p}_1 - \hat{p}_2 \sim N(p_0, SE_{\hat{p}_1 - \hat{p}_2})$. The standard error calculation for the null distribution of $\hat{p}_1 - \hat{p}_2$ depends on the null hypothesis H_0 .

When the null hypothesis is that the difference between proportions is some null value p_0 ($H_0: p_1 - p_2 = p_0$), \hat{p}_1 is used for p_1 and \hat{p}_2 for p_2 in the standard error calculation.

$$\begin{aligned} \text{When } H_0: p_1 - p_2 = p_0, \quad SE_{\hat{p}_1 - \hat{p}_2} &= \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \\ &\approx \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \end{aligned}$$

However, when the null hypothesis is that there is no difference between proportions ($H_0: p_1 - p_2 = 0$), the pooled proportion \hat{p}_{pool} across both groups is used for p_1 and p_2 in the standard error calculation.

$$\begin{aligned} \text{When } H_0: p_1 - p_2 = 0, \quad SE_{\hat{p}_1 - \hat{p}_2} &= \sqrt{\frac{p_{\text{pool}}(1-p_{\text{pool}})}{n_1} + \frac{p_{\text{pool}}(1-p_{\text{pool}})}{n_2}} \\ &\approx \sqrt{\frac{\hat{p}_{\text{pool}}(1-\hat{p}_{\text{pool}})}{n_1} + \frac{\hat{p}_{\text{pool}}(1-\hat{p}_{\text{pool}})}{n_2}} \end{aligned}$$

A confidence interval for the difference in proportions under the null hypothesis is calculated using p_0 as the point estimate and $SE_{\hat{p}_1 - \hat{p}_2}$ under the null hypothesis.

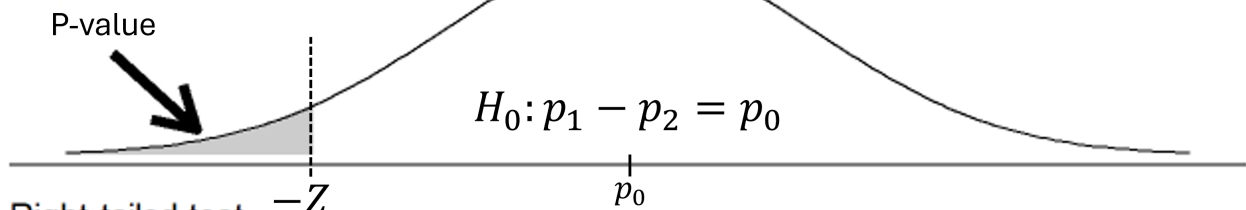
$$p_0 \pm Z^* \times SE_{\hat{p}_1 - \hat{p}_2}$$

The alternate hypothesis of a 2-sample proportion- or Z -test states that the difference in population proportions $p_1 - p_2$ is greater than, less than, or not equal to some null value p_0 .

- $H_A: p < p_0$, left-tailed test (one-sided)
- $H_A: p > p_0$, right-tailed test (one-sided)
- $H_A: p \neq p_0$, two-tailed test (two-sided)

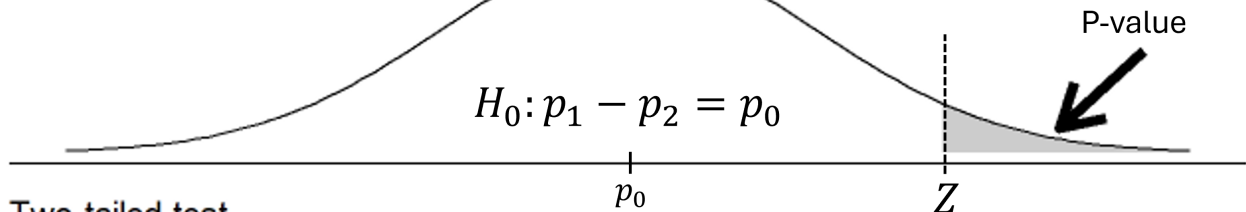
Left-tailed test

$$H_A: p_1 - p_2 < p_0$$



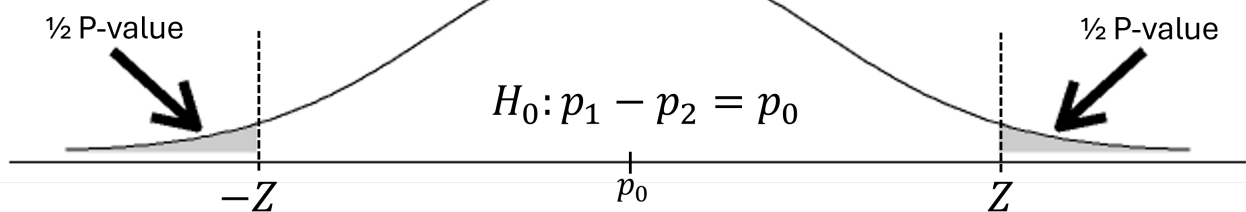
Right-tailed test

$$H_A: p_1 - p_2 > p_0$$



Two-tailed test

$$H_A: p_1 - p_2 \neq p_0$$



The test statistic Z (the Z -statistic) calculates how many standard errors ($SE_{\hat{p}_1 - \hat{p}_2}$) away from the null hypothesis p_0 that the observed statistic $\hat{p}_1 - \hat{p}_2$ is.

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - p_0}{SE_{\hat{p}_1 - \hat{p}_2}} = \begin{cases} \frac{\frac{(\hat{p}_1 - \hat{p}_2) - p_0}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}}{1} & \text{when } H_0: p_1 - p_2 = p_0 \\ \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\frac{\hat{p}_{\text{pool}}(1-\hat{p}_{\text{pool}})}{n_1} + \frac{\hat{p}_{\text{pool}}(1-\hat{p}_{\text{pool}})}{n_2}}} & \text{when } H_0: p_1 - p_2 = 0 \end{cases}$$

The Z -statistic is used to find the probability of your samples having the difference in proportions $\hat{p}_1 - \hat{p}_2$ if the null distribution $N(p_0, SE_{\hat{p}_1 - \hat{p}_2})$ were the “true” distribution in your population.

The probability of the sample statistic $\hat{p}_1 - \hat{p}_2$ under the null hypothesis $p_1 - p_2 = p_0$ is calculated from the Z distribution in R using the function `pnorm()`. This function takes the test statistic Z (called `q` or quantile in R), the `mean` (default = 0), and the standard deviation (`sd`, default = 1) as parameters. This probability is known as the p-value. You do not need to include the `mean` and `sd` parameters when you are using the default standard normal (Z) distribution.

```
# left-tailed hypothesis test
pnorm(-Z)

# right-tailed hypothesis test
pnorm(Z, lower.tail=F)
```

```
# two-tailed hypothesis test
pnorm(-Z)*2
pnorm(Z, lower.tail=F)*2
pnorm(-Z)+pt(Z, lower.tail=F)
```

A small p-value indicates that the probability of taking samples of size n_1 and n_2 with your observed difference in sample proportions $\hat{p}_1 - \hat{p}_2$ from the null sampling distribution $\hat{p}_1 - \hat{p}_2 \sim N(p_0, SE_{\hat{p}_1 - \hat{p}_2})$ is very low.

If the p-value for your Z -statistic is less than your significance threshold α (i.e. the Type I Error Rate), this is considered sufficient evidence that the null hypothesis $H_0: p_1 - p_2 = p_0$ may not be true. Reject the null hypothesis H_0 and accept the alternate hypothesis H_A .

If the p-value for your Z -statistic is greater than α , this is not sufficient evidence against the null hypothesis $H_0: p_1 - p_2 = p_0$. You fail to reject the null hypothesis H_0 .