

Table of Contents

1. Foundations of Inferential Statistics

- 1.1 Standard Normal Distribution (z-scores)
- 1.2 t-Distribution (small sample inference)

2. Sampling & Estimation

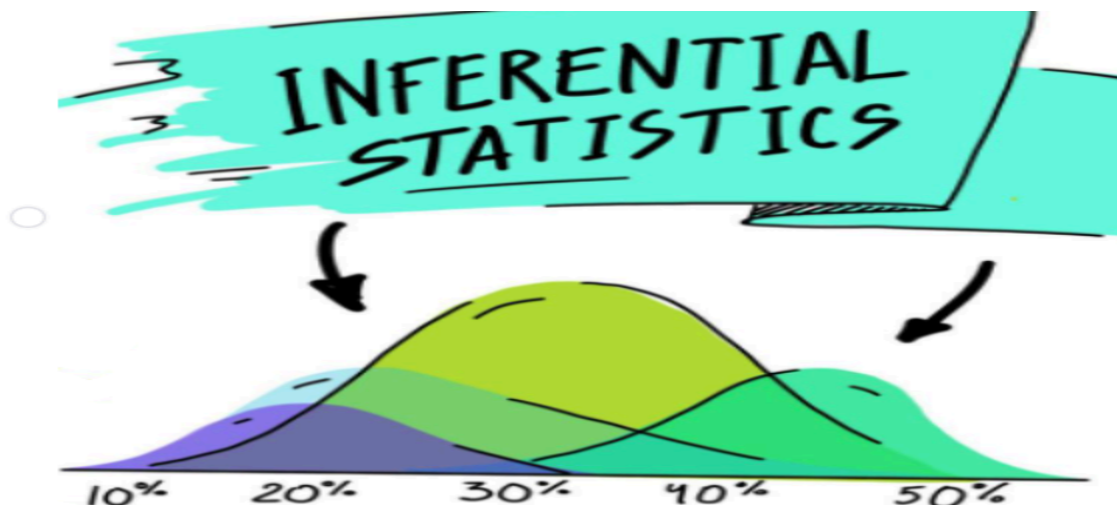
- 2.1 Central Limit Theorem (CLT)
- 2.2 Confidence Intervals (estimating population parameters)

3. Hypothesis Testing

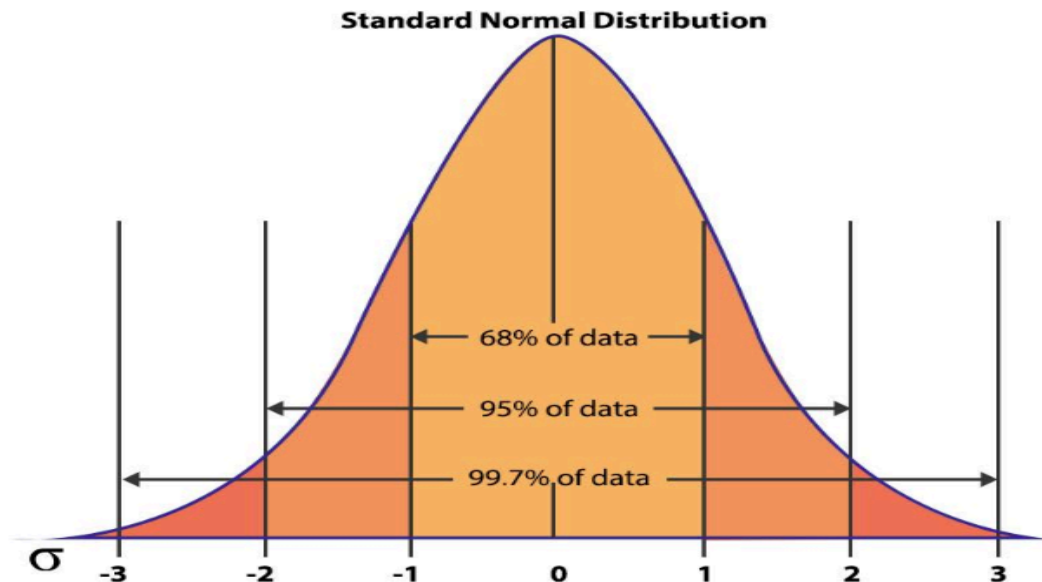
- 3.1 Concepts and Steps (null/alternative, α , p-value)
- 3.2 One-tailed vs. Two-tailed Tests

4. Statistical Tests

- 4.1 z-Test (means and proportions with known variance)
- 4.2 t-Test (one-sample, two-sample, paired samples)
- 4.3 ANOVA
 - One-way ANOVA
 - Two-way ANOVA
- 4.4 Chi-Square Tests
 - Goodness-of-Fit
 - Test of Independence



The Standard Normal (Z) Distribution



The Z-distribution is a normal distribution with a mean of 0 and a standard deviation of 1. It is the theoretical ideal used when we know the true population characteristics.

- **When to Use:** Only when the **population standard deviation (σ) is known** or when the **sample size (n) is very large** (due to the CLT).
- **The Z-Score:** The test statistic calculated from sample data is called the Z-score.

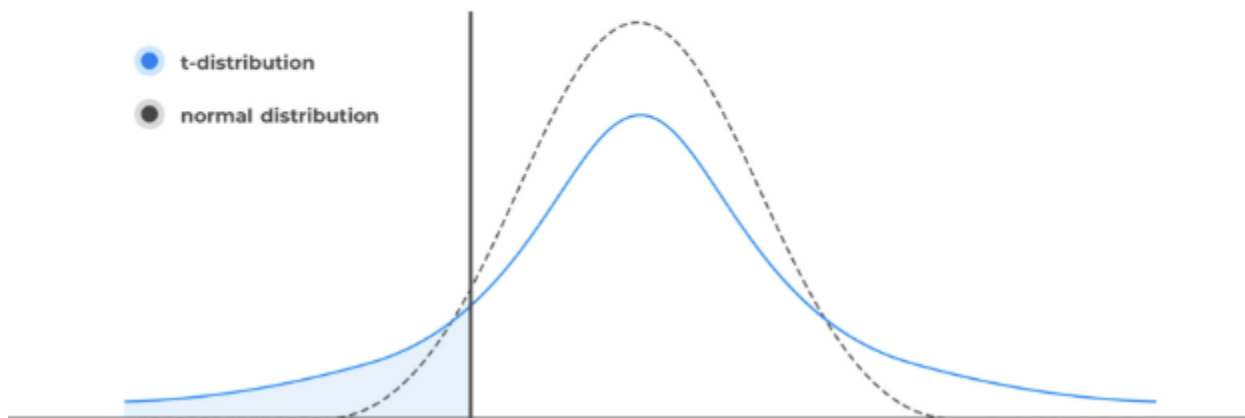
The t-Distribution

The **t-distribution** (or Student's t-distribution) is a family of bell-shaped curves that looks very similar to the Z-distribution but is slightly shorter and has fatter tails.

- **When to Use:** When the **sample size is small ($n < 30$)** AND the **population standard deviation (σ) is unknown**. Since the true σ is rarely known, the t-distribution is used far more often in practical data science than the Z-distribution.
- **Impact of Sample Size:** The shape of the t-distribution is determined by its **degrees of freedom (df)**, which is calculated as $df = n - 1$.
 - For **small samples** (low df), the t-distribution has much fatter tails to account for the higher uncertainty.

- As the **sample size grows** (high df), the t-distribution becomes virtually identical to the Standard Normal (Z) distribution.
- **Why Fatter Tails?** Using the sample standard deviation (s) to estimate the population standard deviation (σ) introduces more error, especially with a small sample. The fatter tails of the t-distribution assign a higher probability to extreme values, reflecting this increased uncertainty.

T-distribution vs Normal Distribution

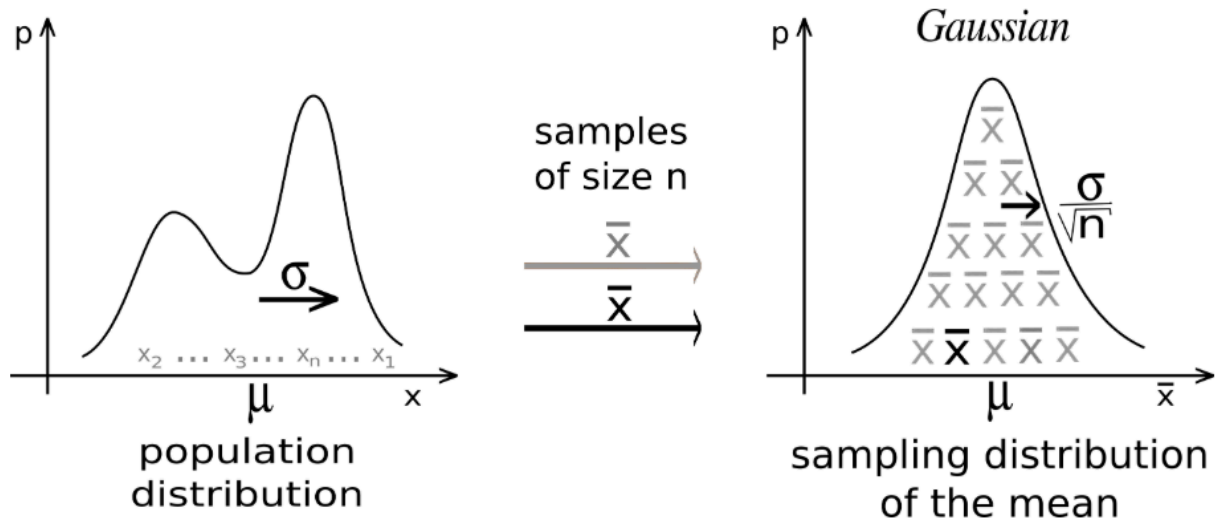


- **The t-Statistic:** The test statistic calculated from sample data is called the t-statistic

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

- The relationship between the sample mean and the population mean is what drives the inferential process, and the t-distribution is the essential tool for dealing with the real-world limitations of sample size and unknown population parameters.

The Central Limit Theorem (CLT)



The **Central Limit Theorem (CLT)** is the most profound and powerful concept in classical statistics. It's the mechanism that allows us to make reliable conclusions about an entire population based solely on a small sample of data.

1. What is the Central Limit Theorem?

The CLT is a statistical principle that describes the behavior of **sample means**. It states that if you draw many random samples of a sufficiently large size from **any** population (even a heavily skewed one), the distribution of those sample means will tend to form a **Normal (Gaussian) distribution**.

This resulting distribution is called the **Sampling Distribution of the Sample Mean**. The CLT guarantees this distribution will be normal, regardless of the original shape of the population data itself.

2. Central Limit Theorem Formula

Sample mean = Population mean = μ

$$\begin{aligned}\text{Sample standard deviation} &= \frac{(\text{Standard deviation})}{\sqrt{n}} \\ &= \frac{\sigma}{\sqrt{n}}\end{aligned}$$

The CLT provides the parameters for this new, normal sampling distribution:

- Mean of the Sampling Distribution ($\mu_{\bar{x}}$): The mean of the sample means is equal to the true mean of the population (μ).
 $\mu_{\bar{x}} = \mu$
- Standard Deviation of the Sampling Distribution (Standard Error, $\sigma_{\bar{x}}$): The standard deviation of the sampling distribution is called the Standard Error of the Mean. It

quantifies the variability of the sample means around the true population mean.

$$\sigma_{\bar{x}} = \sigma / \sqrt{n}$$

where σ is the population standard deviation and n is the sample size.

3. Sample Size and the Central Limit Theorem

The speed and accuracy with which the sampling distribution approaches normality depend entirely on the **sample size (n)**.

- **The Rule of Thumb ($n \geq 30$):** In practice, statisticians generally agree that if the sample size n is **30 or greater**, the sampling distribution will be close enough to a normal distribution to reliably use Z-tests and t-tests.
- **Impact of n :**
 - **Small n :** The sampling distribution will still resemble the original population distribution (e.g., if the population is skewed, the sampling distribution will be slightly skewed).
 - **Large n :** The sampling distribution quickly becomes perfectly symmetrical and bell-shaped, allowing for reliable inference.
 - **Decreasing Spread:** As n increases, the standard error (σ / \sqrt{n}) decreases. This means the sample means cluster more tightly around the true population mean, leading to more precise estimates.

4. Conditions of the Central Limit Theorem

To apply the CLT correctly, three conditions must be met:

1. **Random Samples:** The data must be collected using a random sampling method.
2. **Independence:** The individual data points must be independent of one another. For sampling without replacement, this requires the sample size to be small relative to the population size (e.g., sample size $n \leq 10\%$ of the population size N).
3. **Sample Size:** The sample size must be sufficiently large ($n \geq 30$ is generally the safe minimum).

5. Importance of the Central Limit Theorem

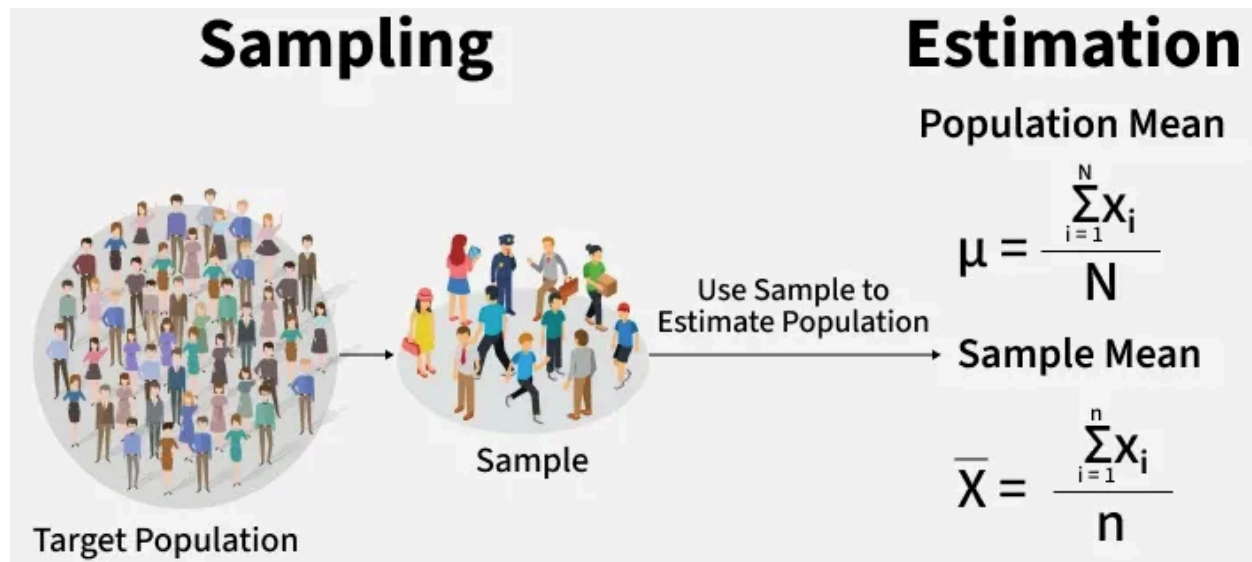
The CLT is the bedrock of inferential statistics for several key reasons:

- **Universality:** It works for virtually any population distribution. This means you don't need to know the shape of the original data to perform inference on its mean.
- **Foundation for Tests:** It enables the use of the **Normal distribution** as a model for the sample mean. Without this, we could not construct reliable confidence intervals or perform Z-tests and t-tests.
- **Quantifying Uncertainty:** It provides the formula for the **Standard Error**, which is the fundamental measure of uncertainty in estimation.

6. Central Limit Theorem Examples

- **E-commerce Conversion Rates:** A company wants to know the average conversion rate of a new ad campaign. The individual user outcomes (convert/don't convert) follow a Bernoulli distribution (highly non-normal). However, if the company takes many samples of 100 users, the distribution of the average conversion rates across those samples will be normal.
- **Customer Service Wait Times (Skewed Data):** Wait times are often heavily right-skewed (most waits are short, a few are very long). If you take random samples of 50 wait times and calculate the mean for each, the distribution of those 50-sample means will be approximately normal. This allows the company to set service level agreements using normal distribution statistics.

Statistical Estimation



Statistical estimation is the process of using sample data to guess or estimate an unknown characteristic (a **parameter**) of the larger population. We do this using two primary methods: the point estimate and the confidence interval.

1. Point Estimates

A **point estimate** is a single, best guess for a population parameter. It is the numerical value of a sample statistic that you use to estimate the corresponding population parameter.

- **Definition:** A single value (a sample statistic) used as the "best guess" for a population parameter.
- **Examples:**

- The **Sample Mean (\bar{x})** is the point estimate for the **Population Mean (μ)**.
- The **Sample Standard Deviation (s)** is the point estimate for the **Population Standard Deviation (σ)**.
- The **Sample Proportion (p^\wedge)** is the point estimate for the **Population Proportion (p)**.
- **Unbiased Estimators:** An estimator is considered **unbiased** if the mean of its sampling distribution is equal to the true value of the parameter being estimated. For instance, the sample mean (\bar{x}) is an unbiased estimator of the population mean (μ). This means that if you were to take many samples, the average of all those sample means would perfectly match the true population mean.

2. Confidence Intervals (CIs)

While a point estimate is the best single guess, it offers no information about how *good* that guess is. A **Confidence Interval (CI)** solves this by providing a range of values within which the true population parameter is likely to fall.

- **Definition and Purpose:** A CI is a range constructed around a point estimate that is likely to contain the true value of the population parameter with a certain level of confidence.
- Formula (General Form):

$$CI = \text{Point Estimate} \pm (\text{Margin of Error})$$
 The Margin of Error is calculated as the critical value (from the Z or t distribution) multiplied by the Standard Error of the estimate.

Interpretation of the Confidence Level

The confidence level (e.g., 90%, 95%, 99%) defines the reliability of the estimation process.

- **Example: 95% Confidence Interval:** If you were to repeat the sampling and calculation process many times, 95% of the confidence intervals constructed would contain the true population parameter.
- **Common Misinterpretation (Crucial Note):** A 95% CI does **NOT** mean there is a 95% probability that the true population parameter falls within the calculated interval. Once the interval is calculated, the parameter is either in it or it is not. The 95% confidence refers to the *methodology* used to construct the interval.

Factors Affecting the Interval Width

The width of the confidence interval is crucial because it indicates the precision of your estimate. A narrower interval is more precise. Three factors determine the width:

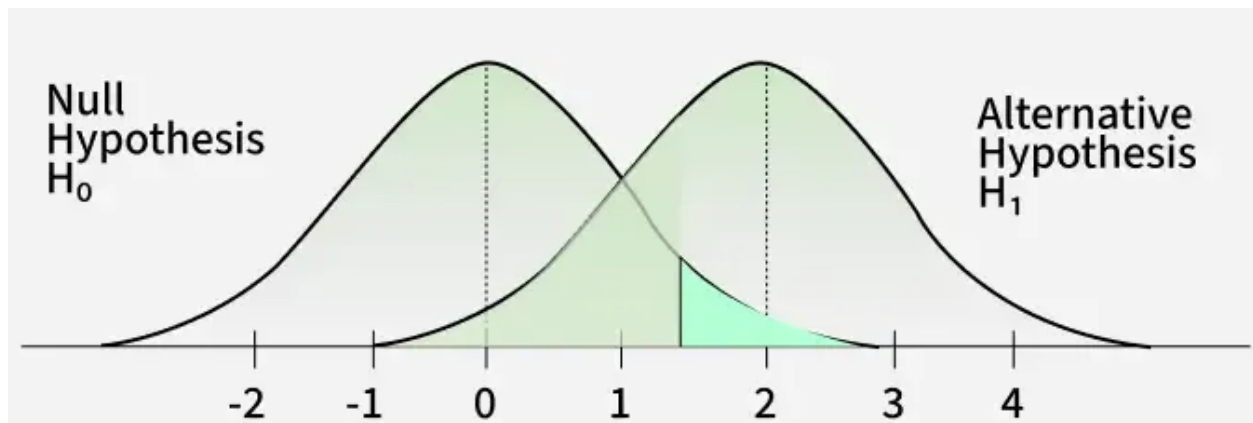
Factor	Change in Factor	Effect on Interval Width	Reason
Confidence Level	Increase (e.g., 90% → 99%)	Wider	To be more certain that you capture the true parameter, you must widen the net.
Sample Size (n)	Increase (e.g., 50 → 500)	Narrower	A larger sample reduces the Standard Error (σ/\sqrt{n}), increasing precision.
Variability (σ or s)	Increase (Higher Standard Deviation)	Wider	If the data is more spread out, the estimate is less precise, requiring a wider interval.

Confidence intervals are a cornerstone of inferential statistics because they allow data scientists to communicate not just their best guess, but also the inherent uncertainty surrounding that guess.

Hypothesis Testing

Hypothesis testing is a formal procedure for determining whether a claim about a population parameter is supported by the evidence in your sample data. It is the core mechanism for drawing actionable, statistical conclusions.

In simpler words, hypothesis testing compares two opposite ideas about a group of people or things and uses data from a small part of that group (a sample) to decide which idea is more likely true. We collect and study the sample data to check if the claim is correct.



Key Terms of Hypothesis Testing

Null Hypothesis (H_0) and Alternative Hypothesis (H_1)

The first step in any hypothesis test is setting up two opposing statements about the population. These hypotheses are mutually exclusive; only one can be true.

1. Null Hypothesis (H_0): The Status Quo

- **Definition:** The statement of no effect, no difference, or no relationship. It is the position assumed to be true until proven otherwise. It always contains a statement of equality (e.g., $=$, \leq , or \geq).
- **Goal:** The test is designed to determine if there is enough evidence in the sample data to **reject** H_0 .
- **Example:** A new drug has no effect on a patient's recovery time. ($H_0: \mu_{\text{new}} = \mu_{\text{old}}$)

2. Alternative Hypothesis (H_1 or H_a): The Claim

- **Definition:** The statement that contradicts the null hypothesis. It is the claim you are often trying to find evidence for (the effect you hope to see). It never contains a statement of equality.
- **Goal:** You hope the evidence allows you to **support** H_1 by rejecting H_0 .
- **Example:** The new drug reduces recovery time. ($H_1: \mu_{\text{new}} < \mu_{\text{old}}$)

Here is one simple, clear example of how to formulate the null and alternative hypotheses, based on a common business scenario.

Example: Testing a Website Change

A company is redesigning its checkout button and wants to know if the new design improves the **conversion rate** (the percentage of users who make a purchase). The historical average conversion rate is **10%**.

Hypothesis	Notation	Explanation
Null Hypothesis (H0)	$H_0:p=0.10$	The new button design has no effect . The true conversion rate remains equal to the historical rate of 10%.
Alternative Hypothesis (H1)	$H_1:p>0.10$	The new button design is better . The true conversion rate is greater than the historical rate of 10%.

I have prepared detailed notes for the introductory part of your **Hypothesis Testing** section. This content is crucial for defining the language and framework used to make statistical decisions.

The P-value

The **P-value** (Probability Value) is the central piece of evidence that the hypothesis test provides.

- **Definition:** The P-value is the probability of observing sample data (or data even more extreme than what you observed) *if the null hypothesis (H0) were actually true*.
- **The Decision Rule:**
 - **If P-value is small (P-value $\leq \alpha$):** The observed data is highly unlikely under H0. You **reject the null hypothesis**.
 - **If P-value is large (P-value $> \alpha$):** The observed data is reasonably likely under H0. You **fail to reject the null hypothesis**.
- **The Key Mistake to Avoid:** The P-value is **NOT** the probability that the null hypothesis is true. It only measures the evidence *against* H0 based on the collected sample.

Significance Level (α)

The **Significance Level (α)** is your predefined threshold for making a decision. You set this value *before* conducting the test.

- **Definition:** α is the maximum acceptable risk of making a **Type I Error** (rejecting H0 when it is actually true).
- **Common Values:** The most common significance level used is **$\alpha=0.05$** (or 5%).
- **Interpretation:**

- If you set $\alpha=0.05$, you are willing to accept a 5% chance of incorrectly rejecting the null hypothesis.
- If your P-value is 0.03 and α is 0.05, you reject H_0 . The chance of being wrong is 3%, which is less than your accepted risk threshold of 5%.

Test Statistic

The **Test Statistic** is a numerical value calculated from your sample data that is used to measure how far your sample results deviate from the assumption made by the Null Hypothesis (H_0).

- **Definition:** A single number that summarizes the relationship between your sample data and the null hypothesis.
- **Purpose:** It transforms your sample data into a standardized value (like a Z-score or t-score) that can be compared against known probability distributions (the Z or t-distribution).
- **Examples:**
 - **Z-statistic:** Used when the population standard deviation (σ) is known or the sample size is very large.
 - **t-statistic:** Used when the population standard deviation (σ) is unknown (the most common scenario).

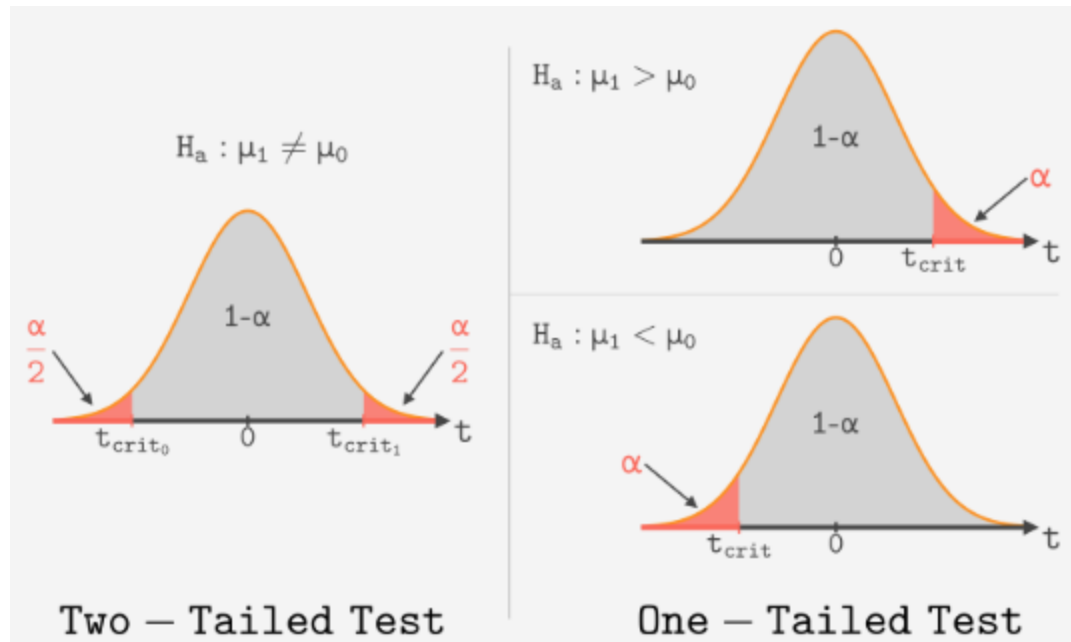
Critical Value

The **Critical Value** is the benchmark or cutoff point used to make the final decision in a hypothesis test. It defines the "region of rejection" on the distribution curve.

- **Definition:** A specific value from the appropriate probability distribution (Z or t) that corresponds to your chosen **Significance Level (α)**.
- **Purpose:** If your calculated **Test Statistic** falls outside this critical value (in the shaded tail of the distribution), it is considered too extreme to have occurred by chance under H_0 , and you **reject the null hypothesis**.

Types of Hypothesis Tests: One-Tailed vs. Two-Tailed

The choice between a one-tailed and two-tailed test is determined by how you formulate your **Alternative Hypothesis (H_1)** and, consequently, your research question. It dictates where the **rejection region** lies on the distribution curve.



1. One-Tailed Test (Directional Test)

A one-tailed test is used when you have a strong prior belief or theoretical reason to expect that the result will deviate from the null hypothesis in **only one specific direction** (either higher or lower, but not both).

Feature	Description	Rejection Region
Research Goal	To confirm a directional change (e.g., <i>increase or decrease</i>).	Located entirely in one tail of the distribution.
Example	Testing if a new drug improves patient recovery time (i.e., time decreases).	Left Tail (looking for a value significantly <i>lower</i> than the null mean).

Types of One-Tailed Tests

Type	Null Hypothesis (H_0)	Alternative Hypothesis (H_1)	Goal / Conclusion
Left-Tailed	$H_0: \mu \geq 50$	$H_1: \mu < 50$	Reject H_0 only if the sample mean is significantly small (i.e., falls in the left tail).
Right-Tailed	$H_0: \mu \leq 50$	$H_1: \mu > 50$	Reject H_0 only if the sample mean is significantly large (i.e., falls in the right tail).

2. Two-Tailed Test (Non-Directional Test)

A two-tailed test is used when you want to see if there is **any difference at all**—meaning the result could be significantly higher **OR** significantly lower than the null hypothesis value.

Feature	Description		Rejection Region
Research Goal	To confirm a difference exists, without predicting the direction.		Split evenly between both the left and right tails of the distribution.
Example	Testing if a marketing strategy affects sales (i.e., sales could go up or down).		Both Tails (looking for a difference that is extreme in either direction).

Type	Null Hypothesis (H_0)	Alternative Hypothesis (H_1)	Goal / Conclusion
Two-Tailed	$H_0 : \mu = 50$	$H_1 : \mu \neq 50$	Reject H_0 if the sample mean is significantly different (either much smaller or much larger) than the hypothesized mean.

Key Distinction: Rejection Region

The primary practical difference lies in how the significance level (α) is distributed:

- **One-Tailed:** The entire α (e.g., 5%) is placed into a single tail. This makes it **easier to reject H_0** in that one direction.
- **Two-Tailed:** The α is split in half ($\alpha/2$) and placed into each tail (e.g., 2.5% in the left tail and 2.5% in the right tail). This requires a more extreme result to reject H_0 in any one direction.

How Hypothesis Testing Works: A Step-by-Step Guide

Hypothesis testing is a formal, six-step procedure used to determine if a claim about a population parameter can be supported by the evidence gathered from a sample.

Step 1: Formulate the Hypothesis (H_0 and H_1)

The process begins by setting up two opposing statements about the population parameter you are investigating.

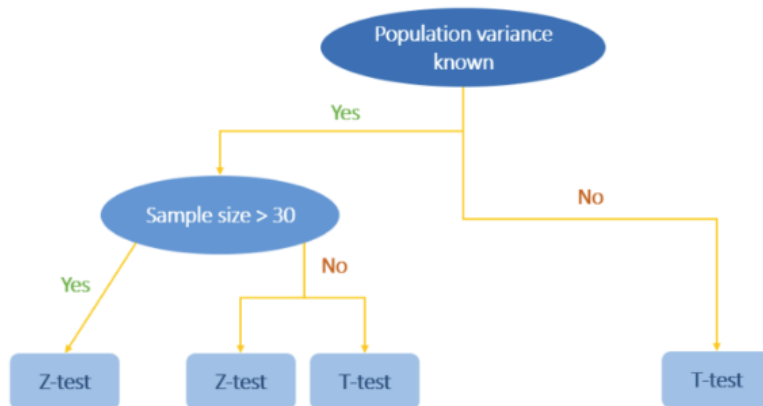
- **Null Hypothesis (H_0):** The statement of **no effect, no difference, or no change**—the status quo. It always assumes equality (e.g., $\mu=50$). This is the hypothesis we assume is true until proven otherwise.
- **Alternative Hypothesis (H_1):** The statement that contradicts H_0 . This is often the claim or effect you are trying to find evidence for (e.g., $\mu \neq 50$, $\mu > 50$, or $\mu < 50$).
 - **Example:** Testing a new algorithm. H_0 : User engagement remains the same ($\mu_{\text{new}} = \mu_{\text{old}}$). H_1 : User engagement has improved ($\mu_{\text{new}} > \mu_{\text{old}}$).

Step 2: Set the Significance Level (α)

The significance level, denoted as α (alpha), is chosen *before* data collection and analysis.

- **Definition:** α is the maximum risk we are willing to accept of making a **Type I Error** (wrongly rejecting a true H_0).
- **Threshold:** It is typically set to **0.05 (or 5%)**. This means we are comfortable with a 5% chance that we conclude the new algorithm works when, in reality, it doesn't.

Step 3: Select the Test and Collect Data



Based on the nature of the data and the research question, the appropriate statistical test is selected.

- **Test Selection:** The choice depends on the variable type, the number of groups, and whether the population variance is known (e.g., **Z-test/T-test** for comparing two means, **Chi-square** for categorical relationships).
- **Data Collection:** Relevant data is gathered (e.g., collecting user engagement metrics both before and after the new algorithm's implementation).

Step 4: Calculate the Test Statistic

The test statistic quantifies the difference between what we observed in the sample data and what we expected to see if the null hypothesis (H_0) were true.

- **Purpose:** To measure the magnitude of the deviation from the null assumption in standardized units.
- **Calculation:** The appropriate formula is applied to the sample data to generate a single number (a Z-score, t-score, or χ^2 value).

T-statistic Example: $t = \frac{\bar{x} - \mu}{s / \sqrt{n}}$

Step 5: Make a Decision (P-value or Critical Value)

The test statistic is used to find either the P-value or the Critical Value, which drives the final decision on the null hypothesis.

Decision Method A: Using the P-value

The **P-value** is the probability of observing the current test statistic (or a more extreme one) if H_0 were true.

- **Decision Rule:**
 - If **P-value $\leq \alpha$** (P-value is small), the data is highly unlikely under H_0 . → **Reject H_0** (The result is statistically significant).
 - If **P-value $> \alpha$** (P-value is large), the data is plausible under H_0 . → **Fail to Reject H_0** .
- **Example:** If P-value is 0.03 and α is 0.05, you reject H_0 because $0.03 < 0.05$.

Decision Method B: Using the Critical Value

The **Critical Value** is the cutoff point on the distribution curve corresponding to α .

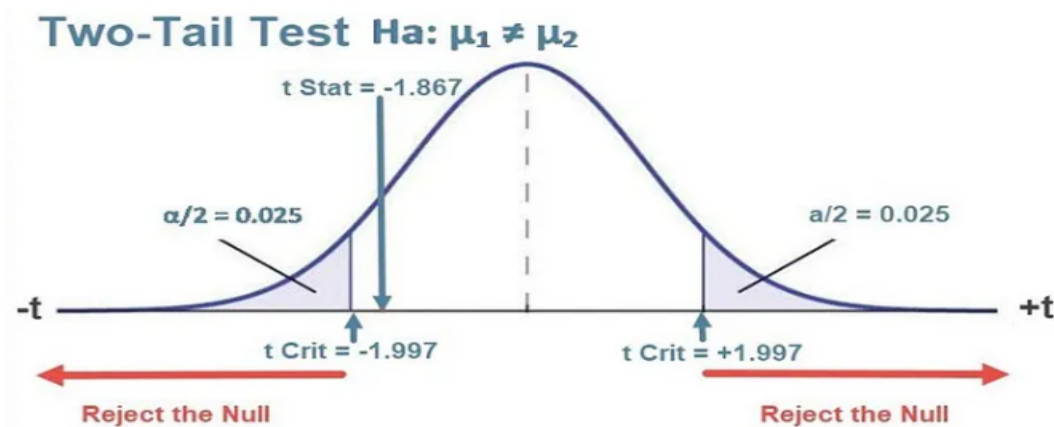
- **Decision Rule:**
 - If the **Test Statistic falls in the Rejection Region** (i.e., is more extreme than the critical value), → **Reject H_0** .

Step 6: Interpret and Conclude

The final step is to translate the statistical decision back into a clear, non-technical conclusion about the original research question.

- **If H_0 is Rejected:** There is sufficient statistical evidence to support the **Alternative Hypothesis (H_1)**. (e.g., "The new algorithm significantly improves user engagement.")
- **If H_0 is Not Rejected:** There is **not enough statistical evidence** to support the Alternative Hypothesis (H_1). (e.g., "We cannot conclude that the new algorithm improves user engagement.")

The T-test: Comparing Two Means



The **T-test** is one of the most widely used inferential tools in data analysis. Its primary purpose is to determine if the means (averages) of two groups are significantly different from each other, or if a single sample mean is different from a known target value.

We use the T-test within the framework of **hypothesis testing**, where we start with the **Null Hypothesis (H0)**: that the group means are equal. The T-test calculates a standardized measure (the T-statistic) that helps us decide whether there is enough evidence to reject this assumption.

Key Assumptions for T-tests

For the results of a T-test to be reliable, your data should meet the following assumptions (though the T-test is relatively robust to minor violations, especially with larger sample sizes):

1. **Independence:** Observations within each group must be independent of one another (e.g., one student's score does not influence another's).
2. **Normality:** The data within each group should be **approximately normally distributed**. This assumption is less critical when the sample size (n) is large ($n \geq 30$) due to the Central Limit Theorem.
3. **Homogeneity of Variances:** When comparing two independent groups, the **variances** (spread) of the two groups should be approximately equal. This ensures that the groups have a similar level of variability.
4. **Absence of Outliers:** Outliers can significantly skew the mean and inflate the standard deviation, strongly influencing the T-test results, especially with small samples.

Types of T-tests

The T-test is divided into two distinct types, based on how the samples are structured.

One-Sample T-Test

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

\bar{X} = observed mean of the sample
 μ = assumed mean
 s = standard deviation
 n = sample size

Two-Sample T-Test

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

\bar{X}_1 = observed mean of 1st sample
 \bar{X}_2 = observed mean of 2nd sample
 s_1 = standard deviation of 1st sample
 s_2 = standard deviation of 2nd sample
 n_1 = sample size of 1st sample
 n_2 = sample size of 2nd sample

One Sample T-Test

The value against which we are comparing is a single value, ie we compare the mean of a sample with a single value to check how much the mean deviates from that single value

Two Sample T-Test

We compare the means and variances of two samples, and we assess how much they differ. A smaller p-value — the probability of observing an extreme value (such as the mean of sample 1) assuming our hypothesis is false — supports the validity of our hypothesis. In other words, a low p-value suggests that our observed data aligns with our hypothesis.

Example Walkthrough: One-Sample T-test

This example demonstrates how to test a sample against a specific target value.

Problem: We want to see if a weight loss camp was effective. We have the weights of 25 participants *after* the camp and want to compare the average post-camp weight (\bar{x} =75 kg)

against the known, historical *target* weight for this population ($\mu=45$ kg).

Parameter	Value	Interpretation
Population Mean (μ)	45 kg	The hypothesized weight <i>before</i> the camp.
Sample Mean (\bar{x})	75 kg	The average weight <i>after</i> the camp.
Sample Size (n)	25	The number of participants.
Sample Std. Dev. (s)	25	The spread of weights after the camp.

- **Hypotheses:**
 - $H_0: \mu=45$ (The camp had no effect, and the average weight is still 45 kg).
 - $H_1: \mu \neq 45$ (The camp *did* have an effect, and the average weight is different from 45 kg).
- **Result (Calculated T-statistic):**

$$t = \frac{75 - 45}{25 / \sqrt{25}} = \frac{30}{5} = 6.0$$

Interpretation: The calculated T-statistic of **6.0** is much larger than the critical value (which is ≈ 2.06 for $\alpha=0.05$ and $df=24$). Since $6.0 > 2.06$, we **reject the Null Hypothesis**. The T-test confirms that the observed change in participants' weights is **statistically significant** and highly unlikely due to random chance.

The Z-test and T-test: Comparing Means

Both the Z-test and T-test are used to compare means within the framework of hypothesis testing. The choice between them depends entirely on **how much you know about the population** and the **size of your sample**.

The most crucial factor is whether the **population standard deviation (σ) is known**.

1. When to Use Which Test

Test	Key Condition	Sample Size (n)	Logic
Z-test	Population Standard Deviation (σ) is KNOWN.	$n > 30$ (Large)	Since σ is known, we use the highly predictable Z-distribution .
T-test	Population Standard Deviation (σ) is UNKNOWN.	$n < 30$ (Small) or $n \geq 30$	Since we must <i>estimate</i> σ with the sample s , we use the t-distribution to account for extra uncertainty.
Practical Reality	The T-test is used far more often in data science because the true population standard deviation (σ) is rarely known.		

2. The Formulas and Test Statistics

Both formulas measure how many standard errors the sample mean (\bar{x}) is away from the hypothesized mean (μ).

Test	Test Statistic Formula	Key Difference
Z-statistic	$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$	Uses the Population Standard Deviation (σ) .
T-statistic	$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$	Uses the Sample Standard Deviation (s) to estimate σ .

3. The T-test: Types and Applications

The T-test is used when the population standard deviation is unknown and is categorized by the relationship between the samples.

T-Test Type	Purpose	Example Application
One-Sample T-test	Compares a single sample mean (\bar{x}) against a known target value (μ).	Does the average wait time in our checkout line (\bar{x}) exceed the company's 5-minute service goal ($\mu = 5$)?
Independent (Two-Sample) T-test	Compares the means of two entirely separate (independent) groups .	Is the average score of students who took the offline class different from those who took the online class ?
Paired T-test	Compares the means of the same group measured twice (dependent samples). Measures the effect of an intervention.	Did a new medication decrease a patient's cholesterol level before versus after treatment?

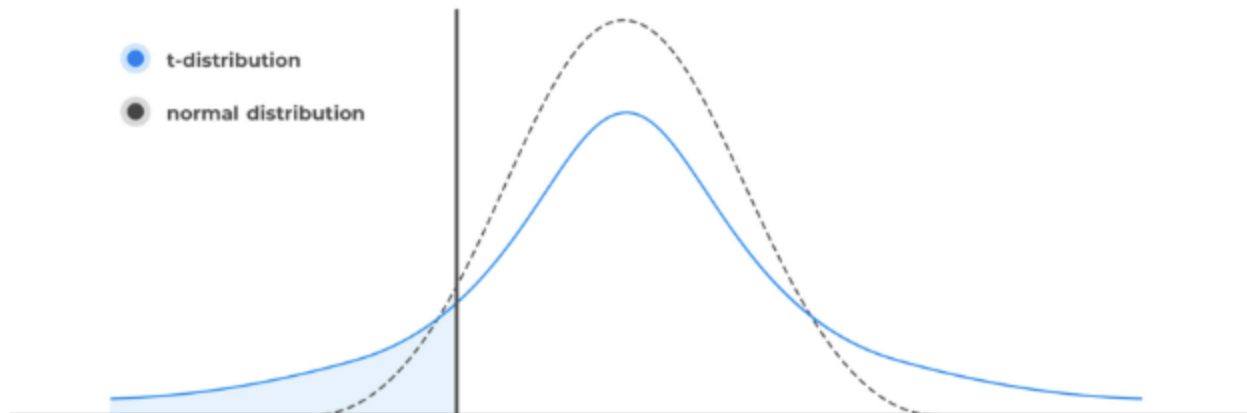
4. Key Assumptions for Both Tests

For the results of any mean-comparison test to be reliable, the data should meet these assumptions:

1. **Independence:** All observations must be independent of one another.
2. **Normality:** The data within each group should be **approximately normally distributed**. (Less critical for $n \geq 30$ due to the Central Limit Theorem).
3. **Homogeneity of Variances (for Independent Tests):** When running an independent two-sample test, the spread (variance) of the two groups should be roughly equal.

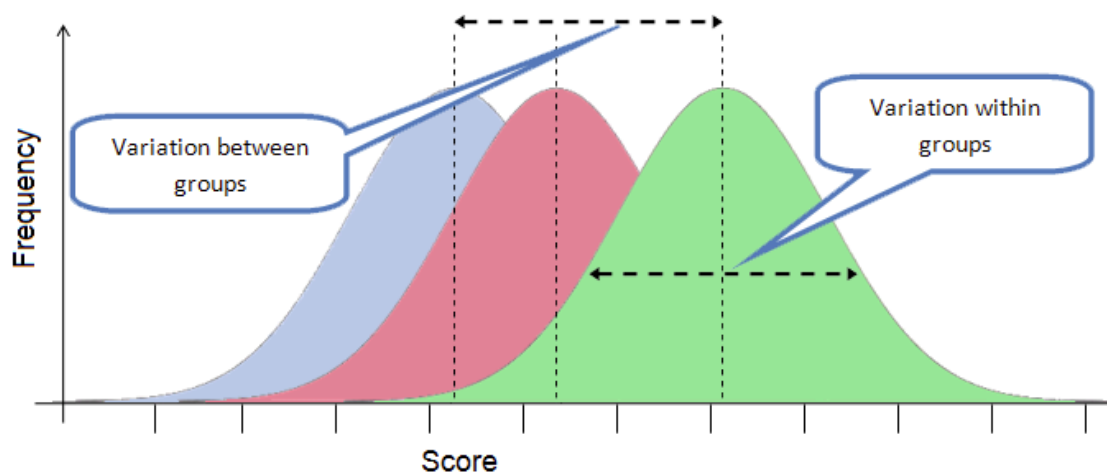
5. Summary: Z vs. T Distribution

The difference between the Z and T distributions highlights the impact of uncertainty.



- The **Z-distribution** is the theoretical ideal, used only when we have near-perfect information (σ is known).
- The **t-distribution** has slightly **fatter tails** than the Z-distribution. These fatter tails assign a higher probability to extreme values, reflecting the increased **uncertainty** introduced when we have to estimate the population standard deviation (σ) using only the sample standard deviation (s). As the sample size grows, the T-distribution sheds its fat tails and becomes identical to the Z-distribution.

ANOVA (Analysis of Variance)



ANOVA is a powerful parametric statistical technique used to determine if there are any statistically significant differences between the means of **three or more independent groups**.

The test is named "Analysis of Variance" because it works by examining the **variability** (variance) in the data to make inferences about the **means**.

- **Hypotheses:**
 - **Null Hypothesis (H0):** All group means are equal. ($\mu_1=\mu_2=\mu_3=\dots$)
 - **Alternative Hypothesis (H1):** At least one group mean is different from the others.

One-Way vs. Two-Way ANOVA

The type of ANOVA used depends on the number of categorical independent variables (factors) being tested.

Type	Number of Factors	Purpose	Example
One-Way ANOVA	One independent variable with ≥ 3 levels/groups.	Tests the impact of one factor on the dependent variable.	Does teaching Method A, B, or C (one factor) affect student scores?
Two-Way ANOVA	Two independent variables (factors).	Tests the main effects of each factor <i>and</i> the interaction effect between them.	Do Teaching Method (Factor 1) AND Class Size (Factor 2) affect student scores?

The Logic of the F-Statistic

The core of ANOVA is the **F-statistic** (or F-ratio), which quantifies how much the groups differ from each other relative to how much the individuals within each group vary.

$$F = \frac{\text{Variance Between Groups (Signal)}}{\text{Variance Within Groups (Noise)}}$$

- **Variance Between Groups (The Signal):** Measures the differences between the means of the various groups. A large difference here suggests a genuine effect caused by the factor (e.g., teaching method).
- **Variance Within Groups (The Noise):** Measures the natural variation among individual observations *inside* each group (due to random chance or individual differences).

Interpretation:

- **If F is large ($F \gg 1$):** The "signal" is much greater than the "noise." This suggests the factor has a significant effect, leading to the **rejection of H_0** .
 - **If F is small ($F \approx 1$):** The "signal" is similar to the "noise." The differences between the group means are likely due to random chance, leading to a **failure to reject H_0** .
-

Assumptions for ANOVA

As a parametric test, ANOVA requires several assumptions to be met for the results to be valid:

1. **Normality:** The dependent variable should be **approximately normally distributed** within each of the groups being compared.
 2. **Independence:** All samples must be selected randomly and be **independent** of one another.
 3. **Homogeneity of Variances (Homoscedasticity):** The variances (spread) of the dependent variable across all groups must be **approximately equal**. This is the **most critical** assumption.
 4. **Continuous Dependent Variable:** The variable you are measuring (e.g., exam score) must be continuous.
-

The Decision Rule

The final decision is made by comparing the calculated F-statistic to a critical value (F_{table}) or, more commonly in modern software, by examining the **P-value**.

Decision Method	Rule	Conclusion
F-Ratio	If $F_{calculated} > F_{table}$	Reject H_0: A significant difference exists between at least one pair of means.
P-value	If P-value $\leq \alpha$ (e.g., 0.05)	Reject H_0: The results are statistically significant.

The Chi-Square (χ²) Test

The **Chi-Square (χ²) Test** is a fundamental non-parametric hypothesis test used primarily for analyzing **categorical data** (data classified into groups or categories). It determines if the observed distribution of data differs significantly from what would be expected under the null hypothesis.

Key Concepts

Terminology	Description
Test Statistic (χ²)	The calculated value that measures the discrepancy between Observed Frequencies (O) and Expected Frequencies (E) .
Observed Frequency (O)	The actual count of data points in each category or cell, gathered from the sample.
Expected Frequency (E)	The theoretical count for each cell, calculated under the assumption that the null hypothesis (e.g., independence) is true.
Contingency Table	The table used to organize the raw counts of two categorical variables before the test is performed.
Degrees of Freedom (df)	Determines the shape of the χ² distribution. Calculated as $df = (\text{Rows} - 1) \times (\text{Columns} - 1)$ for a Test of Independence.

Formula for the Test Statistic:

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

Two Main Types of Chi-Square Tests

The χ² test is applied in two distinct ways, depending on the research question:

1. Chi-Square Test of Independence

- **Purpose:** To determine if there is a statistically significant **association** between **two categorical variables**.
- **Hypotheses:**
 - **Null (H0):** The two variables are **independent** (no relationship exists).
 - **Alternative (H1):** The two variables are **dependent** (a relationship exists).

- **Example:** Is a person's **Favorite Color** (Categorical Variable 1) independent of their **Preferred Ice Cream Flavor** (Categorical Variable 2)?
- **Calculation Focus:** The Expected Frequencies (E) are calculated based on the idea that the probability of a cell is the product of its row total probability and its column total probability.

2. Chi-Square Goodness-of-Fit Test

- **Purpose:** To determine if the **observed frequency distribution** of a single categorical variable matches a **hypothesized or known expected distribution**.
 - **Hypotheses:**
 - **Null (H0):** The observed distribution matches the expected distribution.
 - **Alternative (H1):** The observed distribution does not match the expected distribution.
 - **Example:** A casino expects a roulette wheel to land on black, red, and green in a 50:47:3 ratio. Does a sample of 100 spins fit this expected distribution?
-

Steps and Decision Logic

The Chi-Square test follows the standard hypothesis testing procedure:

1. **Formulate Hypotheses:** Define H0 (Independence/Match) and H1 (Dependence/Mismatch).
 2. **Organize Data:** Create a Contingency Table with **Observed Frequencies (O)**.
 3. **Calculate Expected Frequencies (E):** Determine what the cell counts should be if H0 were true.
 4. **Calculate χ^2 Statistic:** Apply the formula $\chi^2 = \sum \frac{(O-E)^2}{E}$.
 5. **Determine Degrees of Freedom (df):** Calculate $df = (r-1)(c-1)$.
 6. **Make a Decision:** Compare the calculated χ^2 value to the critical value or use the P-value.
 - **If P-value $\leq \alpha$ (e.g., 0.05):** The difference between O and E is too large to be random. **Reject H0.** (The variables are associated.)
 - **If P-value $> \alpha$: Fail to Reject H0.** (Insufficient evidence of an association.)
-

Assumptions and Data Science Applications

Key Assumptions

1. **Independence of Observations:** Each observation must be independent of all others.
2. **Categorical Data:** Both variables must be categorical (nominal or ordinal).

3. **Expected Frequencies:** Crucially, every cell in the contingency table should have an expected frequency (E) of at least 5. If not, the test's validity is compromised, and alternatives like Fisher's Exact Test should be considered.

Applications in Data Science

- **A/B Testing & Feature Evaluation:** Used to check if the observed clicks or conversions between Test Group A and Group B are statistically significant, rather than just random variation.
- **Machine Learning (Feature Selection):** The Chi-Square test is used during feature preprocessing to quickly assess whether a categorical feature is independent of the target variable. If H_0 (independence) is rejected, the feature is highly correlated with the target and is therefore a good candidate for the model.
- **NLP:** Can be used to evaluate word usage, testing if the frequency of certain words in one corpus differs from another.