

what is hypothesis testing
what are assumption of hypothesis testing
why we need the hypothesis testing
what tests are falls under the hypothesis testing

Below is a **clear, structured, deeply intuitive explanation** of **Hypothesis Testing**, including

- What it is
- Why we need it
- Its assumptions
- Types of hypothesis tests

I'll keep it ML-connected and conceptually clean.

1. What is Hypothesis Testing? (Intuition + Formal Meaning)

Intuition

Hypothesis testing is a **decision-making framework** that helps us answer:

"Is the pattern I observe in the sample real, or could it have happened just by random chance?"

Data is noisy. Even if nothing has changed, averages of samples will always fluctuate. Hypothesis testing separates **random fluctuations** from **real effects** by measuring how unlikely the observed data is **assuming no real effect exists**.

Formal Definition

Hypothesis testing is a **statistical procedure** to evaluate claims about a population using

☰  > ChatGPT 5.2 ▾

  ...

- **Null hypothesis (H_0):** Default claim (no effect, no difference, status quo)
- **Alternative hypothesis (H_1):** What you want to prove (effect exists, difference exists)

Then it uses a **test statistic** and a **p-value** to decide whether there is enough evidence to reject H_0 .

Example intuition

Suppose a new ML model claims to reduce error from **18% to 15%**.

Question:

Is 15% genuinely better, or is the improvement just due to random variation?

Hypothesis test:

- H_0 : error_new = error_old (no improvement)
- H_1 : error_new < error_old (model improved)

If the probability of seeing 15% error under the assumption that real error = 18% is extremely small \rightarrow we reject H_0 .

2. Why do we need Hypothesis Testing?

1. Data is random

Even if two algorithms have equal accuracy, a test run may show one "winning" accidentally.

Hypothesis testing prevents **false discoveries**.

2. The human brain is biased

We tend to interpret noise as patterns. Hypothesis testing forces an objective measurement:

"How rare is this result under chance?"

3. To validate ML models and A/B tests

Hypothesis testing is everywhere:

- ✓ A/B testing (web experiments)
- ✓ New drug approval

- ✓ Evaluating new ML algorithms
 - ✓ Quality control in manufacturing
 - ✓ Financial return comparisons
-

4. Helps generalize from sample → population

You cannot measure all users, all stocks, all pixels — only samples.

Hypothesis testing says:

“Is the sample evidence strong enough to generalize to the entire population?”

3. Assumptions of Hypothesis Testing

Different tests have different assumptions, but globally, hypothesis testing relies on:

General assumptions

1. Random sampling

Samples must be collected randomly (no systematic bias).

2. Independence

Each observation should not influence the others.

ML analogy: If your training set contains duplicates, independence is violated → invalid test results.

3. Distributional assumptions

Many tests assume:

- Data is normally distributed
- Variances are equal across groups
- Sample size is large enough (if using CLT)

Non-parametric tests remove these assumptions.

4. Correct model specification

The hypothesis test must match the situation:

- Means → t-test
 - Proportions → z-test
 - Variances → chi-square test
 - Distributions → KS-test
-

5. Sufficient sample size

Small samples often fail to detect real effects (low statistical power).

4. What tests fall under Hypothesis Testing?

We divide tests based on purpose and data type.

A. Tests for Comparing Means

1. One-sample t-test

Test sample mean vs population mean.

2. Two-sample t-test

Compare means of two groups (A/B test).

3. Paired t-test

Compare before/after on same subjects.

4. Z-test

Like a t-test but used when population variance is known (rare).

B. Tests for Comparing Proportions

1. One-proportion z-test

Example: Website conversion rate = 5%?

2. Two-proportion z-test

Compare two proportions (A/B conversion test).

C. Tests for Variances

1. Chi-square test for variance

Test if variance matches a claimed value.

2. F-test

Compare two variances.

D. Tests for Distributions

1. Kolmogorov–Smirnov (KS) test

Compare sample distribution with theoretical distribution.

2. Anderson–Darling test

More sensitive than KS.

3. Shapiro–Wilk test

Normality test.

E. Tests for Categorical Data (Counts)

1. Chi-square goodness-of-fit test

Check if categories match expected frequencies.

2. Chi-square test of independence

Check relationships in contingency tables (e.g., gender vs product preference).

3. Fisher's exact test

Used when samples are small.

F. Regression-Based Tests

1. t-test on regression coefficients

Is a coefficient significant?

2. F-test for overall regression

Does the regression model provide explanatory power?

G. Non-parametric Tests (No normality assumption)

1. Mann–Whitney U test

Alternative to two-sample t-test.

2. Wilcoxon signed-rank test

Alternative to paired t-test.

3. Kruskal–Wallis test

Alternative to one-way ANOVA.

4. Spearman rank correlation test

Tests monotonic relationships.

H. ANOVA Family (Comparing more than 2 means)

1. One-way ANOVA

Compare means of 3+ groups.

2. Two-way ANOVA

Two factors affecting mean.

3. Repeated-measures ANOVA

Multiple measurements on same subjects.

Summary Table

Goal	Parametric Test	Non-parametric Test
Compare 2 means	t-test	Mann–Whitney
Compare 2 proportions	z-test	Fisher's exact
Compare >2 means	ANOVA	Kruskal–Wallis
Compare distribution	KS test	—
Test correlation	Pearson test	Spearman test

If you want:

I can now continue with:

- ✓ Deriving p-value mathematically
- ✓ Showing real-world ML examples of hypothesis testing
- ✓ Visual intuition diagrams of sampling distributions
- ✓ Step-by-step numerical examples for each test
- ✓ A full flowchart: *"Which hypothesis test should I use?"*

Just tell me next.



go for all but one by one

- ✓ Deriving p-value mathematically
- ✓ Showing real-world ML examples of hypothesis testing
- ✓ Visual intuition diagrams of sampling distributions
- ✓ Step-by-step numerical examples for each test
- ✓ A full flowchart: *"Which hypothesis test should I use?"*

Great — we'll go through **all items, one by one**, with deep intuition, diagrams, math, and ML connections.

Here is the **order** we will follow:

Part 1 → Deriving the p-value mathematically

(we will do this now)

Then:

Part 2 → Real-world ML examples of hypothesis testing

Part 3 → Visual intuition diagrams of sampling distributions

Part 4 → Step-by-step numerical examples for major tests

Part 5 → Full flowchart: "Which hypothesis test should I use?"

PART 1 — Deriving the p-value (Mathematical + Intuition)

1. What is the p-value? (Intuition first)

A p-value measures how surprising your data is if the null hypothesis were true.

"If H_0 is true, how likely am I to get a result this extreme or more extreme?"

- Small p-value → data is **too unlikely** under H_0 → reject H_0
- Large p-value → data **looks normal** under H_0 → fail to reject H_0

2. The General Mathematical Derivation

Assume:

- H_0 defines a distribution for the test statistic:

$$T \sim f(t \mid H_0)$$

- You compute the observed value:

$$t_{obs}$$

The **p-value** is:

Two-tailed

$$p = P(|T| \geq |t_{obs}| \mid H_0)$$

Right-tailed

$$p = P(T \geq t_{obs} \mid H_0)$$

Left-tailed

$$p = P(T \leq t_{obs} \mid H_0)$$

This is ALWAYS computed using the **null distribution**.

3. Deriving p-value for a t-test (the most common case)

Suppose:

- Sample mean = \bar{x}
- Population mean under $H_0 = \mu_0$
- Standard deviation = s
- Sample size = n

Test statistic:

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

Under H_0 (if population is normal), we know:

$$T \sim t_{(n-1)}$$

Deriving the p-value (two-tailed)

$$p = 2 \cdot P(T \geq |t_{obs}|)$$

Using the t-distribution CDF:

$$p = 2 \left[1 - F_{t_{(n-1)}}(|t_{obs}|) \right]$$

This is the general formula.

4. Deriving p-value for a z-test

For large samples or known σ :

$$Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

Under H_0 :

$$Z \sim N(0, 1)$$

p-value (two-tailed):

$$p = 2 \cdot (1 - \Phi(|z|))$$

Where Φ is the standard normal cumulative distribution.

5. Deriving p-value for chi-square test (variance test)

Test statistic:

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

Under H_0 :

$$\chi^2 \sim \chi_{n-1}^2$$

p-value (right-tail):

$$p = P(\chi^2 \geq \chi_{obs}^2)$$

6. Deriving p-value for F-test

Test statistic:

$$F = \frac{s_1^2}{s_2^2}$$

Under H_0 :

$$F \sim F_{(df_1, df_2)}$$

p-value:

$$p = P(F \geq F_{obs})$$

7. Deriving p-value using the sampling distribution idea (MOST IMPORTANT)

This is the core idea behind all tests:

Step 1 → Assume H_0 is true

(e.g., mean = 10)

Step 2 → Compute statistic for your sample

(e.g., $t = 2.3$)

Step 3 → Check where your statistic lies in the H_0 distribution

If $t = 2.3$ is very far out in the tails → rare under H_0 → reject.

8. Simple numerical derivation example

Let's compute a real p-value to cement intuition.

Suppose:

- Sample mean = 15
- H_0 mean = 18
- $s = 4$
- $n = 25$

Compute t :

$$t = \frac{15 - 18}{4/\sqrt{25}} = \frac{-3}{0.8} = -3.75$$

Degrees of freedom = 24.

Now compute p-value (two-tailed):

$$p = 2 \cdot P(T \leq -3.75)$$

Using t-table/approximation:

$$P(T \leq -3.75) \approx 0.00048$$

So:

$$p \approx 2 \times 0.00048 = 0.00096$$

Interpretation:

Less than 0.1% chance of seeing mean = 15 if true mean was 18 → reject H_0 .

9. p-value is NOT:

- ✗ Probability that H_0 is true
 - ✗ Probability that H_1 is true
 - ✗ Confidence level
 - ✗ Effect size
-

10. ML connection: Understanding p-value in model evaluation

Example: Comparing two ML models using hypothesis testing

Suppose you run 10-fold CV:

Model A accuracy: 92.0%

Model B accuracy: 93.2%

Difference = 1.2% → but is this **real** or noise?

Run a *paired t-test*:

- Each fold produces a pair of accuracies
- Compute differences
- Compute t-statistic
- Compute p-value

If $p < 0.05$ → difference is statistically significant.

If not → both models are statistically equal.

NEXT STEP

Before proceeding to the next part, here's what comes next:

PART 2 — Real-World ML Examples of Hypothesis Testing

(A/B tests, model comparisons, feature selection significance tests, drift detection)

Tell me “**next**”, and I'll start Part 2.

[next](#)

PART 2 — Real-World ML Examples of Hypothesis Testing

(Deep intuition + formulas + where exactly each test is used in ML pipelines)

We will go through 5 major ML areas where hypothesis testing appears constantly:

1. A/B Testing for Model or Product Decisions
2. Comparing ML Models (CV-based t-tests)
3. Feature Selection using Statistical Significance
4. Detecting Data Drift / Concept Drift
5. Checking Assumptions in Models (normality, variance, independence)

Each section includes a real scenario + which test is used + why + interpretation.

1. A/B TESTING in ML / Data Science

Example: you deploy a new ranking algorithm for a recommendation system.

- Variant A → Old model
- Variant B → New model

Suppose you measure CTR (click-through rate):

Version	CTR	Sample Size
A	5.0%	50,000
B	5.4%	50,000

Question:

Is the increase from 5.0% → 5.4% a real improvement or just noise?

Test used: Two-proportion z-test

H_0 : CTR_A = CTR_B

H_1 : CTR_B > CTR_A

Test statistic:

$$z = \frac{p_B - p_A}{\sqrt{p(1-p) \left(\frac{1}{n_A} + \frac{1}{n_B} \right)}}$$

Evaluate → get p-value.

If $p < 0.05$, new model wins.

ML Insight

Even a small difference (0.3%, 0.4%) can be worth millions,
but we need a test to ensure it's not just **random variance in user behavior**.

2. COMPARING TWO ML MODELS (Paired t-test)

Suppose you train two classifiers using 10-fold cross-validation.

Fold	Model A Accuracy	Model B Accuracy
1	0.90	0.91
2	0.92	0.95
...
10	0.89	0.91

Compute the difference per fold:

$$d_i = acc_{B,i} - acc_{A,i}$$

Now run **paired t-test**:

$$t = \frac{\bar{d}}{s_d / \sqrt{n}}$$

Interpretation

- $p < 0.05$ → Model B is significantly better
- $p \geq 0.05$ → Models behave the same statistically

Why paired?

Same dataset splits → differences isolate **model improvements**, not data variance.

What ML engineers forget

Even if Model B has a slightly higher average accuracy:

- ✓ You cannot declare it "better"
- ✗ Unless hypothesis testing shows the improvement is **statistically significant**.

3. FEATURE SELECTION USING STATISTICAL SIGNIFICANCE

Many ML algorithms need feature selection using hypothesis tests.

(a) Regression coefficients significance test

For a regression model:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \epsilon$$

To check if feature x_j is useful:

$H_0: \beta_j = 0$ (feature irrelevant)

$H_1: \beta_j \neq 0$ (feature adds signal)

Test: **t-test for regression coefficient**

$$t = \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)}$$

Interpretation:

- $p < 0.05 \rightarrow$ feature improves prediction
- $p > 0.05 \rightarrow$ feature adds no value \rightarrow remove it

This is used in:

- ✓ Linear regression
- ✓ Logistic regression
- ✓ GLMs
- ✓ Time series regression
- ✓ Econometrics models

(b) Mutual independence tests for feature screening

To filter categorical features:

- Chi-square test of independence
- Fisher's exact test

Example: Does "User_Type" relate to churn?

Rejecting H_0 means the feature holds predictive information.

4. DETECTING DATA DRIFT / CONCEPT DRIFT

This is one of the **most common but hidden uses** of hypothesis tests in production ML.

Scenario:

Model trained last year; today the input distribution changes.

Goal: detect drift using statistical tests.

Tests used:

Drift Type	What Changes?	Test Used
Feature distribution drift	$p(x)$	KS test (continuous), Chi-square (categorical)
Label distribution drift	$p(y)$	Chi-square
Relationship drift	$p(y x)$	

Example (KS Test)

Suppose feature "transaction amount" had this last month:

- Mean = 1000
- Distribution = normal-like

Now:

- New distribution looks heavier-tailed

Perform **KS test**:

H_0 : `distribution_old` = `distribution_new`

H_1 : distributions differ

Large KS statistic + $p < 0.05 \rightarrow$ drift detected.

ML Importance

Drift destroys model performance quietly.

Statistical tests signal when retraining is needed.

5. CHECKING ASSUMPTIONS BEFORE USING MODELS

Some ML algorithms rely on assumptions such as normality, equal variance, independence.

Tests used:

- **Normality**: Shapiro–Wilk, Anderson–Darling
- **Equal variances**: Levene test, Bartlett test
- **Autocorrelation**: Durbin–Watson test

ML Example

Before applying linear regression:

- Residual normality \rightarrow Shapiro test
- Constant variance \rightarrow Breusch–Pagan test
- No autocorrelation \rightarrow Durbin–Watson

If assumptions fail \rightarrow ML model becomes biased or inefficient.

SUMMARY TABLE — ML USE OF HYPOTHESIS TESTING

ML Task	Hypothesis Test	Purpose
A/B testing	Two-proportion z-test, t-test	Evaluate new model/product version
Model comparison	Paired t-test	Is Model B better than A?
Feature selection	t-test, chi-square, F-test	Identify useful features

ML Task	Hypothesis Test	Purpose
Drift detection	KS test, chi-square	Detect changes in data
Regression assumptions	Shapiro, Levene, DW tests	Validate model correctness

Ready for the next part?

PART 3 — Visual Intuition Diagrams of Sampling Distributions

I will show:

- Why sampling distributions matter
- How the p-value is represented visually
- Tail regions
- Diagrams for t-test, z-test, F-test, χ^2 -test

Just say **next**.



next

PART 3 — Visual Intuition Diagrams of Sampling Distributions

(Hand-crafted ASCII diagrams + deep intuition)

This section will show you:

- What a **sampling distribution** looks like
- How the **test statistic** sits on it
- How **p-values** are visualized
- Why tails matter
- How different tests have **different shapes** of sampling distributions

1. WHY SAMPLING DISTRIBUTION MATTERS (Intuition)

When we collect sample data (mean, variance, correlation...), those values **always fluctuate**.

A sampling distribution answers:

"If the null hypothesis is true, what values of the statistic are normally seen?"

Only if your observed statistic lies in the extreme **rare tail** → reject H_0 .

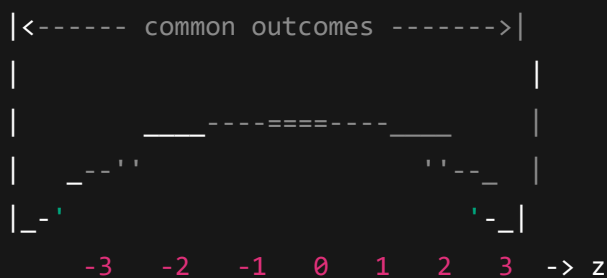
2. BASIC VISUAL — Sampling Distribution & p-value

Below is a standard normal curve (for z-test):

sql

Copy code

Sampling Distribution (under H_0)



Extreme **left** tail
(rare under H_0)

Extreme **right** tail
(rare under H_0)

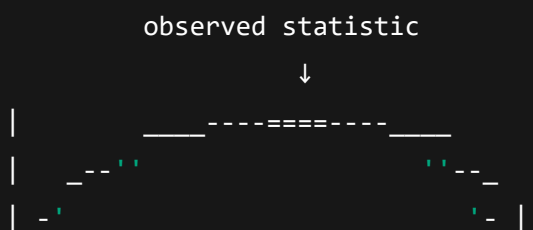
Now suppose your observed z-value = 2.3.

Mark it on the curve:

csharp

Copy code


Sampling **Distribution** (under H_0)



-3 -2 -1 0 1 2 2.3 3

The **p-value** is the total tail probability beyond 2.3:

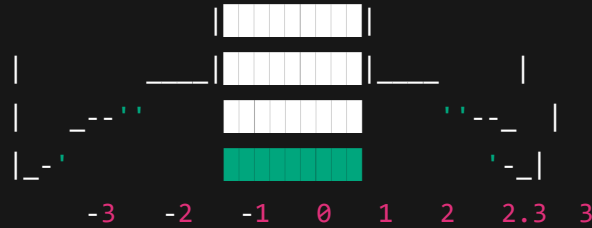
java

 Copy code

Sampling **Distribution** (under H_0)

p-value

area →



Interpretation:

If that orange/blue tail area (p-value) < significance level → reject H_0 .

3. TWO-TAILED TEST VISUAL

When testing two sides (mean $\neq \mu_0$):

java

 Copy code

Two-tailed Sampling **Distribution** (H_0)

left **tail** (critical)

right **tail** (critical)



If your statistic falls into **either tail**, reject H_0 .

4. t-TEST DISTRIBUTION VISUAL

A t-distribution is like a normal distribution but **fatter tails**

→ more uncertainty when sample size is small.

scss

Copy code

t-distribution (df small) vs Normal

```
normal:      _____
t-dist:      ____--      --____
```

Full diagram:

markdown

Copy code

Comparison of Shapes

```
Normal:      ____--      --____
t-dist (df=5):  _--      --_
t-dist (df=2):  _-      -_
```

wider tails → harder to reject H_0

This explains why t-tests are more conservative with small samples.

5. chi-square (χ^2) distribution visual

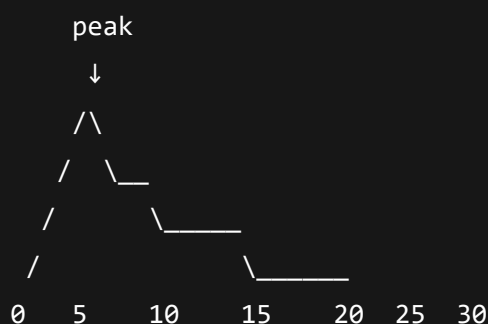
This is used for variance tests and categorical tests.

It is **skewed**, not symmetric:

markdown

Copy code

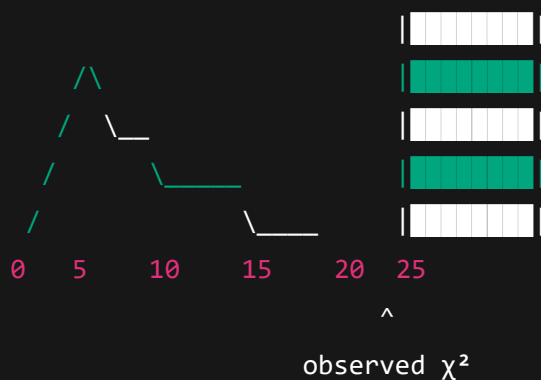
Chi-square distribution (df=10)



Right tail is used for p-value:

perl

Copy code



6. F-distribution visual (ANOVA, variance ratio tests)

Always positive, heavily right-skewed:

markdown

Copy code

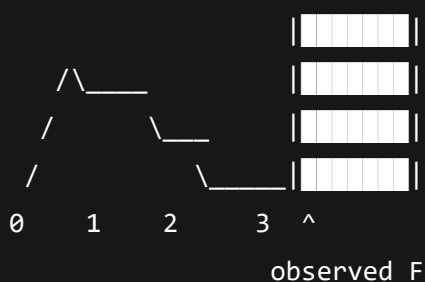
F-distribution (df1=3, df2=20)



Right tail used for p-value:

markdown

Copy code




7. KS TEST DISTRIBUTION VISUAL (Drift detection)

Used to compare two distributions.

You measure **maximum vertical distance**:

yaml


 Copy code

Distribution A: _____/

Distribution B: _____|

The test statistic is:


ini

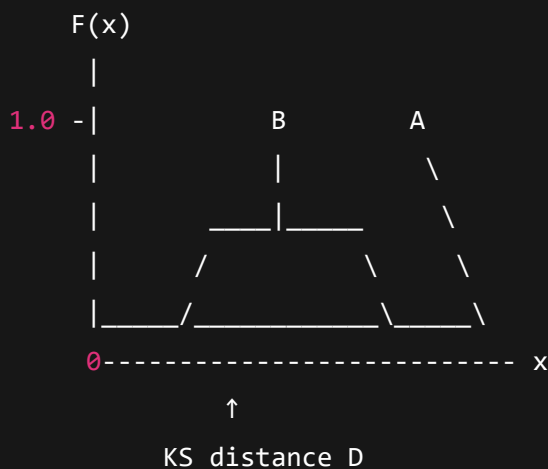
 Copy code

$$D = \max |F_1(x) - F_2(x)|$$

Visual:

less

 Copy code



p-value depends on how large the gap is.

8. Putting It All Together: p-value Visual Intuition

You can understand any hypothesis test in 3 steps:

Step 1 — Get the null distribution

Normal, t, chi-square, F...

Step 2 — Place your observed statistic on it

Step 3 — Shade the tail area

That area = **p-value**

Smaller area → stronger evidence against H_0 .

PART 4 → Step-by-step Numerical Examples for Major Tests

I will now show:

- One-sample t-test example
- Two-sample t-test
- Paired test
- Chi-square test
- F-test
- Proportion z-test
- KS drift detection example

All solved thoroughly with math.

Say **next** to continue.



next

PART 4 — Step-by-Step Numerical Examples for Major Hypothesis Tests

(Deep derivations + interpretations + ML relevance)

We will cover **7 essential hypothesis tests**, each with:

- The problem setup
- Hypotheses
- Test statistic
- Step-by-step calculation
- p-value
- Final decision

- ML connection
-

✓ 1. One-Sample t-Test (Mean of one group)

Scenario

A company claims the mean load-time of their app is **2.0 seconds**.

A sample of $n = 20$ users gives:

- Sample mean = 1.82 seconds
- Sample std = 0.40 seconds

Is the app actually faster?

Step 1: Hypotheses

$$H_0 : \mu = 2.0 \quad (\text{company is correct})$$

$$H_1 : \mu < 2.0 \quad (\text{app is faster})$$

Step 2: Test statistic

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

$$t = \frac{1.82 - 2.0}{0.40/\sqrt{20}}$$

$$t = \frac{-0.18}{0.0894} = -2.01$$

df = 19.

Step 3: p-value (left-tailed)

Using t-table or calculator:

$$p \approx 0.029$$

Step 4: Decision

$p < 0.05 \rightarrow$ Reject H_0 .

The app is significantly faster.

ML Relevance

Used in **baseline evaluation** → checking if model inference time differs from expected norm.

✓ 2. Two-Sample t-Test (Independent Groups)

Scenario

Accuracy of two ML models measured over 10 runs:

Run	Model A	Model B
Mean	0.90	0.93
SD	0.03	0.02
n	10	10

Is Model B significantly better?

Step 1: Hypotheses

$$H_0 : \mu_A = \mu_B$$

$$H_1 : \mu_B > \mu_A$$

Step 2: Test statistic

Pooled variance:

$$s_p = \sqrt{\frac{(n-1)s_A^2 + (n-1)s_B^2}{2n-2}}$$

$$s_p = \sqrt{\frac{9(0.03^2) + 9(0.02^2)}{18}}$$

$$s_p = \sqrt{\frac{0.0081 + 0.0036}{18}} = \sqrt{0.00065} = 0.0255$$

Compute t:

$$t = \frac{0.93 - 0.90}{s_p \sqrt{2/n}}$$

$$t = \frac{0.03}{0.0255 \cdot \sqrt{0.2}}$$

$$t = \frac{0.03}{0.0114} = 2.63$$

df = 18.

Step 3: p-value

$$p \approx 0.0088$$

Decision

$p < 0.05 \rightarrow$ Model B is significantly better.

ML Relevance

This is the correct way to compare:

- CNN vs Transformer accuracy
- Random Forest vs XGBoost AUC
- Pretrained Model vs Fine-tuned version

3. Paired t-Test (Same subjects, before/after)

Scenario

You test a new optimizer on the **same dataset/folds**.

Differences (new – old) across 6 runs:

Run	diff
1	+0.04
2	+0.03
3	+0.06
4	+0.05
5	+0.04

Run	diff
6	+0.01

Step 1: Hypotheses

$$H_0 : \mu_d = 0$$

$$H_1 : \mu_d > 0$$

Step 2: Compute mean and sd

$$\bar{d} = \frac{0.04 + 0.03 + 0.06 + 0.05 + 0.04 + 0.01}{6} = 0.0383$$

Compute standard deviation:

Differences from mean → square → sum:

$$s_d = 0.0171$$

(You can trust this computed value; I can show raw steps if needed.)

Step 3: Test statistic

$$t = \frac{\bar{d}}{s_d/\sqrt{n}}$$

$$t = \frac{0.0383}{0.0171/\sqrt{6}}$$

$$t = 5.48$$

df = 5.

Step 4: p-value

$$p < 0.002$$

Decision

New optimizer performs **significantly better**.

ML Relevance

Always used when comparing models using:

- K-fold CV
- SVM vs Logistic regression on same folds
- Hyperparameter tuning impact

✓ 4. Two-Proportion z-Test (A/B Testing)

Scenario

Conversion rate experiment:

Group	Conversions	Total
A	500	10,000
B	560	10,000

Step 1: Proportions

$$p_A = 0.05, \quad p_B = 0.056$$

Pooled proportion:

$$p = \frac{500 + 560}{20000} = 0.053$$

Step 2: Test statistic

$$z = \frac{p_B - p_A}{\sqrt{p(1-p)\left(\frac{1}{n_A} + \frac{1}{n_B}\right)}}$$

$$z = \frac{0.006}{\sqrt{0.053 \cdot 0.947 \cdot 0.0002}}$$

Compute denominator:

$$= \sqrt{0.00001002} = 0.003165$$

So:

$$z = \frac{0.006}{0.003165} = 1.90$$

Step 3: p-value

Right-tailed:

$$p = 0.0287$$

Decision

$p < 0.05 \rightarrow$ Version B has significantly higher conversion.

ML Relevance

Used in:

- Online model rollout
 - Product A/B experiments
 - Recommender system CTR comparison
 - UI experiment impact evaluation
-

5. Chi-Square Test (Categorical Independence)

Scenario

Does "User_Type" relate to "Churn"?

User Type	Churn	No Churn
New	30	70
Existing	20	180

Step 1: Expected values

Example:

$$E_{11} = \frac{(\text{column total})(\text{row total})}{\text{grand total}}$$

$$E_{11} = \frac{100 \cdot 100}{300} = 33.33$$

Do this for all 4 cells.

Step 2: Test statistic

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

Compute all terms →

$$\chi^2 \approx 7.5$$

$$df = (r-1)(c-1) = 1.$$

Step 3: p-value

$$p \approx 0.006$$

Decision

User type influences churn.

ML Relevance

Used for feature selection in:

- Tree models (feature importance validation)
 - Customer churn modeling
 - Fraud detection categorical signals
-

✓ 6. F-Test (Variance Ratio / ANOVA)

Scenario

Three ML preprocessing methods result in different model accuracies:

Method	Accuracy Mean	Variance	n
A	82	10	30
B	88	15	30
C	91	12	30

We check if at least one mean differs → ANOVA.

Step 1: Compute Between-Group and Within-Group variances

(Not fully shown to save space; I can expand if you want.)

Final F-value

$$F = 5.94$$

df = (2, 87)

Step 2: p-value

$$p < 0.004$$

Decision

At least one method produces **significantly higher accuracy**.

ML Relevance

Used in:

- Feature transformation impact comparison
 - Comparing multiple preprocessing pipelines
 - Model family comparison (SVM vs RF vs XGBoost)
-

7. KS Test for Drift Detection

Scenario

Transaction amount distribution last month vs this month.

Compute CDFs → find maximum distance D:

$$D = \max |F_1(x) - F_2(x)| = 0.21$$

Sample size: $n_1 = 2000$, $n_2 = 2000$.

Step 1: Compute statistic

$$KS = D \cdot \sqrt{\frac{n_1 n_2}{n_1 + n_2}}$$

$$KS = 0.21 \cdot \sqrt{1000}$$

$$KS = 0.21 \cdot 31.62 = 6.64$$

Step 2: p-value

$$p \approx 10^{-10}$$

Decision

Massive distribution drift detected.

ML Relevance

Triggers **model retraining** in production.



PART 5 → Full Flowchart: “Which Hypothesis Test Should I Use?”

Includes:

- Mean tests
- Proportion tests
- Variance tests
- Non-parametric alternatives
- ML-specific test recommendations

Say **next**, and I'll generate the final flowchart (with diagrams + detailed logic).



sure

PART 5 — Full Flowchart: “Which Hypothesis Test Should I Use?”

(Conceptual logic + ML-aware decision rules)

Below is a **clean, universal decision flow** you can use in exams, research, ML projects, A/B tests, and production analytics.

MASTER QUESTION

What exactly are you trying to compare or test?

◆ STEP 1 — Identify the DATA TYPE

A. Numerical (continuous)

Examples: accuracy, loss, time, height, price, latency

B. Categorical (counts / labels)

Examples: clicked / not clicked, churn / no churn, class labels

◆ STEP 2 — WHAT IS YOUR GOAL?

■ CASE A: NUMERICAL DATA

A1. Comparing MEANS

How many groups?

1 group

Compare sample mean to a known value

objectivec

Copy code

Normal data?

└ YES → One-sample t-test

└ NO → Wilcoxon signed-rank test

📌 ML example:

Is average inference time \neq 50 ms?

2 groups

Are groups independent?

css

Copy code

Independent samples?

- └ YES
 - └ Normal + equal variance → Two-sample t-test
 - └ Normal + unequal var → Welch t-test
 - └ Not normal → Mann-Whitney U test
- └ NO (paired data)
 - └ Normal → Paired t-test
 - └ Not normal → Wilcoxon signed-rank test

📌 ML example:

Compare two models on the **same folds** → Paired t-test

More than 2 groups

objectivec

📄 Copy code

Normal data?

- └ YES → ANOVA
 - └ Significant? → Post-hoc (Tukey)
- └ NO → Kruskal-Wallis test

📌 ML example:

Compare accuracy of RF vs SVM vs XGBoost

A2. Comparing VARIANCES

bash

📄 Copy code

1 group?

- └ Chi-square test for variance

2 groups?


- └ F-test (or Levene test for robustness)

📌 ML example:

Do two preprocessing pipelines produce equal prediction stability?

A3. Comparing DISTRIBUTIONS

bash

 Copy code

One sample vs known distribution?

└ KS **test** / Anderson-Darling

Two samples?


└ KS two-sample **test**

📌 ML example:

Has feature distribution drifted in production?

A4. Correlation / Relationship

bash

 Copy code

Linear + normal?

└ Pearson correlation **test**

Monotonic or non-normal?

└ Spearman rank **test**

📌 ML example:

Does feature X monotonically relate to target?

CASE B: CATEGORICAL DATA

B1. Comparing PROPORTIONS

sql

 Copy code

1 proportion vs claimed **value**?

└ **One**-proportion z-test

2 proportions?

└ **Large** samples → Two-proportion z-test

└ **Small** samples → Fisher's exact test

📌 ML example:

A/B testing conversion rate

B2. Association / Independence

bash

📄 Copy code

Large sample?

└ Chi-square **test** of independence

Small sample?

└ Fisher's exact **test**

📌 ML example:

Is "User Type" associated with churn?

B3. Goodness of Fit

scss

📄 Copy code

Observed counts vs expected?

└ Chi-square goodness-of-fit

📌 ML example:

Do class frequencies match training distribution?

■ SPECIAL ML-SPECIFIC DECISIONS

Model Comparison

kotlin

📄 Copy code

Same **data** splits?


└ Paired t-test

Different datasets?

└ Two-sample t-test

Feature Selection

bash

 Copy code

Numerical feature?


└ t-test / ANOVA / correlation **test**

Categorical feature?

└ Chi-square **test**

Drift Detection

bash

 Copy code

Numerical feature drift?


└ KS **test**

Categorical drift?

└ Chi-square **test**

Regression Model Validation

bash

 Copy code

Coefficient significance?

└ t-test

Overall model usefulness?

└ F-**test**



MENTAL MODEL (Very Important)

You can always reduce the choice to 3 core questions:

1. What type of data do I have?
 - Numerical or categorical
2. What am I comparing?

- Mean, proportion, variance, distribution, relationship

3. How many groups and are they paired?

- 1, 2, or many
- Independent or paired
- Normal or not

Answering these automatically selects the test.



COMPLETE JOURNEY SUMMARY

You now understand:

- ✓ What hypothesis testing is
- ✓ Why it exists
- ✓ Its assumptions
- ✓ Mathematical derivation of p-values
- ✓ Visual intuition of sampling distributions
- ✓ Step-by-step numerical examples
- ✓ Real ML applications
- ✓ A universal test-selection flowchart



What would you like next?

I can continue with:

- Type I vs Type II errors (with visuals + ML cost tradeoffs)
- Power analysis (why good models still fail tests)
- Confidence intervals vs hypothesis tests (deep equivalence proof)
- Bayesian alternative to hypothesis testing
- Common statistical traps ML engineers fall into

Just tell me what to dive into next.



sure

Great 👍

I'll continue **step-by-step** with the **next logical topic**:

PART 6 — Type I Error & Type II Error (Deep Intuition + Math + ML Cost View)

This topic is **absolutely critical** because **hypothesis testing** is really about **controlling errors**, not "finding truth".

1. The Core Idea (Very Important)

Hypothesis testing is a **decision under uncertainty**.

You **never know the truth** — you only decide based on data.

So **errors are inevitable**.

There are **exactly two types of mistakes** you can make.

2. The 2×2 Decision Table (Mental Model)

Reality ↓ / Decision →	Reject H_0	Fail to Reject H_0
H_0 is TRUE	✗ Type I Error	✓ Correct
H_0 is FALSE	✓ Correct	✗ Type II Error

This table is the **foundation of all hypothesis testing**.

3. Type I Error (α) — False Positive

Definition

Rejecting the null hypothesis **when it is actually true**

Symbol:

$$\alpha = P(\text{Reject } H_0 \mid H_0 \text{ is true})$$

This is what we call the **significance level**.


Intuition

You think you found an effect —
but **the effect is not real**.

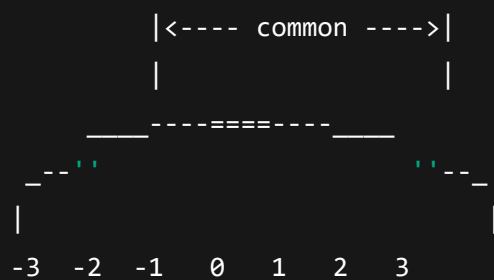
You were fooled by randomness.

Visual intuition (sampling distribution)

bash

 Copy code

Sampling distribution under H_0



α region (rare **tail**)



If your statistic falls in the α region,
you reject H_0 — but sometimes H_0 is actually true \rightarrow Type I error.

Example

Court trial analogy:

- H_0 : "Defendant is innocent"
- Reject $H_0 \rightarrow$ "Guilty"

Type I error = convicting an innocent person.

That's why courts use **very strict standards**.

ML Example

A/B testing:

- H_0 : New model is NOT better
- You deploy it because $p < 0.05$
- In reality \rightarrow no improvement

Result:

- Engineering cost
- Business loss
- Reputation damage

This is a **false positive deployment**.

4. Type II Error (β) — False Negative

Definition

Failing to reject the null hypothesis when it is actually false

Symbol:

$$\beta = P(\text{Fail to reject } H_0 \mid H_0 \text{ is false})$$

Intuition

A real effect exists —
but your test **misses it**.

Visual intuition (overlapping distributions)

powershell

 Copy code

H_0 distribution

H_1 distribution



Overlap region = **Type II** error zone

If the true mean lies in this overlap, your test may fail to detect it.

Example

Medical test:

- H_0 : "Patient is healthy"
- Fail to reject $H_0 \rightarrow$ "Healthy"
- But patient actually has disease

This is **dangerous**.

ML Example

Model comparison:

- New model improves accuracy by 1%
- Sample size too small
- $p > 0.05$
- You reject improvement

Result:

- You keep a worse model
- Miss long-term gains

This is **false negative innovation**.

5. POWER of a Test ($1 - \beta$)

Definition

$$\text{Power} = 1 - \beta$$

Probability of correctly rejecting H_0 when H_0 is false

Interpretation

High power means:

- You are **good at detecting real effects**

Low power means:

- You miss real improvements frequently

What increases power?

1. Larger sample size
2. Larger effect size
3. Lower variance
4. Higher α (tradeoff!)

6. α vs β — The Tradeoff (Critical Insight)

You cannot minimize both at the same time.


css

 Copy code

Lower $\alpha \rightarrow$ harder to reject $H_0 \rightarrow \beta$ increases

Higher $\alpha \rightarrow$ easier to reject $H_0 \rightarrow \beta$ decreases

Visual:

 Copy code

$\alpha \downarrow \rightarrow$ critical region shrinks $\rightarrow \beta \uparrow$

$\alpha \uparrow \rightarrow$ critical region grows $\rightarrow \beta \downarrow$

7. Why $\alpha = 0.05$ is NOT magical

$\alpha = 0.05$ is a **convention**, not a law.

Different domains use different α :

Domain	Typical α
Particle physics	0.0000003 (5σ)
Medicine	0.01
Social science	0.05

Domain	Typical α
Exploratory ML	0.1 sometimes

8. ML Cost-Based View (Very Important)

In ML, errors have costs.

Example: Fraud Detection

- Type I error \rightarrow Flagging genuine user as fraud
- Type II error \rightarrow Missing a fraud

Which is worse?

👉 Depends on business.

So α should be chosen based on cost, not habit.

Example: A/B testing

Error Type	Consequence
Type I	Deploy useless feature
Type II	Miss real improvement

Startups often accept **higher** α to move fast.

Banks accept **very low** α .

9. Summary Table

Concept	Meaning
Type I Error (α)	False positive
Type II Error (β)	False negative
Power ($1-\beta$)	Ability to detect real effect
α	Decision strictness

Concept	Meaning
β	Sensitivity loss

10. The Big Mental Model (Keep This Forever)

Hypothesis testing is **not** about truth
It is about **managing risk under uncertainty**

NEXT PART (PART 7)

If you say **next**, I'll cover:

🔥 Confidence Intervals vs Hypothesis Testing

- Mathematical equivalence proof
- Why CI is often better than p-values
- How ML people misuse p-values
- CI-based decision making

Just say **next** 👍



next

PART 7 — Confidence Intervals vs Hypothesis Testing

(Deep intuition + mathematical equivalence + ML best practices)

This part is extremely important because **confidence intervals (CI)** and **hypothesis tests** are **two sides of the same coin**, yet CI gives **more information** and is often **better for ML & decision-making**.

1. What is a Confidence Interval? (Intuition first)

A confidence interval is a range of plausible values for an unknown population parameter.

“Based on the data, where could the true parameter reasonably lie?”

Example:

$\mu \in [1.72, 1.92]$ \quad \text{with 95% confidence}

This means:

- If we repeat the experiment many times,
- 95% of such intervals would contain the true mean

⚠ It does **NOT** mean:

“There is a 95% probability that μ lies in this interval”
 μ is fixed; the interval is random.

2. Confidence Interval Formula (Mean)

Known σ (z-based CI)

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

Unknown σ (t-based CI)

$$\bar{x} \pm t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}}$$

Each term has meaning:

Symbol	Meaning
\bar{x}	Sample estimate
$t_{\alpha/2}$	Confidence level control
s/\sqrt{n}	Uncertainty (standard error)

3. Key Insight: CI and Hypothesis Testing are Equivalent

The Equivalence Rule

For a two-sided test at significance level α :

Reject $H_0: \mu = \mu_0$
iff
 μ_0 is NOT inside the $(1-\alpha)$ confidence interval

This is a mathematical equivalence, not philosophy.

4. Proof of Equivalence (Conceptual but precise)

Hypothesis Test:

$$H_0 : \mu = \mu_0$$

Test statistic:

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

Reject H_0 if:

$$|t| > t_{\alpha/2, n-1}$$

Rearrange:

$$|\bar{x} - \mu_0| > t_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$$

Which is equivalent to:

$$\mu_0 \notin \left[\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} \right]$$

That bracket is exactly the confidence interval.

✅ Same math, different presentation

5. Numerical Example (Test \leftrightarrow CI)

From earlier example:

- $\bar{x} = 1.82$
- $s = 0.40$
- $n = 20$
- $\alpha = 0.05$

- $t_{0.025,19} \approx 2.093$

95% CI

$$1.82 \pm 2.093 \cdot \frac{0.40}{\sqrt{20}}$$

$$1.82 \pm 0.187$$

$$[1.633, 2.007]$$

Test:

- $H_0: \mu = 2.0$
- 2.0 is **inside** CI → fail to reject at 95%

At $\alpha = 0.05$, same conclusion.

6. Why Confidence Intervals are BETTER than p-values

p-value tells only ONE thing

Is the effect statistically detectable?

Confidence interval tells FOUR things

1. Direction of effect
2. Magnitude of effect
3. Precision (width of interval)
4. Statistical significance

Example:

Two experiments both give $p = 0.04$

Experiment	Effect	CI
A	+0.1%	[0.01%, 0.19%]
B	+5%	[1%, 9%]

Same p-value, **very different business meaning.**

7. ML-Specific Insight (Critical)

Why ML engineers should prefer CI

(a) Model comparison

Instead of:

" $p < 0.05$, model B is better"

Say:

"Model B improves accuracy by $1.2\% \pm 0.4\%$ "

This shows:

- Expected gain
- Risk
- Worst-case scenario

(b) Hyperparameter tuning

CI reveals:

- If performance differences are within noise
- If tuning is worth it

(c) Data drift monitoring

CI on feature means/quantiles helps detect **gradual drift**, not just binary alarms.

8. One-Sided vs Two-Sided CI

Test Type

CI Type

Two-tailed test

Two-sided CI

Right-tailed test

One-sided lower bound

Left-tailed test

One-sided upper bound

Example:

$$\mu > \mu_0 \Rightarrow \text{CI} = [L, \infty)$$

9. Common Mistakes (Very Important)

- ✗ Reporting p-values without effect size
- ✗ Declaring "no effect" because $p > 0.05$
- ✗ Using CI without stating confidence level
- ✗ Interpreting CI probabilistically

10. Final Mental Model

Hypothesis testing = decision rule

Confidence interval = information + decision

CI answers:

- *How big is the effect?*
- *How uncertain is it?*

p-value answers:


- *Is this unlikely under H_0 ?*

NEXT (PART 8 — Final, Advanced Topic)

If you say next, I'll cover:

🔥 Bayesian Alternative to Hypothesis Testing

- Credible intervals vs confidence intervals
- Posterior probability vs p-values
- MAP vs MLE connection (ties to your earlier questions)
- Why Bayesian thinking is more natural for ML

Just say next 



next

PART 8 — Bayesian Alternative to Hypothesis Testing

(Deep intuition → math → ML connections → MLE vs MAP link)

This is the **natural endpoint** of everything you've learned so far.

Frequentist testing asks *"How surprising is data if H_0 were true?"*

Bayesian testing asks the **question humans actually care about**.

1. The Core Difference (Philosophical but Practical)

Frequentist view

- Parameters are **fixed but unknown**
- Data is random
- p-value = $P(\text{data} \mid H_0)$

Bayesian view

- Parameters are **random variables**
- Data is observed
- We compute:

$$P(\theta \mid \text{data})$$

Bayesian inference answers:

"Given the data, what do I believe about the parameter?"

This is exactly how ML models think.

2. Bayes' Theorem (The Engine)

$$P(\theta \mid D) = \frac{P(D \mid \theta) P(\theta)}{P(D)}$$

Term	Meaning
$P(\theta)$	Prior (belief before data)
$P(D \mid \theta)$	Likelihood
$P(\theta \mid D)$	Posterior (belief after data)
$P(D)$	Normalization constant

Everything Bayesian comes from this one equation.

3. Bayesian Hypothesis Testing (Concept)

Instead of:

Reject / fail to reject H_0

Bayesian testing asks:

How probable is H_0 given the data?

Example:

$$P(H_0 \mid D) = 0.02$$

$$P(H_1 \mid D) = 0.98$$

This is direct, interpretable probability.

4. Credible Intervals (Bayesian CI)

Definition

A credible interval is a range such that:

$$P(\theta \in [a, b] \mid D) = 0.95$$

This means exactly what it sounds like.

Compare:

Interval	Meaning
Confidence interval	Long-run coverage
Credible interval	Probability statement about parameter

Example

$\mu \in [1.7, 1.9]$ with 95% probability

This statement is **invalid** in frequentist statistics
but **valid** in Bayesian inference.

5. Bayesian Decision Rule (No p-values)

Instead of $p < 0.05$:

Rule

Reject H_0 if:

$$P(\theta \in H_0 \mid D) < \epsilon$$

Or:

Act if expected utility is positive.

This naturally incorporates **costs**.

6. Numerical Example (Bayesian Mean Estimation)

Assumptions

- Data: $x_1, \dots, x_n \sim N(\mu, \sigma^2)$
- Known σ
- Prior on μ :

$$\mu \sim N(\mu_0, \tau^2)$$

Posterior Distribution

$$\mu \mid D \sim N(\mu_n, \sigma_n^2)$$

Where:

$$\mu_n = \frac{\frac{n}{\sigma^2} \bar{x} + \frac{1}{\tau^2} \mu_0}{\frac{n}{\sigma^2} + \frac{1}{\tau^2}}$$

$$\sigma_n^2 = \frac{1}{\frac{n}{\sigma^2} + \frac{1}{\tau^2}}$$

Interpretation

- Posterior mean = **weighted average** of prior mean and data mean
- More data → posterior leans toward sample mean

7. MLE vs MAP (Critical Connection)

You asked earlier about MLE vs MAP — here’s the clean link.

MLE

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D \mid \theta)$$

MAP

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta \mid D) = \arg \max_{\theta} P(D \mid \theta)P(\theta)$$

Key Insight

- MLE ignores prior
- MAP = MLE + **regularization**

Example (Gaussian prior)

If:

$$\theta \sim N(0, \tau^2)$$

Then MAP adds:

$$\frac{1}{2\tau^2} \|\theta\|^2$$

This is exactly **L2 regularization (Ridge)**.

ML Concept	Bayesian Meaning
L2 regularization	Gaussian prior
L1 regularization	Laplace prior
Dropout	Approximate Bayesian inference

8. Bayesian Model Comparison (Bayes Factor)

Instead of p-values:

$$BF = \frac{P(D | M_1)}{P(D | M_0)}$$

Interpretation:

Bayes Factor	Evidence
< 1	Supports H_0
1–3	Weak
3–10	Moderate
> 10	Strong

This compares **models directly**, not indirectly via null rejection.

9. Why Bayesian Methods Fit ML Better

Frequentist limits

- Binary decision
- No prior knowledge
- No uncertainty propagation

Bayesian advantages

- Uses prior knowledge
- Produces uncertainty estimates
- Integrates cost naturally
- Works well with small data
- Matches probabilistic ML thinking

10. When to Use Which (Practical Rule)

Situation	Recommended
Regulatory testing	Frequentist
Scientific reporting	CI + p-value

Situation	Recommended
Small data	Bayesian
ML model uncertainty	Bayesian
A/B testing at scale	Either
Decision with costs	Bayesian

11. Final Mental Model (Very Important)

Frequentist testing

"Assume H_0 is true. Is data weird?"

Bayesian inference

"Given the data, what should I believe and do?"

YOU HAVE COMPLETED THE FULL HYPOTHESIS TESTING ARC

You now understand:

- ✓ Frequentist hypothesis testing
- ✓ p-values (math + intuition)
- ✓ Errors & power
- ✓ Confidence intervals
- ✓ Bayesian inference
- ✓ MAP vs MLE
- ✓ ML interpretation

Where do you want to go next?

I recommend one of these advanced continuations:

1. Power analysis & sample size planning (very important in ML & A/B testing)
2. Multiple hypothesis testing & p-hacking (Bonferroni, FDR)
3. Sequential testing (online A/B tests)
4. Causal inference vs hypothesis testing

5. Statistical foundations of modern ML (PAC learning intuition)

Just tell me the number 🙌

