

Applied Hypothesis Testing

Advanced Module A0

Digital Finance

Lesson 2: Hypothesis Testing and Statistical Inference

Topics Covered

- Hypothesis testing framework
- Type I and Type II errors
- p-values and significance
- One-sample t-test
- Two-sample t-test
- Paired t-test

Advanced Topics

- One-way ANOVA
- Post-hoc tests
- Power analysis
- Effect sizes
- Multiple comparisons
- Practical examples in R

Focus on practical application and interpretation

Key Components

- Null hypothesis (H_0)
- Alternative hypothesis (H_1 or H_a)
- Test statistic
- Significance level (α)
- p-value
- Decision rule



Hypothesis Testing Framework

The foundation of statistical inference

Two ways to be wrong:

- **Type I Error (α) = False Alarm**
 - Reject true H_0
 - Like convicting an innocent person
 - Controlled by significance level
- **Type II Error (β) = Missed Detection**
 - Fail to reject false H_0
 - Like letting a guilty person go free
 - Related to statistical power



Error Types Diagram

Balance between error types is crucial for study design

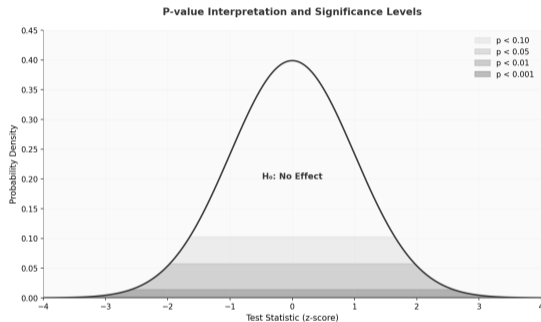
p-value Definition

- Probability of observing data at least as extreme as the actual observed results
- Calculated under the assumption that H_0 is true
- **NOT** the probability that H_0 is true

Interpretation

- $p < 0.05$: “Significant”
- $p < 0.01$: “Highly significant”
- $p < 0.001$: “Very highly significant”

p-values measure evidence against H_0 , not effect size

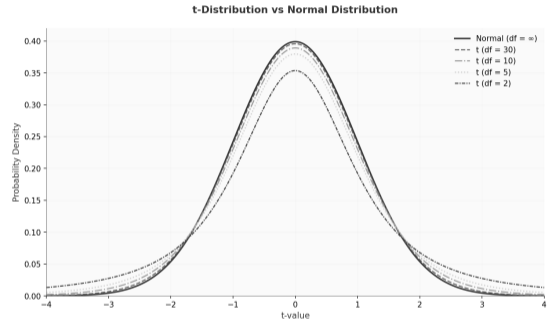


When to Use t-tests

- Comparing means
- Small sample sizes
- Unknown population variance
- Approximately normal data

The t-distribution

- Similar to normal distribution
- Heavier tails
- Approaches normal as df increases



William Sealy Gosset ("Student"), 1908

Hypotheses

- $H_0: \mu = \mu_0$
- $H_1: \mu \neq \mu_0$ (two-tailed)
- $H_1: \mu > \mu_0$ (one-tailed)
- $H_1: \mu < \mu_0$ (one-tailed)

Test Statistic (“How surprising is our data?”)

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{\text{observed} - \text{expected}}{\text{typical variation}}$$

where:

- \bar{x} = sample mean
- μ_0 = hypothesized mean
- s = sample standard deviation
- n = sample size

Foundation for comparing a sample to a known value

Assumptions

- Random sampling
- Independent observations
- Approximately normal distribution
- Or large sample size ($n > 30$)

Degrees of Freedom

$$df = n - 1$$

Decision Rule

- Reject H_0 if $|t| > t_{critical}$
- Or if p-value $< \alpha$

Example: Testing Mean Height

```
1 # Sample data
2 heights <- c(172, 168, 175,
3             170, 165, 178,
4             169, 171, 173,
5             167)
6
7 # Population mean = 170 cm
8 # One-sample t-test
9 t.test(heights,
10        mu = 170,
11        alternative = "two.sided")
12
13 # Output:
14 # t = 0.949, df = 9
15 # p-value = 0.367
16 # 95% CI: [168.2, 173.4]
17 # Sample mean: 170.8
```

Always check assumptions before interpreting results

Interpretation

- Sample mean: 170.8 cm
- Test statistic: $t = 0.949$
- p-value = $0.367 > 0.05$
- Cannot reject H_0
- No evidence that mean differs from 170 cm

Confidence Interval

- 95% CI: [168.2, 173.4]
- Contains $\mu_0 = 170$
- Consistent with hypothesis test

Two-sample t-test: Theory

Purpose

- Compare means of two independent groups
- Test for differences between populations

Hypotheses

- $H_0: \mu_1 = \mu_2$
- $H_1: \mu_1 \neq \mu_2$

Test Statistic (equal variances)

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Most common test for comparing two groups

Pooled Standard Deviation

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

Assumptions

- Independent samples
- Normal distributions
- Equal variances (or use Welch's test)

Degrees of Freedom

- Equal variance: $df = n_1 + n_2 - 2$
- Welch's test: complex formula

Two-sample t-test: R Example

```
1 # Treatment groups
2 control <- c(23, 25, 21,
3             22, 26, 20)
4 treatment <- c(28, 30, 32,
5               29, 31, 27)
6
7 # Test for equality of variances
8 var.test(treatment, control)
9 # F = 1.08, p = 0.94
10
11 # Two-sample t-test
12 t.test(treatment, control,
13        var.equal = TRUE)
14
15 # Output:
16 # t = 5.99, df = 10
17 # p-value = 0.00013
18 # Mean difference: 7.17
```

Always report effect size alongside p-values

Results Interpretation

- Control mean: 22.8
- Treatment mean: 29.5
- Difference: 6.7 units
- p-value < 0.001
- Strong evidence against H_0
- Treatment effect is significant

Effect Size

```
1 # Cohen's d
2 d <- (mean(treatment) -
3      mean(control)) / s_p
4 # d = 2.45 (large effect)
```

When to Use

- Before-after measurements
- Matched pairs
- Repeated measures
- Cross-over designs

Approach

- Calculate differences: $d_i = x_{1i} - x_{2i}$
- Test if mean difference = 0
- Reduces to one-sample t-test

Paired design removes between-subject variability

Test Statistic

$$t = \frac{\bar{d} - 0}{s_d / \sqrt{n}}$$

where:

- \bar{d} = mean of differences
- s_d = standard deviation of differences
- n = number of pairs

Advantages

- Controls for individual differences
- More powerful than independent samples
- Smaller sample size needed

Paired t-test: R Example

```
1 # Weight loss study
2 before <- c(82, 78, 85,
3            90, 76, 88)
4 after <- c(78, 76, 80,
5           84, 73, 82)
6
7 # Paired t-test
8 t.test(before, after,
9        paired = TRUE)
10
11 # Output:
12 # t = 5.77, df = 5
13 # p-value = 0.0022
14 # Mean difference: 5.17 kg
15
16 # Effect size
17 diff <- before - after
18 cohen_d <- mean(diff) /
19          sd(diff)
20 # d = 2.36 (large)
```

Data Structure

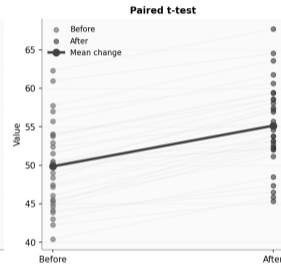
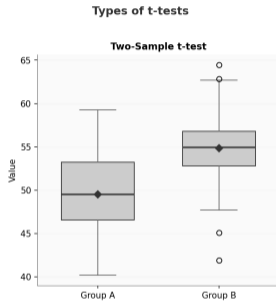
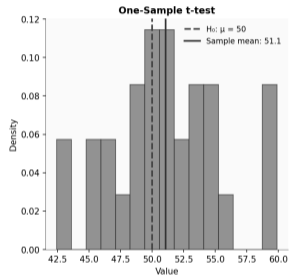
Subject	Before	After
1	82	78
2	78	76
3	85	80
4	90	84
5	76	73
6	88	82

Interpretation

- Average weight loss: 5.17 kg
- p-value = 0.002
- Significant weight reduction
- Large effect size ($d = 2.36$)

Paired design is ideal for before-after comparisons

Types of t-tests: Summary



Choose the appropriate test based on

study design

Analysis of Variance (ANOVA): Concept

Purpose

- Compare means of 3+ groups
- Extension of two-sample t-test
- Test overall difference

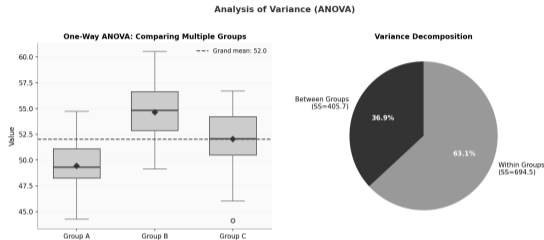
Hypotheses

- $H_0: \mu_1 = \mu_2 = \dots = \mu_k$
- H_1 : At least one mean differs

Key Idea (“Signal vs Noise”)

- **Signal:** Between-group variance (are groups different?)
- **Noise:** Within-group variance (random variation)
- **F-ratio:** Signal / Noise
- Large F = groups really differ

ANOVA tests for any difference among groups



Sum of Squares

- **SST** (Total): Total variability
- **SSB** (Between): Group differences
- **SSW** (Within): Random error
- $SST = SSB + SSW$

Mean Squares

- $MSB = SSB / (k-1)$
- $MSW = SSW / (N-k)$

F-ratio measures signal-to-noise ratio

F-statistic

$$F = \frac{MSB}{MSW} = \frac{\text{Between-group variance}}{\text{Within-group variance}}$$

Degrees of Freedom

- $df_{\text{between}} = k - 1$
- $df_{\text{within}} = N - k$
- $df_{\text{total}} = N - 1$

Decision

- Large $F \rightarrow$ reject H_0
- Compare to F-distribution

One-Way ANOVA: R Example

```
1 # Three teaching methods
2 method_A <- c(78, 82, 75,
3              80, 77)
4 method_B <- c(85, 88, 83,
5              87, 84)
6 method_C <- c(72, 74, 70,
7              73, 71)
8
9 # Combine data
10 scores <- c(method_A,
11            method_B,
12            method_C)
13 method <- factor(rep(
14   c("A", "B", "C"),
15   each = 5))
16
17 # ANOVA
18 model <- aov(scores ~ method)
19 summary(model)
```

ANOVA tells us differences exist, not which groups differ

ANOVA Table

Source	df	SS	MS	F
Method	2	338.5	169.3	31.7
Error	12	64.0	5.3	

p-value < 0.001

Interpretation

- $F(2,12) = 31.7$
- $p < 0.001$
- Strong evidence of differences
- Need post-hoc tests

Group Means

- Method A: 78.4
- Method B: 85.4
- Method C: 72.0

Key Assumptions

- 1 Independence of observations
- 2 Normality within groups
- 3 Homogeneity of variances

Checking Assumptions

- **Normality:** Q-Q plots, Shapiro-Wilk test
- **Equal variances:** Levene's test, Bartlett's test
- **Independence:** Study design

Robustness

- ANOVA fairly robust to normality violations
- Sensitive to unequal variances
- Very sensitive to dependence

Alternatives

- Welch's ANOVA (unequal variances)
- Kruskal-Wallis test (non-parametric)
- Transformation of data

Always check assumptions before interpreting ANOVA

The Problem

- ANOVA: overall difference
- Which groups differ?
- Multiple pairwise tests
- Inflated Type I error

Example: 4 groups

- $\binom{4}{2} = 6$ comparisons
- If $\alpha = 0.05$ per test
- Family-wise error rate:

$$1 - (1 - 0.05)^6 = 0.265$$

- 26.5% chance of false positive!

Balance between Type I error control and statistical power

Solutions

- **Bonferroni:** $\alpha_{adj} = \alpha/m$
- **Tukey HSD:** All pairwise
- **Dunnnett:** Compare to control
- **Scheffé:** All contrasts

Trade-offs

- Control Type I error
- Reduced power
- Conservative corrections

```
1 # Tukey HSD test
2 TukeyHSD(model)
3
4 # Output (abbreviated):
5 #      diff      p adj
6 # B-A      7.0  0.0003
7 # C-A     -6.4  0.0008
8 # C-B    -13.4  0.0001
9
10 # Bonferroni correction
11 pairwise.t.test(scores,
12                 method,
13                 p.adj = "bonf")
14
15 # Dunnett test (A as control)
16 library(multcomp)
17 dunnett <- glht(model,
18                linfct = mcp(
19                  method = "Dunnett"))
20 summary(dunnett)
```

Choose post-hoc test based on research questions

Tukey HSD Results

Comparison	Difference	95% CI	p-adj
B - A	7.0	[3.8, 10.2]	0.0003
C - A	-6.4	[-9.6, -3.2]	0.0008
C - B	-13.4	[-16.6, -10.2]	0.0001

Interpretation

- All pairwise differences significant
- Method B > Method A > Method C
- Adjusted for multiple comparisons
- Family-wise error controlled at 0.05

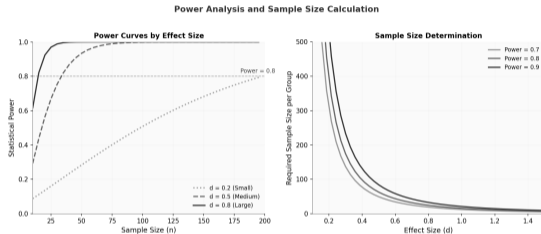
Definition (“Your ability to find what’s really there”)

- Power = $1 - \beta = 1 - P(\text{missed detection})$
- Probability of detecting a real effect
- Want power ≥ 0.80 (80% chance to detect)

Factors Affecting Power

- Effect size (larger \rightarrow more power)
- Sample size (larger \rightarrow more power)
- Significance level (α)
- Variability (less \rightarrow more power)

Power analysis essential for study design



```
1 library(pwr)
2
3 # Two-sample t-test power
4 # Effect size d = 0.8
5 # n per group = 30
6 pwr.t.test(n = 30,
7           d = 0.8,
8           sig.level = 0.05,
9           type = "two.sample")
10 # Power = 0.86
11
12 # Sample size for power=0.8
13 pwr.t.test(d = 0.5,
14           power = 0.8,
15           sig.level = 0.05,
16           type = "two.sample")
17 # n = 64 per group
```

Always conduct power analysis before data collection

ANOVA Power

```
1 # 3 groups, effect size f=0.4
2 pwr.anova.test(k = 3,
3              f = 0.4,
4              sig.level = 0.05,
5              power = 0.8)
6 # n = 21 per group
7
8 # Detectable effect
9 pwr.anova.test(k = 3,
10              n = 15,
11              sig.level = 0.05,
12              power = 0.8)
13 # f = 0.47
```

Planning Studies

- Estimate effect size from pilot/literature
- Calculate required sample size
- Consider practical constraints

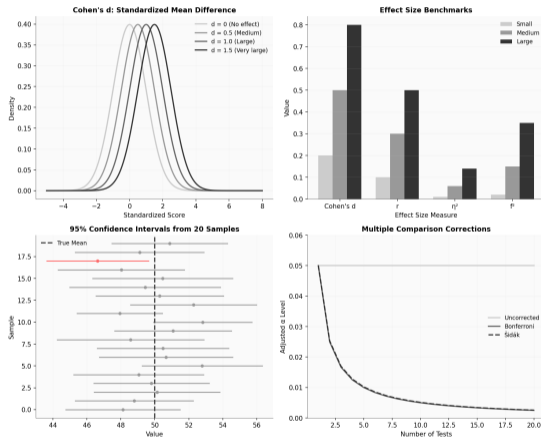
Why Effect Sizes? (“How big is the effect?”)

- p-value = “Is it real?” (yes/no)
- Effect size = “How big?” (magnitude)
- With large n, tiny effects become significant
- Effect sizes are comparable across studies

Common Measures

- **Cohen's d**: Mean difference in SD units (small=0.2, medium=0.5, large=0.8)
- **r**: Correlation coefficient
- η^2 : Variance explained (ANOVA)

Effect Sizes and Statistical Considerations



Always report effect sizes with hypothesis tests

Study Design

- Compare 3 treatments
- Primary outcome: Pain score
- $n = 20$ per group
- $\alpha = 0.05$

```
1 # Load data
2 data <- read.csv("trial.csv")
3
4 # Check assumptions
5 bartlett.test(pain ~ treatment,
6               data = data)
7 # p = 0.34 (equal var OK)
8
9 # ANOVA
10 model <- aov(pain ~ treatment,
11             data = data)
12 summary(model)
13 # F(2,57) = 8.42, p = 0.0006
```

Complete analysis includes assumptions, effect size, and power

```
1 # Effect size
2 library(effectsize)
3 eta_squared(model)
4 # eta2 = 0.228 (large)
5
6 # Post-hoc tests
7 TukeyHSD(model)
8 # Drug A vs Placebo: p=0.001
9 # Drug B vs Placebo: p=0.023
10 # Drug A vs Drug B: p=0.142
11
12 # Power analysis (post-hoc)
13 pwr.anova.test(k = 3,
14               n = 20,
15               f = 0.54,
16               sig.level = 0.05)
17 # Power = 0.84
```

Conclusions

- Both drugs superior to placebo
- No difference between drugs
- Adequate power achieved

Before Analysis

- Define hypotheses clearly
- Choose appropriate test
- Set significance level
- Conduct power analysis
- Plan for multiple comparisons

During Analysis

- Check all assumptions
- Use appropriate corrections
- Calculate effect sizes
- Document all decisions

Reporting Results

- State hypotheses tested
- Report test statistics
- Include exact p-values
- Provide effect sizes
- Show confidence intervals
- Interpret in context

Common Pitfalls

- p-hacking
- Ignoring assumptions
- Multiple testing without correction
- Confusing significance with importance

Good statistical practice ensures reproducible research

Key Concepts

- Hypothesis testing framework
- Type I and Type II errors
- p-values and their interpretation
- Effect sizes complement p-values

Tests Covered

- One-sample t-test
- Two-sample t-test
- Paired t-test
- One-way ANOVA
- Post-hoc comparisons

Important Reminders

- Always check assumptions
- Consider study design
- Plan sample size with power analysis
- Report complete results
- Statistical vs practical significance

Next Topics

- Two-way ANOVA
- Non-parametric tests
- Multiple regression
- Mixed models

Statistical inference is the foundation of empirical research