

STA 35C: Statistical Data Science III

Lecture 18: Multiple Hypotheses Testing

Dogyoon Song

Spring 2026, UC Davis

Announcement

Homework 5 is due tomorrow

Midterm 2 on Fri, May 15 (1:10 pm–2:00 pm in class)

- **Arrive early:** The exam starts at 1:10 pm and ends at 2:00 pm sharp
- **One hand-written cheat sheet:** Letter-size (8.5"×11"), double-sided, brief formulas/notes
- **Calculator:** A simple (non-graphing) scientific calculator is allowed
- **No other materials** beyond the single cheat sheet (no textbooks, etc.)
- **SDC accommodations:** Confirm scheduling with AES online ASAP

Advice for preparation:

- Primary coverage: Lectures 12–19
- Core concepts from earlier may be assumed, especially regression and model assessment
- Two practice midterms and brief solution keys are posted on course webpage
- Office hours next week:
 - Instructor: Wed, 3:30–4:30 pm
 - TA: Tue 12:30–2:30pm

Agenda

Recap: Challenges & motivation for multiple hypotheses testing

- Issues arising with multiple tests
- Real-world concerns: p-hacking & data dredging

Today: How to control false positives

- **Family-wise error rate (FWER)**
 - Definition & intuition
 - Controlling FWER: Bonferroni correction & Holm's step-down
- **False discovery rate (FDR)**
 - Definition & intuition
 - Controlling FDR: Benjamini–Hochberg procedure

Recap: Multiple testing

Single-hypothesis test:

- Typically set up H_0 , and gather data to reject it if there is significant evidence
- Type I error = false positive; Type II error = false negative
- Each test has Type I error at most α (e.g. 0.05)

Modern data analysis: multiple tests simultaneously

- E.g. Testing thousands of predictors or biomarkers to discovery significant ones
- If m is large, false rejections can occur easily by chance even when all nulls are true
- If all m nulls are true and each test is at level α , the expected number of false positives is $m \times \alpha$

Key challenge: Address and adjust for the inflation of false positives as m grows

Articles warning about misused statistical significance

PLOS Medicine

Browse Publish About

OPEN ACCESS

ESSAY

Why Most Published Research Findings Are False

John P. A. Ioannidis

Published: August 30, 2005 • <https://doi.org/10.1371/journal.pmed.0020124>

The Economist

Weekly edition The world in brief War in the Middle East War in Ukraine United States The world economy Business

Briefing | Unreliable research

Trouble at the lab

Scientists like to think of science as self-correcting. To an alarming degree, it is not

Oct 18th 2013

nature

Explore content v About the journal v Publish with us v Subscribe

[nature](#) > [comment](#) > [article](#)

COMMENT | 20 March 2019

Scientists rise up against statistical significance

Valentin Amrhein, Sander Greenland, Blake McShane and more than 800 signatories call for an end to hyped claims and the dismissal of possibly crucial effects.

By [Valentin Amrhein](#), [Sander Greenland](#) & [Blake McShane](#)

Figure: Many reproducibility crises trace back to undisclosed multiple testing or selective reporting. Proper adjustments can help mitigate these issues.

Related issues: p -hacking and data dredging

Real danger: Searching for “significant” results in many ways until something “works”

- Repeatedly testing different hypotheses/subgroups
- Eventually, some test may yield $p < 0.05$ *by chance*
- Example: trying many outcomes, subgroups, transformations, or model specifications, then reporting only the smallest p -value

Outcome: Spurious “discoveries”

- Published claims may fail to replicate
- True findings can be overshadowed by noise

Conclusion: Systematic multiple-testing corrections are crucial, especially for large m

Recall: Hypothesis testing as classification

A single hypothesis test makes a decision about whether to reject H_0 :

- **Goal:** Discover a "real phenomenon" (H_1) or conclude not enough evidence (H_0)
 - H_0 is true \iff no real effect
 - H_0 is false \iff there is a real effect (H_1)
 - We claim a "discovery" when H_0 is rejected
- **Analogy to classification:** Depending on the evidence in the data,
 - Reject $H_0 \iff$ classify $\hat{H} = 1$
 - Fail to reject $H_0 \iff$ classify $\hat{H} = 0$

| | H_0 is true (" H=0 ") | H_0 is not true (" H=1 ") |
|--|--------------------------------|------------------------------------|
| Reject H_0 (" $\hat{H}=1$ ") | Type-I error (FP) | Correct (TP) |
| Not reject H_0 (" $\hat{H}=0$ ") | Correct (TN) | Type-II error (FN) |

$\implies \Pr(\text{Type I error}) = \Pr(\text{reject a true null}) \approx \frac{FP}{FP+TN}$ empirically among true nulls

- By setting a threshold α , we want to control $\Pr(\text{Type I error})$ below α

Family-wise error rate (FWER): Definition

Single test: $\Pr(\text{Type I error}) = \Pr(\text{reject a true null})$

Multiple tests (m hypotheses):

$$\begin{aligned}\text{FWER} &= \Pr(\text{reject at least one true } H_0) \\ &= \Pr(\# \text{ Type-I error} \geq 1),\end{aligned}$$

i.e. the probability of *any* false positive among m tests

If all m null hypotheses are true, the tests are independent, and each are at level α :

$$\text{FWER} = 1 - (1 - \alpha)^m,$$

- When $m = 1$, $\text{FWER} = 1 - (1 - \alpha)^m = 1 - (1 - \alpha) = \alpha$
- Grows quickly with m
 - E.g. $m = 20$, $\alpha = 0.05 \implies \text{FWER} \approx 0.64 \gg 0.05$

Visualization: FWER grows as m increases

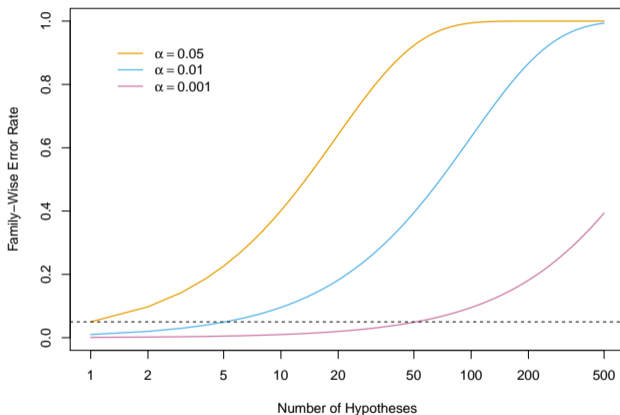


Figure: FWER vs. number of tests m (log scale) for $\alpha = 0.05$ (orange), 0.01 (blue), 0.001 (purple). The dashed line is 0.05. For $m = 50$ and target FWER=0.05, each test must be at $\alpha = 0.001$ [JWHT21, Figure 13.2].

Testing multiple hypotheses at level α

Suppose we test m hypotheses $H_{0,1}, \dots, H_{0,m}$, all at level α , obtaining confusion matrix:

| | H_0 is true | H_0 is not true |
|------------------|---------------|-------------------|
| Reject H_0 | N_{FP} | N_{TP} |
| Not reject H_0 | N_{TN} | N_{FN} |

- $N_{FP}, N_{TP}, N_{TN}, N_{FN}$ are random variables that sum to m
- Roughly, we expect $N_{FP} \approx m \cdot p_{FP}$; when all m nulls are true, $N_{FP} \approx m \cdot \alpha$

If these m tests are independent,

- Probability of *at least one* false positive $\approx 1 - (1 - \alpha)^m$
- For $m = 20, \alpha = 0.05$, that probability is $\approx 64\%$

Goal of FWER control: Keep the probability of making even one false positive small:

$$\text{FWER} = \Pr(N_{FP} \geq 1) \leq \alpha$$

The Bonferroni correction for FWER control

Key idea: Use the union bound over the true null hypotheses:

$$\text{FWER} = \Pr \left(\sum_{j=1}^m \{\text{Reject } H_j\} \right) \leq \sum_{j=1}^m \Pr (\{\text{Reject } H_j\})$$

- Each test is done at level $\alpha/m \implies \Pr (\{\text{Reject } H_j\}) \leq \alpha/m \implies \text{FWER} \leq \alpha$

The Bonferroni method (Bonferroni correction):

- For each hypothesis H_1, \dots, H_m , reject H_j if $p_j < \frac{\alpha}{m}$

Pros & Cons:

- **Pros:** Simple & widely used
- **Cons:** Often *very conservative* \implies few rejections (=discoveries) & lower power¹

¹Power = TPR = the fraction of false null hypotheses that are successfully rejected

Example: Bonferroni correction

Example

Let $m = 6$ hypotheses with p-values:

$$p_1 = 0.0018, \quad p_2 = 0.009, \quad p_3 = 0.021, \quad p_4 = 0.034, \quad p_5 = 0.045, \quad p_6 = 0.070.$$

At $\alpha = 0.05$, threshold = $\frac{\alpha}{m} = \frac{0.05}{6} \approx 0.00833$.

Reject H_j if $p_j < 0.00833$.

Hence:

$$p_1 = 0.0018 < 0.00833 \implies \text{reject } H_1,$$

but $p_2 = 0.009 > 0.00833$ and the rest are larger. So Bonferroni rejects only H_1 .

Conclusion: 1 rejection using Bonferroni, whereas naive $p < 0.05$ would reject 5 of them (p_1, \dots, p_5).

Holm's step-down procedure for FWER control

Holm's method refines Bonferroni to be less conservative:

Holm's method

- 1 Specify α , the level at which to control the FWER
- 2 Compute the p -values for the m null hypotheses, H_{01}, \dots, H_{0m}
- 3 Sort p -values so that $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(m)}$
- 4 Define

$$L = \min \left\{ j : p_{(j)} > \frac{\alpha}{m + 1 - j} \right\}$$

- 5 Reject all null hypotheses H_{0j} for which $p_j < p_{(L)}$

Properties:

- Ensures $\text{FWER} \leq \alpha$
- Rejects at least as many hypotheses as Bonferroni

Example: Holm's step-down procedure

Example

Step 1: Set $\alpha = 0.05$

Step 2: $p_1 = 0.0018, p_2 = 0.009, p_3 = 0.021, p_4 = 0.034, p_5 = 0.045, p_6 = 0.070$.

Step 3: Sort p -values $p_{(1)} = 0.0018, p_{(2)} = 0.009, p_{(3)} = 0.021, p_{(4)} = 0.034, p_{(5)} = 0.045, p_{(6)} = 0.070$.

Step 4: Find $L = 3$ because

$$p_{(1)} = 0.0018 ? 0.0018 \leq \frac{0.05}{6 + 1 - 1} = \frac{0.05}{6} \approx 0.00833 \quad \Rightarrow \text{reject } H_{(1)}, \text{ continue}$$

$$p_{(2)} = 0.009 ? 0.009 \leq \frac{0.05}{6 + 1 - 2} = \frac{0.05}{5} = 0.01 \quad \Rightarrow \text{reject } H_{(2)}, \text{ continue}$$

$$p_{(3)} = 0.021 ? 0.021 \leq \frac{0.05}{6 + 1 - 3} = \frac{0.05}{4} = 0.0125? \text{ No} \quad \Rightarrow \text{stop; } L = 3$$

Step 5: We reject $H_{(1)}, H_{(2)}$ total 2 rejections. The rest are not rejected.

Conclusion: Holm's method rejects 2, whereas Bonferroni rejected only 1.

Visualization: Bonferroni vs. Holm

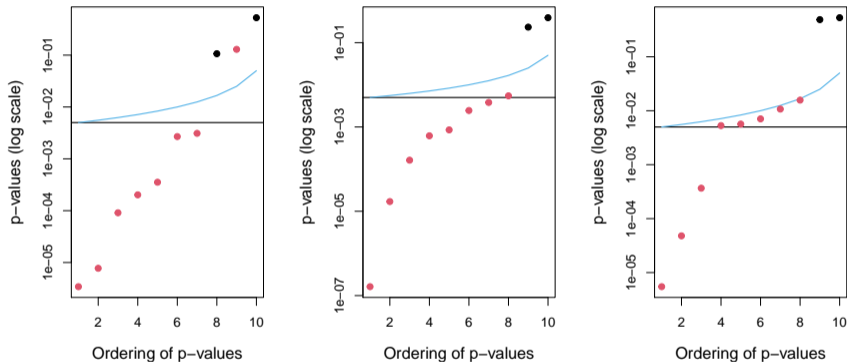


Figure: Each panel shows sorted p -values from a separate simulation of $m = 10$ null hypotheses, with the two true nulls in black and the others in red. Controlling the FWER at 0.05, Bonferroni rejects all points below the **black** line, while Holm rejects all below the **blue** line. The gap between these lines indicates the additional hypotheses Holm rejects but Bonferroni does not. In the middle panel, Holm rejects one more null than Bonferroni; in the right panel, it rejects five more [JWHT21, Figure 13.3].

Pop-up quiz: Bonferroni vs. Holm

Suppose $m = 4$, $\alpha = 0.05$, and the sorted p -values are

$$p_{(1)} = 0.010, \quad p_{(2)} = 0.015, \quad p_{(3)} = 0.030, \quad p_{(4)} = 0.200.$$

Question: How many hypotheses are rejected by Bonferroni and by Holm's method?

- A) Bonferroni rejects 1; Holm rejects 2.
- B) Bonferroni rejects 2; Holm rejects 1.
- C) Both reject 1.
- D) Both reject 2.

Answer: A. Bonferroni uses $0.05/4 = 0.0125$, so it rejects only $p_{(1)}$. Holm compares $0.010 \leq 0.0125$, then $0.015 \leq 0.05/3$, then stops at $0.030 > 0.05/2$, so it rejects two.

Illustration: Power vs. FWER trade-off

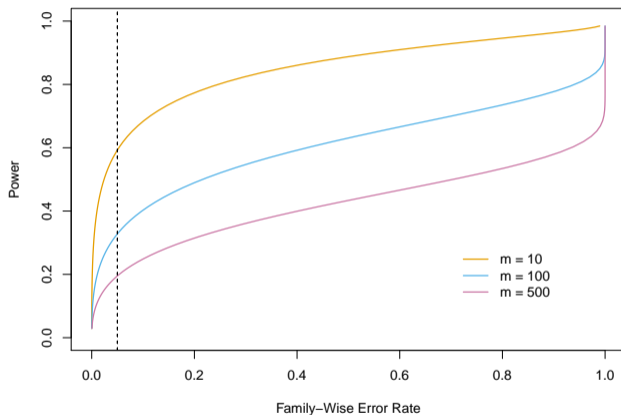


Figure: In a simulation with 90% of m nulls true, the power is displayed against FWER. Colors of the curves: $m = 10$ (orange), $m = 100$ (blue), $m = 500$ (purple). Larger m reduces power. The vertical dashed line marks FWER=0.05 [JWHT21, Figure 13.5].

FWER control may not suffice

FWER demands *no* false rejections with probability at least $1 - \alpha$:

- Very stringent if m is large
- Tends to reduce power (fewer true positives found)

In modern “exploratory” studies:

- We may tolerate a small fraction of false positives to discover more true ones
- This leads to the *false discovery rate (FDR)* approach

| | H_0 is true | H_0 is not true |
|------------------|-------------------|--------------------|
| Reject H_0 | Type-I error (FP) | Correct (TP) |
| Not reject H_0 | Correct (TN) | Type-II error (FN) |

New goal: Find more true discoveries while controlling the proportion of false discoveries.

False discovery rate (FDR): Definition and motivation

Motivation: Controlling FWER can be too conservative for large m

Instead: Control the fraction of rejected hypotheses that are *false positives*

$$\text{FDP} = \frac{\# \text{ false positives}}{\# \text{ total rejections}} = \frac{\# \text{ FP}}{\# \text{ FP} + \# \text{ TP}}$$

- Controlling FDP is impossible because we never know which H_{0j} are true/false

False discovery rate (FDR) = $\mathbb{E}[\text{FDP}]$

- Allow up to fraction q of false positives *on average* among the “claimed discoveries”
- The choice of q is context- and dataset-dependent (no gold standard like $\alpha = 0.05$)

Properties:

- Accept a small fraction of false positives, in exchange for more total discoveries
- Typically yields more rejections (“discoveries”) than FWER-based methods

Controlling FDR: Benjamini–Hochberg procedure

Benamini-Hochberg procedure

- 1 Specify q , the level at which to control the FDR
- 2 Compute the p -values for the m null hypotheses, H_{01}, \dots, H_{0m}
- 3 Sort p -values so that $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(m)}$

4 Define

$$L = \max \left\{ j : p_{(j)} < \frac{qj}{m} \right\}$$

- 5 Reject all null hypotheses H_{0j} for which $p_j \leq p_{(L)}$

Result:

- Under standard conditions, ensures $\text{FDR} \leq q$, but not necessarily small FWER
- Typically more powerful, yielding more rejections, than Bonferroni/Holm if m is large

Example: Benjamini–Hochberg procedure

Example

Step 1: Set $q = 0.05$

Step 2: $p_1 = 0.0018, p_2 = 0.009, p_3 = 0.021, p_4 = 0.034, p_5 = 0.045, p_6 = 0.070$.

Step 3: Sort p -values $p_{(1)} = 0.0018, p_{(2)} = 0.009, p_{(3)} = 0.021, p_{(4)} = 0.034, p_{(5)} = 0.045, p_{(6)} = 0.070$.

Step 4: Find $L = 3$ because

$$j = 1: 0.0018 \leq 0.05 \times \frac{1}{6} \approx 0.0083? \checkmark$$

$$j = 2: 0.009 \leq 0.05 \times \frac{2}{6} \approx 0.0167? \checkmark$$

$$j = 3: 0.021 \leq 0.05 \times \frac{3}{6} = 0.025? \checkmark$$

$$j = 4: 0.034 \leq 0.05 \times \frac{4}{6} \approx 0.0333? \text{ No } (0.034 > 0.0333)$$

Step 5: Reject $H_{(1)}, H_{(2)}, H_{(3)}$.

Conclusion: BH rejects 3, while Holm rejects 2, Bonferroni rejects 1.

Visual comparison: Bonferroni vs. Benjamini–Hochberg

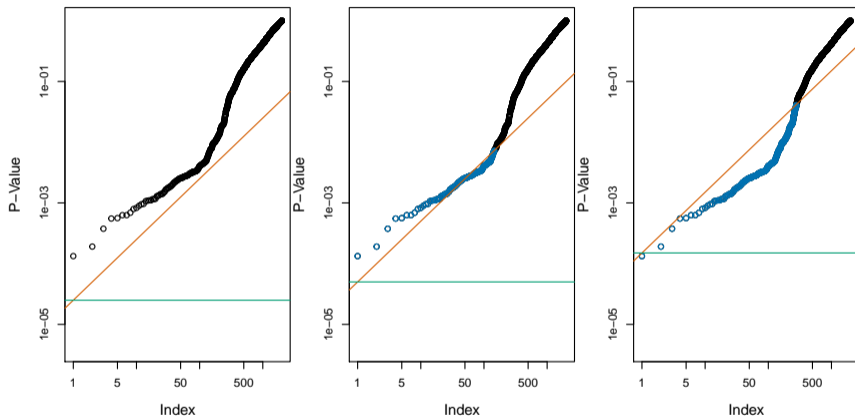


Figure: Panels: same set of $m = 2000$ sorted p-values for the Fund dataset. Green lines: thresholds for FWER control (Bonferroni) at $\alpha = 0.05, 0.1, 0.3$ (left to right). Orange lines: thresholds for FDR control (Benjamini–Hochberg) at $q = 0.05, 0.1, 0.3$ (left to right). E.g., When the FDR is controlled at $q = 0.1$, 146 nulls are rejected (center, blue points). At $q = 0.3$, 279 nulls are rejected (right, blue points) [JWHT21, Figure 13.6].

Pop-up quiz: Comparing FDR vs. FWER

Suppose there are m hypotheses to test simultaneously. **Which statement best captures the difference between FDR and FWER?**

- A) FDR forces the probability of zero false positives to stay below α , whereas FWER allows a small fraction q .
- B) FDR aims to keep $\mathbb{E}[\text{fraction of false positives among rejections}] \leq q$, while FWER demands $\Pr(\text{at least one false positive}) \leq \alpha$.
- C) Under FDR control, no false positives are allowed once you discover enough true positives.
- D) FDR only works for independent tests, but FWER can handle correlated tests without adjustments.

Answer: B.

FDR control (e.g., Benjamini–Hochberg) allows a certain fraction of false positives *on average*, whereas FWER control (e.g., Bonferroni/Holm) requires the chance of *any* false positive be controlled below α .

Wrap-up

- Multiple testing increases the chance of false positives; ordinary $p < 0.05$ rules are not enough when many hypotheses are tested.
- **FWER** controls the probability of making false positives:
 - Strictly ensures $\Pr(N_{\text{FP}} \geq 1) \leq \alpha$, e.g., via Bonferroni/Holm
 - Conservative for large m , leading to fewer rejections & reduced power
- **FDR** controls the expected fraction of false positives among rejections:
 - Controls $\text{FDR} = \mathbb{E}[\text{FDP}] \leq q$, e.g., via Benjamini–Hochberg
 - Typically yields more rejections than FWER, especially when m is large
- **Practical consideration:**
 - Use FWER for strict confirmatory analyses needing minimal Type I error
 - Use FDR for exploratory, large-scale studies, tolerating some false positives to gain more discoveries

References



Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani.

An Introduction to Statistical Learning: with Applications in R, volume 112 of *Springer Texts in Statistics*.

Springer, New York, NY, 2nd edition, 2021.