

# The Benjamini-Hochberg Procedure: Controlling the False Discovery Rate in Multiple Hypothesis Testing

Authored by  
**Mohammed loot**

November 8, 2025

## RECOMMENDED CITATION

Mohammed loot (2025). *The Benjamini-Hochberg Procedure: Controlling the False Discovery Rate in Multiple Hypothesis Testing*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=13837>

The core of modern empirical science relies heavily on [statistical hypothesis testing](#), a methodical approach used to validate or reject conjectures based on observed data. However, inherent in this methodology is the ever-present risk of drawing an incorrect conclusion. Specifically, when we execute a single statistical test, there is a defined probability that the resulting [p-value](#)--the measure of evidence against the null hypothesis--will fall below the predetermined significance threshold (often set at the conventional alpha,  $\alpha = 0.05$ ). This can occur purely by chance, even if the underlying [null hypothesis](#) is entirely true. This scenario constitutes a [Type I error](#), commonly referred to in this context as a **false discovery**.

To illustrate this fundamental challenge, consider a straightforward experiment in biology: a team of botanists wants to determine if a new fertilizer causes a specific plant species to grow taller than a benchmark mean height of 10 inches. The classical hypotheses formulated for this single test are:

**H<sub>0</sub>:**  $\mu = 10$  inches (The true mean height is 10 inches, the fertilizer has no effect.)

**H<sub>A</sub>:**  $\mu > 10$  inches (The true mean height is greater than 10 inches.)

Even if the fertilizer has absolutely no effect (meaning H<sub>0</sub> is true), random chance dictates that the collected [sample](#) of plants might, by luck, include specimens that are unusually tall. If the statistical calculations yield a p-value below the critical alpha ( $\alpha$ ), the researchers are compelled to reject the null hypothesis. In this outcome, the researcher claims a "significant result"--a discovery that the fertilizer works--when in reality, the null hypothesis was true all along. This erroneous rejection is the practical definition of a **false discovery**, and its accumulation can severely compromise the reliability of scientific literature.

## The Exponential Challenge of Multiple Comparisons

While the rate of a **false discovery** in a single, isolated test is generally controllable--typically capped at 5% using  $\alpha = 0.05$ --this risk does not remain stable when researchers move from one test to many. The probability of encountering at least one false positive escalates dramatically when numerous statistical tests are conducted simultaneously, a severe methodological concern known as the [multiple comparisons problem](#). For instance, if a diligent researcher performs 100 independent statistical tests, and assumes the null hypothesis is true for every single one of them, the laws of probability dictate that we should statistically anticipate approximately 5 of those 100 tests (5% of 100) to result in a Type I error, leading to 5 expected **false discoveries**.

This issue is profoundly magnified by the scale of contemporary research. Technological advancements in fields such as high-throughput genomics, advanced neuroscience imaging, and large-scale machine learning applications now allow scientists to execute hundreds, thousands, or even millions of statistical comparisons within a single study. A prime example is genetic

association studies, where medical researchers routinely screen tens of thousands of genes to identify those whose expression levels correlate with a specific disease. If 10,000 such tests are run, even using an extremely conservative significance level of  $\alpha = 0.01$  would still statistically predict 100 false positive results. Without systematic adjustment, these spurious findings consume vast resources in follow-up research based purely on statistical noise.

Consequently, the necessity for robust error control methods becomes paramount. Traditional and highly stringent methods, such as the [Bonferroni correction](#), are designed to control the [Family-Wise Error Rate \(FWER\)](#)--the probability of making at least one false discovery among the entire collection of tests. While effective at minimizing error, FWER corrections are often criticized for being excessively conservative. This high stringency results in a significant loss of [statistical power](#), meaning they dramatically increase the likelihood of a Type II error (a false negative, or missing a true effect). The field required a more nuanced and powerful approach: one that controls the expected proportion of false discoveries made among all the null hypotheses that are rejected. This delicate balance is achieved through the control of the [False Discovery Rate \(FDR\)](#), which is most effectively managed by the seminal **Benjamini-Hochberg Procedure**.

## Introducing the Benjamini-Hochberg Procedure

The **Benjamini-Hochberg Procedure** (BH procedure), introduced by Yoav Benjamini and Yosef Hochberg in 1995, is a sequential, step-up method specifically engineered to control the [False Discovery Rate \(FDR\)](#). The key distinction between FDR control and FWER control is crucial: FWER attempts to ensure that the probability of making \*any\* false positive is low, whereas FDR controls the expected proportion of rejected null hypotheses that turn out to be false. Practically speaking, if a researcher utilizes the BH procedure and sets the acceptable FDR threshold at 10% (or  $Q = 0.10$ ), they are statistically ensuring that, on average, no more than 10% of their declared significant findings are likely to be incorrect.

This approach offers substantial benefits, particularly in high-dimensional data environments like bioinformatics, neuroimaging, or high-throughput drug screening. In these domains, scientists often accept that a small, controlled number of false positives is a necessary trade-off if it means avoiding the catastrophic error of missing many true discoveries (Type II errors). By controlling the FDR instead of the FWER, the [Benjamini-Hochberg Procedure](#) maximizes [statistical power](#) compared to stricter FWER corrections, thereby successfully navigating the critical challenge of balancing error control with the ability to identify genuine effects.

## Step-by-Step Implementation of the Benjamini-Hochberg Procedure

The application of the **Benjamini-Hochberg Procedure** is systematic and relies fundamentally on the ranking of p-values derived from the full set of statistical tests conducted. The procedure

establishes an adjusted critical value for each test based on its rank, ensuring that only the most compelling results are declared significant under the chosen FDR threshold.

**Step 1: Conduct Tests and Collect P-values.** First, the researcher executes all  $m$  statistical tests relevant to the study. For each test, the corresponding raw p-value must be calculated and recorded, resulting in a set of values denoted as  $P_1, P_2, \dots, P_m$ .

**Step 2: Rank the P-values.** All  $m$  p-values are then arranged in strict ascending order, from the smallest value to the largest. A sequential rank,  $i$ , is assigned to each value. The smallest p-value receives a rank of  $i=1$ , the next smallest receives  $i=2$ , and the largest p-value receives the final rank of  $i=m$ .

**Step 3: Calculate the BH Critical Value.** For every ranked p-value, a corresponding Benjamini-Hochberg critical value must be calculated. This calculation provides the unique, rank-adjusted significance threshold for that specific test, taking into account the multiplicity of tests performed. The formula used to calculate this critical value is:

$$(i / m) \times Q$$

Where the terms are defined as:

$i$  = the rank of the current p-value (1, 2, ...,  $m$ ).

$m$  = the total number of statistical tests performed.

$Q$  = the pre-selected maximum [False Discovery Rate](#) (e.g., 0.05 or 0.20).

**Step 4: Identify the Significance Cutoff.** The critical step involves iterating backward, starting from the largest p-value (rank  $m$ ) and moving toward the smallest (rank 1). The researcher identifies the largest p-value,  $P_i$ , that satisfies the condition of being less than or equal to its unique calculated BH critical value,  $(i / m) \times Q$ . This specific p-value establishes the final significance benchmark for the entire test set.

**Step 5: Declare Significance.** All p-values that are less than or equal to the p-value identified in Step 4 are declared statistically significant findings. This set of rejected null hypotheses includes the benchmark p-value itself and every p-value ranked smaller than it.

This procedure is characteristically sequential and inherently non-increasing in stringency. Since the critical value  $(i / m) \times Q$  linearly increases with the rank  $i$ , the procedure allows less stringent criteria for higher-ranked (larger) p-values. This is contingent on the fact that the most significant (smallest) p-values have already passed their respective, stricter tests. This sequential dependency is precisely what provides the BH procedure with its superior [statistical power](#) when contrasted with simpler adjustments like the Bonferroni correction.

## Practical Example: Applying the Procedure to Health Research

To demonstrate the practical utility of the BH procedure, consider a health research study where investigators examine 20 distinct lifestyle and biomarker variables for potential association with a specific chronic disease. This generates 20 separate statistical tests and 20 p-values. For this exploratory analysis, the researchers strategically choose a maximum acceptable **False Discovery Rate** (Q) of 20% (Q = 0.20), acknowledging the exploratory nature of the initial screen.

The initial step requires ranking the resulting p-values from the smallest (rank 1) to the largest (rank 20). The illustration below depicts the ranked p-values for the 20 variables investigated.

Variable	P-Value	Rank
Variable #12	0.001	1
Variable #17	0.034	2
Variable #3	0.035	3
Variable #11	0.039	4
Variable #1	0.094	5
Variable #10	0.106	6
Variable #4	0.145	7
Variable #13	0.145	8
Variable #16	0.223	9
Variable #9	0.241	10
Variable #18	0.315	11
Variable #20	0.338	12
Variable #7	0.432	13
Variable #14	0.478	14
Variable #6	0.497	15
Variable #8	0.566	16
Variable #15	0.654	17
Variable #2	0.876	18
Variable #5	0.905	19
Variable #19	0.965	20

Given that the total number of tests ( $m$ ) is 20 and the chosen FDR (Q) is 0.20, the formula for the Benjamini-Hochberg critical value for any rank  $i$  simplifies to:  $(i / 20) \times 0.20$ . This unique critical value establishes the specific benchmark against which the observed p-value for that rank must be evaluated.

The subsequent table displays the calculated Benjamini-Hochberg critical value alongside each ranked p-value, facilitating the comparison.

Variable	P-Value	Rank	$(i/m)*Q$
Variable #12	0.001	1	0.010
Variable #17	0.034	2	0.020
Variable #3	0.035	3	0.030
Variable #11	0.039	4	0.040
Variable #1	0.094	5	0.050
Variable #10	0.106	6	0.060
Variable #4	0.145	7	0.070
Variable #13	0.145	8	0.080
Variable #16	0.223	9	0.090
Variable #9	0.241	10	0.100
Variable #18	0.315	11	0.110
Variable #20	0.338	12	0.120
Variable #7	0.432	13	0.130
Variable #14	0.478	14	0.140
Variable #6	0.497	15	0.150
Variable #8	0.566	16	0.160
Variable #15	0.654	17	0.170
Variable #2	0.876	18	0.180
Variable #5	0.905	19	0.190
Variable #19	0.965	20	0.200

The critical cutoff point is determined by working backward from the largest rank (Variable #1, rank 20). We must identify the first instance where the observed p-value is less than or equal to its corresponding BH critical value. In this dataset, that crucial intersection occurs at Variable #11. Variable #11 has an observed p-value of 0.039, and its calculated B-H critical value is 0.040. Since 0.039 is less than or equal to 0.040, Variable #11 is declared statistically significant.

Following the final rule of the procedure, this cutoff variable and all tests exhibiting smaller p-values (i.e., those ranked 1 through 10) are definitively considered significant findings. The finalized set of significant results is highlighted in the table below.

	Variable	P-Value	Rank	$(i/m)*Q$
Significant	Variable #12	0.001	1	0.010
	Variable #17	0.034	2	0.020
	Variable #3	0.035	3	0.030
	Variable #11	0.039	4	0.040
	Variable #1	0.094	5	0.050
	Variable #10	0.106	6	0.060
	Variable #4	0.145	7	0.070
	Variable #13	0.145	8	0.080
	Variable #16	0.223	9	0.090
	Variable #9	0.241	10	0.100
	Variable #18	0.315	11	0.110
	Variable #20	0.338	12	0.120
	Variable #7	0.432	13	0.130
	Variable #14	0.478	14	0.140
	Variable #6	0.497	15	0.150
	Variable #8	0.566	16	0.160
	Variable #15	0.654	17	0.170
	Variable #2	0.876	18	0.180
	Variable #5	0.905	19	0.190
	Variable #19	0.965	20	0.200

It is instructive to note the sequential power of the BH procedure here. For example, Variable #17 (p-value 0.012, B-H critical value 0.014) and Variable #3 (p-value 0.021, B-H critical value 0.030) individually pass their rank-specific tests. However, their significance is ultimately derived from the fact that they have smaller p-values than the designated cutoff variable (Variable #11). If Variable #11 had failed the test (e.g., if its p-value was 0.041), the cutoff would have shifted to the next largest p-value that passed its individual test, thereby demonstrating the dynamic, sequential nature of FDR control.

## Strategic Selection of the False Discovery Rate (Q)

The selection of the maximum acceptable **False Discovery Rate (Q)** is not merely a statistical exercise; it represents the most important non-mathematical decision when deploying the [Benjamini-Hochberg Procedure](#). The chosen value of Q directly quantifies the researcher's tolerance for committing Type I errors within the set of hypotheses declared significant. To maintain the integrity and objectivity of the analysis, this threshold must be established *a priori*--before any data collection or statistical testing begins--thus preventing bias based on observed results.

The appropriate Q value is intensely dependent on the specific research context and the practical consequences associated with both false positives and false negatives. Researchers frequently utilize the BH procedure during the initial, exploratory phase of their work, aiming to cast a wide

investigative net to capture potential leads that will later be subjected to more rigorous, focused validation. If these subsequent validation tests are inexpensive, simple, or non-invasive, a higher  $Q$  (e.g., 0.10 or 0.20) may be justified. A higher FDR threshold maximizes [statistical power](#) and the likelihood of detecting true effects, operating under the assumption that the few resulting false discoveries will be quickly and cheaply eliminated during the low-cost validation stages.

Conversely, if the cost of basing further action on a **false discovery** is prohibitively high--such as initiating the development of an extremely expensive pharmaceutical based on a spurious genetic marker, or implementing a major, irreversible public policy change--then the researcher must mandate a much lower FDR (e.g.,  $Q = 0.01$  or  $0.05$ ). Furthermore, in fields where the risk of missing a critical, potentially life-altering discovery (a Type II error) is exceptionally high, researchers might deliberately accept a higher FDR to maximize sensitivity. Thus, the determination of  $Q$  is a crucial strategic judgment that meticulously balances the maximization of statistical discovery power against the practical implications and consequences of committing errors in that specific domain.

## Additional Resources for Further Study

For readers seeking a more profound understanding of the theoretical foundations and advanced applications of multiple testing corrections, including detailed comparisons between FWER and FDR methodologies, the following resources and topics are highly recommended:

Reviewing the original Benjamini and Hochberg (1995) paper that first proposed the [False Discovery Rate](#) control method.

Studying the precise mathematical and practical distinction between the [Family-Wise Error Rate \(FWER\)](#) and the False Discovery Rate (FDR).

Investigating the concept of the  $q$ -value, which serves as an adaptation of the  $p$ -value specifically utilized for robust FDR control.