

General use of multiple testing corrections in life sciences could boost replicability

Otília Menyhárt^{1,2,3}, Balázs Györfy^{1,2,3,4}

¹ Cancer Biomarker Research Group, Institute of Molecular Life Sciences, Hungarian Research Network, Budapest, 1117, Hungary

² National Laboratory for Drug Research and Development, Budapest, 1117, Hungary

³ Department of Bioinformatics, Semmelweis University, Budapest, 1094, Hungary

⁴ Department of Biophysics, Medical School, University of Pecs, Pecs, 7624, Hungary



Many studies fail to adjust for multiple hypothesis testing, resulting in irreproducible findings. We created a tool (www.multipletesting.com) to automate these corrections.

BACKGROUND

- Most published research findings may be false.
- Testing multiple hypotheses increases false positives, and uncorrected p-values can obscure real biological insights.
- Many studies, including clinical trials and epidemiology, fail to adjust for multiple testing.

Prevalence of adjustments (reviews)

Study	Year published	Number of studies with multiple comparisons	Proportion of studies with adjustments
Cohen	2010		22%
Tyler et al.	2011	>20	5.80%
Stacey et al.	2012	538	14%
Baron et al.	2013		40%
Wason et al.	2014		51%
Gewandter	2014	33	45%
Chalkidou	2014	15	1 study
Vickerstaff et al.	2015	60	25%
Kirkham et al.	2015	140	10%
Dworkin et al.	2016	29	21%
Benjamini and Cohen	2017		20%
Brand	2021	89	2 studies altogether
Nevin	2022	38	11% final reports, 7% protocols
Pike	2022	28	48%

➔ For 1000 variables, assume that all H0 are true, i.e. no real difference between control and experimental groups

p - values without correction:

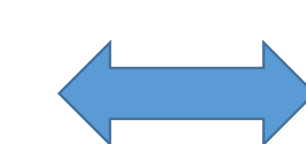
p = 0.05 results in 50 false positive

p = 0.01 results in 10 false positives

		actual	
		positive	negative
predicted	positive	true positive	false positive
	negative	false negative	true negative

BONFERRONI vs. FALSE DISCOVERY RATE?

Bonferroni



FDR

The original p value

$$\text{Bonferroni-corrected } p \text{ value} = \frac{\alpha}{n}$$

The number of tests performed

- strong control of type I error
- effective when a small number of hypotheses are tested
- high risk of false negatives
- drastically lowers statistical power

- controls the expected proportion of false positives
- higher power: more suitable for large datasets
- useful for exploratory research where some true effects are expected among many hypotheses.

when testing 20 variables

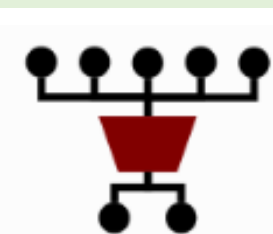
adjusted p 0.05/20 = 0.0025 - too stringent threshold!

FDR= 0.05 allows 5% of the discovered biomarkers to be false positives, which is more forgiving than the Bonferroni!

GOALS

to provide an automated interface for scientists to apply corrections for multiple hypothesis testing

multipletesting.com



Multiple Testing Correction

A tool for life science researchers for multiple hypothesis testing correction

What is multipletesting.com is useful for?

Conducting multiple statistical tests increases the likelihood that a significant proportion of associations will be false positives, clouding real discoveries. Several strategies exist to overcome the problem of multiple hypothesis testing. Our multiple testing correction tool provides the five most frequently used adjustment tools to solve the problem of multiple hypothesis testing, including the Bonferroni, the Holm (step-down), the Hochberg (step-up) corrections, and allows to calculate the False Discovery Rate (FDR) and q-values.

Using this multiple testing calculator is straightforward and user-friendly. It has never been easier to adjust p-values! Check out the list of possibilities for multiple hypothesis testing!

Multiple Testing

START

PERFORM MULTIPLE HYPOTHESIS TESTING USING A LIST OF P VALUES

PUBLICATION

READ OUR GUIDE TO MULTIPLE HYPOTHESIS TESTING

Analysis

Please enter (copy-paste) your p-values into the allotted space and select the relevant correction method(s). For more information please refer to our paper.

Step 1: Enter list of p-values:

Step 2: Compute following tests:

Set significance threshold at:

p = 0.05

p = 0.01

p = 0.001

user set:

Bonferroni

Holm

Hochberg

FDR

FDR = 10%

FDR = 5%

FDR = 1%

user set:

q value

RESULTS

First significant p-value (values over these thresholds are not significant):

Bonferroni

0.00166

Holm

0.00166

Hochberg

0.00365

FDR

0.00365

q values:

0.025186081

0.027689517

CONCLUSIONS

- our tool allows the **immediate application** of the most commonly used multiple testing correction methods
- **easy data upload**: copy and paste p-values
- allows **comparisons across adjustment methods**
- suitable for **all scientific disciplines**

Scan to see the site! ➔



Acknowledgments

This project was supported by the National Research, Development, and Innovation Office, Hungary (PharmaLab, Grant No.: RRF-2.3.1-21-2022-00015). Otília Menyhárt was supported by the Janos Bolyai Scholarship of the Hungarian Academy of Sciences and the Hungarian Scientific Research Fund (Grant No.: OTKA FK147194). The support of ELIXIR Hungary (www.bioinformatics.hu) is acknowledged. The authors declare no conflict of interest.

