

1

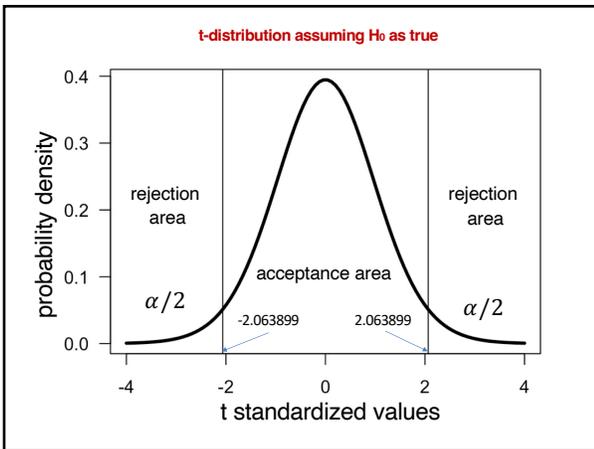
Why do we conduct ANOVA rather than perform multiple pairwise tests of means?

BIOL 422 & 680, Pedro Peres-Neto, Biology, Concordia University

A pedagogical guide for understanding the issues underlying Multiple hypothesis testing

Why should we not trust results obtained from multiple statistical tests?

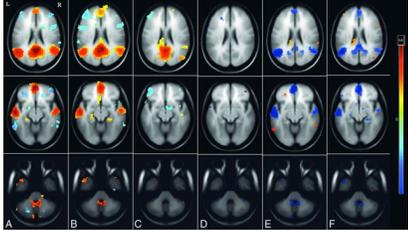
2



3

Examples of really huge numbers of multiple tests

Compare signal changes between task and no-task conditions using a t-test across more than 250,000 voxels (3D brain pixels).



Seizure Frequency Can Alter Brain Connectivity: Evidence from Resting-State fMRI

R.D. Bhardwaj, S. Sinha, R. Pandu, K. Rajhavanra, L. George, G. Chaitanya, A. Gupta, and P. Satishchandra

19

How to avoid inflated false positives (type I errors) due to multiple testing? Or the so-called family-wise error rate (FWER)

There is a large number of specific (e.g., Tukey-test for comparing two the difference between two means) and general procedures; the latter applying to any statistical test as they are used to control for multiple tests by correcting P-values.

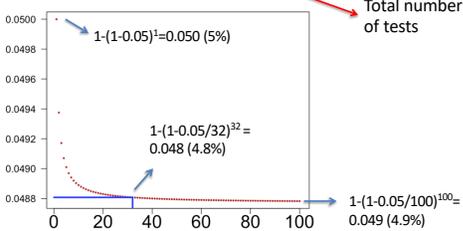
There are many commonly used procedures to correct for FWER; here we will review two (very commonly-used) general procedures:

- 1) Bonferroni correction (simplest): it controls the family Type I error.
- 2) False Discovery Rate (FDR; very much used these days): it controls the false discovery rate.

20

Bonferroni correction

Carlo Emilio Bonferroni developed the correction; modern use credited to Olive Dunn

$$\alpha_{Bonferroni} = \alpha/m = 0.05/32 = 0.0015625$$


Instead of using the original pre-established (desired) α , use α adjusted by the number of test instead to assure a family-wise (type I) error rate (FWER).

21

Bonferroni correction (common table presentation)

comparison	uncorrected P (t test)	Bonferroni P (t test)
control vs eyes	0.0029	0.0088
control vs knee	0.9418	1.0000
knee vs eyes	0.0044	0.0132

The Tukey test or Tukey's HSD (honest significant difference) usually taught in Intro stats

1) is a solution to correct for single two-sample t-tests.

2) It works well for small number of pairwise comparisons but not large.

25



26

When you claim discovery,
how often are you wrong?

To learn from many signals (p -values), it's like medical diagnosis: stricter criteria reduce false alarms but cause more real cases to be missed.

27

False Discovery Rates - FDR (or false positive rate)
When you claim discovery, how often are you wrong?

A strict control on Type I error comes at a cost (a *trade-off*): as we reduce the chance of false positives, we dramatically increase the chance of false negatives (Type II errors).

In other words, we protect ourselves from being wrong, but we also prevent ourselves from discovering real effects. This is why Bonferroni-type corrections are often said to reduce the power of discovery.

This situation is called a *trade-off* because improving one goal necessarily worsens another.

When we tighten our criteria to reduce false positives, we make it harder for results to be declared significant. This protects us from being wrong, but it also causes real effects to be missed more often.

We cannot minimize both errors at the same time, so reducing one inevitably increases the other.

28

False Discovery Rates - FDR (or false positive rate)
When you claim discovery, how often are you wrong?

To learn from many signals (p-values), it's like medical diagnosis: stricter criteria reduce false alarms but cause more real cases to be missed.

Bonferroni asks a stricter question than FDR: "What is the probability of making at least one false positive?"

FDR asks a more pragmatic and often more relevant question: "Among the results I am declaring as significant (my discoveries), what proportion are likely to be false positives (Type I error) do to multiple testing?"

The philosophy behind Bonferroni-type corrections is different and much stricter.

Rather than evaluating tests one by one, FDR asks whether the collection of reported discoveries is mostly correct or substantially contaminated by false positives, by estimating the proportion of false positives among them.

29

False Discovery Rates - FDR (or false positive rate)
versus Bonferroni (strict rules)

When conservation resources are limited, acting on false alarms for one population can reduce protection for many others—making strict control of false positives essential.

30

False Discovery Rates - FDR (or false positive rate) versus Bonferroni (strict rules)

To learn from many signals (p-values), it's like medical diagnosis: stricter criteria reduce false alarms but cause more real cases to be missed.

Goal: identify candidate genes potentially involved in a biological response, which will later be validated using independent experiments

Rather than evaluating tests one by one, FDR asks whether the collection of reported discoveries is mostly correct or substantially contaminated by false positives, by estimating the proportion of false positives among them.

31

Bonferroni versus FDR (quick contrast)

Number of significant tests after adjustment

Bonferroni = 0
FDR = 0

FDR logic: controlling the False Discovery Rate does not guarantee findings; it guarantees that any findings we choose to report are not expected to be heavily contaminated by false positives.

Bonferroni = 2
FDR = 1200

32

False Discovery Rates is widely used!

Methods in Ecology and Evolution

Methods in Ecology and Evolution 2011, 2, 278-282 doi: 10.1111/j.2041-210X.2010.00661.x

Using false discovery rates for multiple comparisons in ecology and evolution

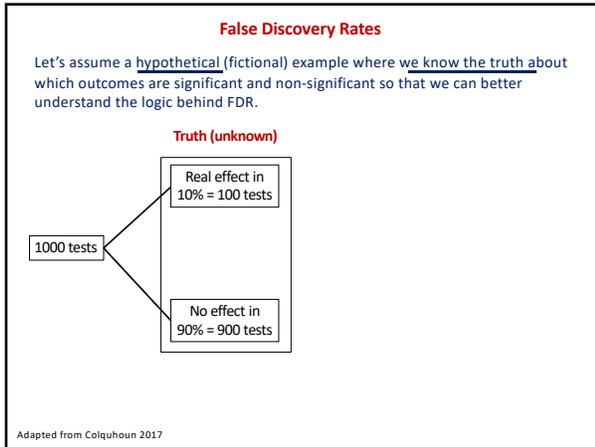
Nathan Pike*
Department of Zoology, University of Oxford, South Parks Road, Oxford OX1 3PS, UK

Statistical significance for genomewide studies

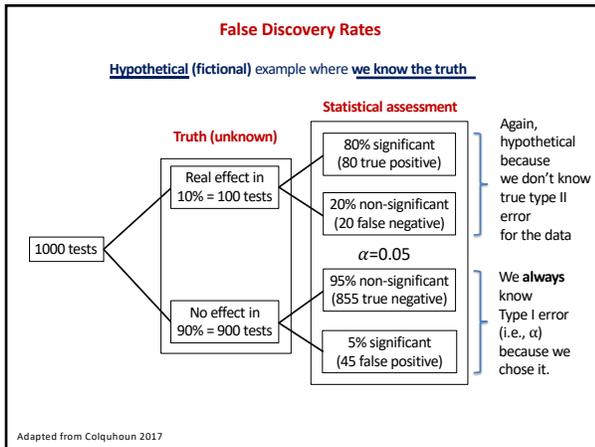
John D. Storey^{1*} and Robert Tibshirani²
¹Department of Biostatistics, University of Washington, Seattle, WA 98195; and ²Departments of Health Research and Policy and Statistics, Stanford University, Stanford, CA 94305

9440-9445 | PNAS | August 5, 2003 | vol. 100 | no. 16

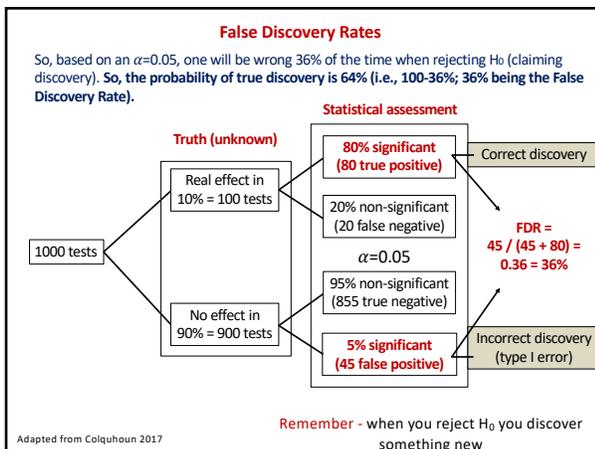
33



34



35



36

False Discovery Rates.- we are interested in positives (i.e., significant P-values) because those are “discoveries” – some likely wrong and some correct

Based on an $\alpha=0.05$, in this case, we will be wrong 36% of the time if we reject H_0 (claiming discovery). So, the probability of true discovery (reject a false H_0) is 64%.

The goal is to reduce the FDR to say 0.05 instead of keeping it at 0.36! So that the true discovery is higher (0.95 = 95%)

How to estimate FDR based on real data where we don't know the truth about false positives and negative as in this example?

Statistical assessment

80% significant (80 true positive)	Correct knowledge
20% non-significant (20 false negative)	
$\alpha=0.05$	
95% non-significant (855 true negative)	
5% significant (45 false positive)	Incorrect knowledge (type I error)

FDR =
 $45 / (45 + 80) =$
 $0.36 = 36\%$

Remember - when you reject H_0 you discover something new

37

False Discovery Rates.- we are interested in positives (i.e., significant P-values) because those are “discoveries” – some likely wrong and some correct

Interpretation: 36% is **not** the probability that any single result is wrong.

It is a **rate describing the whole set of discoveries:** if you were to repeat the study many times, about **36% (in average) of the results you call significant would be a false positive, on average.**

Statistical assessment

80% significant (80 true positive)	Correct knowledge
20% non-significant (20 false negative)	
$\alpha=0.05$	
95% non-significant (855 true negative)	
5% significant (45 false positive)	Incorrect knowledge (type I error)

FDR =
 $45 / (45 + 80) =$
 $0.36 = 36\%$

Remember - when you reject H_0 you discover something new

38

FDR then requires an estimate of the number of true positives!

Required knowledge (Step 1): Understand that when samples (e.g., control versus treatment) come from the same population (H_0 is true), the frequency distribution of P-values is flat (uniform).

≤ 0.05

> 0.05

Frequency distribution of infinite P-values generated by testing the difference between two samples (t-test) taken from the same population.

Each bin contains exactly 5% of P-values

Proportion of p-values

P-values

39

FDR then requires an estimate of the number of true positives!

Required knowledge (Step 1): Understand that when samples or groups (e.g., control versus treatment) come from the same population (i.e., H_0 is true), the frequency distribution of P-values is flat (uniform).

```

vector.pvalues <- matrix(0,1000)
for (i in 1:10000){
  x1 <- rnorm(20,5,2)
  x2 <- rnorm(20,5,2)
  vector.pvalues[i] <-
    t.test(x1, x2, alternative = "two.sided", var.equal = FALSE)$p.value
}
hist(vector.pvalues,ylim=c(0,1000),col="firebrick")
    
```

How to estimate FDR based on real data where we don't know the truth about false positives and negative as in this example?

40

FDR then requires an estimate of the number of true positives!

A simulation to show that when samples (e.g., control versus treatment) come from the same population (H_0 is true), the frequency distribution of P-values is flat (uniform).

Frequency distribution of 10,000 P-values generated by testing the difference between two samples (t-test) taken from the same population.

Each bin contains about 5% of P-values

41

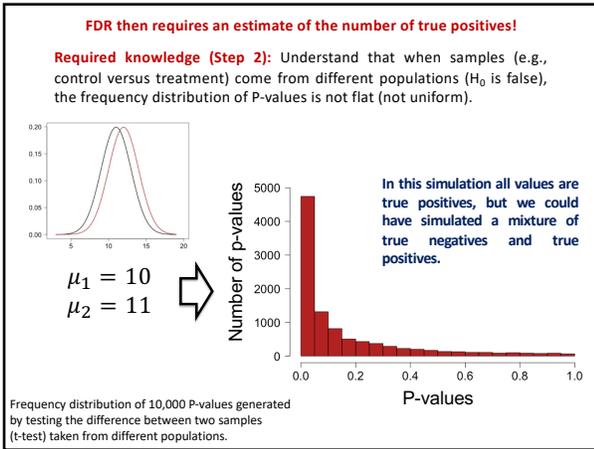
FDR then requires an estimate of the number of true positives!

Required knowledge (Step 2): Understand that when samples (e.g., control versus treatment) come from different populations (H_0 is false), the frequency distribution of P-values is not flat (not uniform).

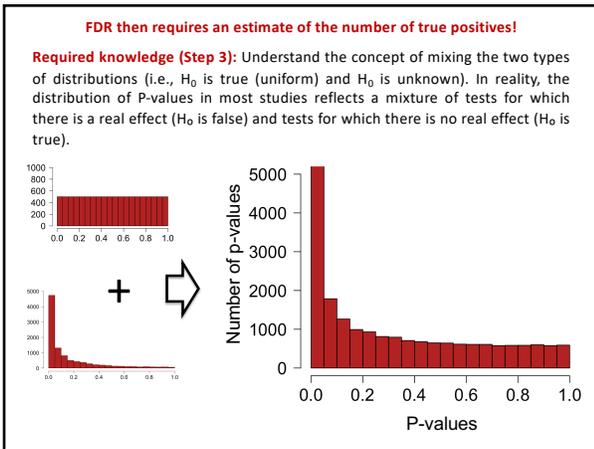
```

vector.pvalues <- matrix(0,1000)
for (i in 1:10000){
  x1 <- rnorm(20,10,2)
  x2 <- rnorm(20,11,2)
  vector.pvalues[i] <-
    t.test(x1, x2, alternative = "two.sided", var.equal = FALSE)$p.value
}
hist(vector.pvalues,ylim=c(0,1000),col="firebrick")
    
```

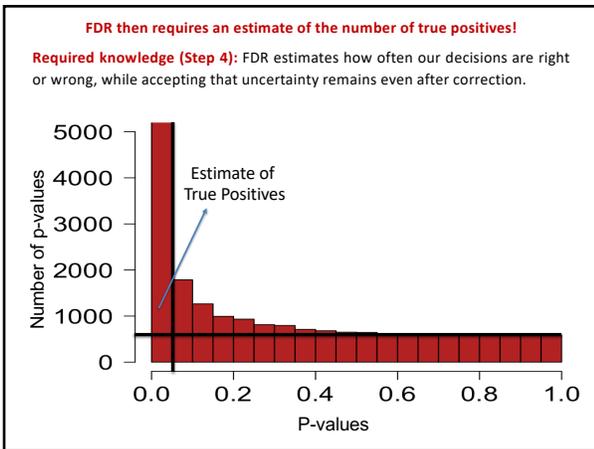
42



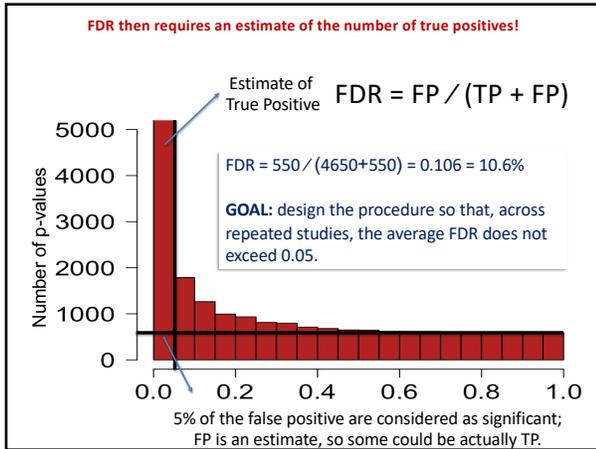
43



44



45



49

False Discovery Rates - FDR (or false positive rate)
How much did you learn that was false positive?

The are different types of FDR procedures and the one by Benjamini-Hochberg is likely the most commonly used! To correct the P-values based on the BH-FDR procedure, the calculation is conditional on previous P-values. R does it for you!!

Whatever the FDR is, the goal is to design the procedure so that, across repeated studies, the average FDR does not exceed 0.05.

50

Step 5: Adjust probabilities based on the FDR principle (NOT CRITICAL TO KNOW)

Rank (i)	Original P-value	BH multiplier (m / i)	BH-adjusted P-value
1	0.01	10.00	0.10
2	0.11	5.00	0.55
3	0.21	3.33	0.70
4	0.31	2.50	0.78
5	0.41	2.00	0.82
6	0.51	1.67	0.85
7	0.61	1.43	0.87
8	0.71	1.25	0.89
9	0.81	1.11	0.90
10	0.91	1.00	0.91

No significant p-value based on the FDR logic: controlling the False Discovery Rate does not guarantee findings; it guarantees that any findings we choose to report are not expected to be heavily contaminated by false positives.

51

Journal of Research on Educational Effectiveness, 5: 189-211, 2012
Copyright © Taylor & Francis Group, LLC
DOI: 10.1080/19347757.2011.614213

 **Some Bayesian dissent**

METHODOLOGICAL STUDIES

Why We (Usually) Don't Have to Worry About Multiple Comparisons

Andrew Gelman
Columbia University, New York, New York, USA

Jennifer Hill
New York University, New York, New York, USA

Masamichi Yajima
University of California, Los Angeles, Los Angeles, California, USA

Main issues from a Bayesian perspective (my summary):

- 1) FWER (family wise error, e.g., Bonferroni) is the general goal and this is an issue because it puts sole emphasis on Type I error (even FDR in many ways);
- 2) issues with dependent tests;
- 3) FDR good for very large number of tests but Bayesians may not recommend it for small numbers.

Bottom line: journals will request multiple testing and routine procedures are easier to implement and "articulate" than Bayesian ones. So...for the majority of scientists, Type I error is a really BIG ISSUE and needs to be dealt with using appropriate adjustments!

52

What should be corrected for?

- Variance and multiple t tests?
- All tests in a paper?
- All tests across all papers within a journal issue?
- All test across all papers within a year
- The world is the limit!

Look into this blog (*Why you don't need to adjust your alpha level for all tests you'll do in your lifetime*):

<http://daniellakens.blogspot.com/2016/02/why-you-dont-need-to-adjust-you-alpha.html>

I don't necessarily agree with everything in there, but good food for thought!

53