

Multiple Hypothesis Testing

BMI 206

| <u>P-VALUE</u> | <u>INTERPRETATION</u> |
|----------------|--|
| 0.001 | HIGHLY SIGNIFICANT |
| 0.01 | |
| 0.02 | |
| 0.03 | |
| 0.04 | SIGNIFICANT |
| 0.049 | |
| 0.050 | OH CRAP. REDO CALCULATIONS. |
| 0.051 | ON THE EDGE OF SIGNIFICANCE |
| 0.06 | |
| 0.07 | HIGHLY SUGGESTIVE, SIGNIFICANT AT THE $P < 0.10$ LEVEL |
| 0.08 | |
| 0.09 | |
| 0.099 | HEY, LOOK AT THIS INTERESTING SUBGROUP ANALYSIS |
| ≥ 0.1 | |

<https://xkcd.com/1478/>

What is a p-value?

- **p-value:** the probability of obtaining a result at least as extreme as observed if H_0 is true.
 - Null hypothesis (H_0) is usually: chance/no effect
- $P < 0.05$ does not necessarily indicate a **meaningful** difference.
- $P > 0.05$ does not necessarily indicate no **meaningful** difference.

Outcomes of One Test

| | | Reject | Do Not Reject |
|-------|-------|---------------------|---------------------|
| H_0 | True | False Positive (FP) | True Negative (TN) |
| | False | True Positive (TP) | False Negative (FN) |

Type I vs. Type II error

| | | Reject | Do Not Reject |
|-------|-------|--------------------|--------------------|
| H_0 | True | Type I Error | True Negative (TN) |
| | False | True Positive (TP) | Type II Error |

Life Hack: the boy who cried wolf



1. Caused a **type I error**:
 - townspeople thought there was a wolf when there was not (False Positive)
2. Then caused a **type II error**:
 - townspeople thought there was no wolf when there was (False Negative)

Controlling Errors in One Test

| | | Reject | Do Not Reject |
|-------|-------|---------------------|---------------------|
| H_0 | True | False Positive (FP) | True Negative (TN) |
| | False | True Positive (TP) | False Negative (FN) |

Significance Level (α) = P(FP)

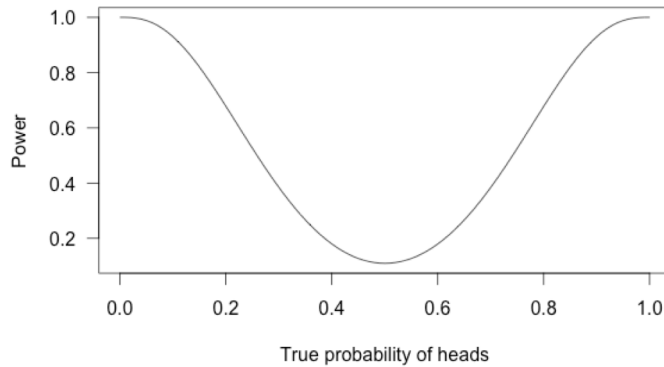
Power = $1 - P(\text{FN}) = 1 - \beta$

Statistical Power

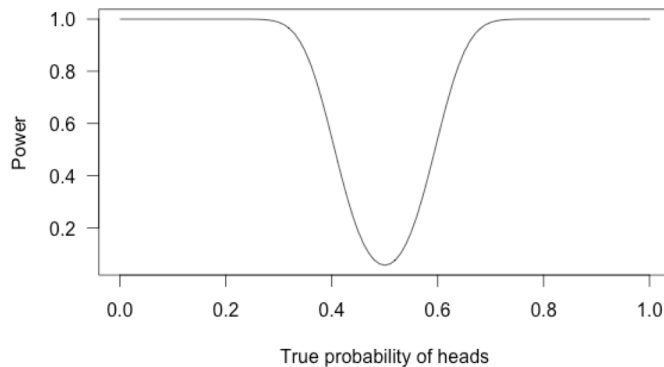
- The **power** of a test is the probability of rejecting a false null hypothesis ($1 - P(\text{FP})$)
- Power varies based on the **effect size** and the **sample size**.

Statistical Power

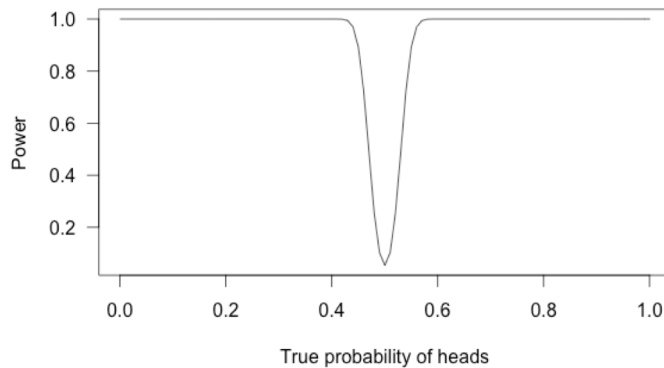
10 Flips



100 Flips



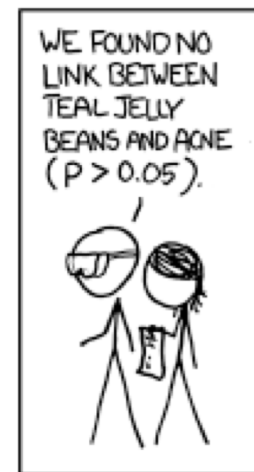
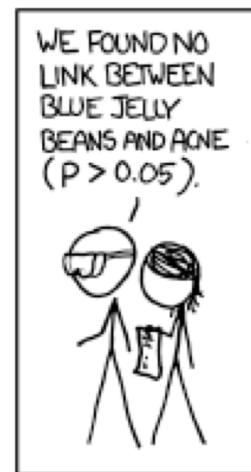
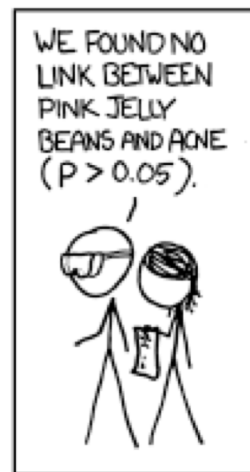
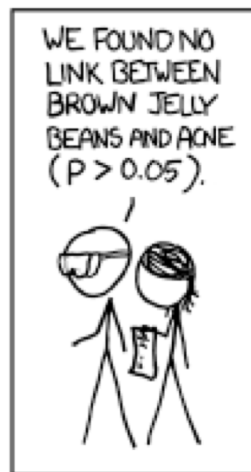
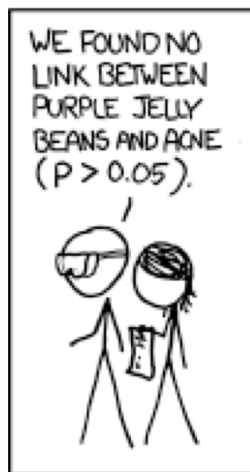
1000 Flips

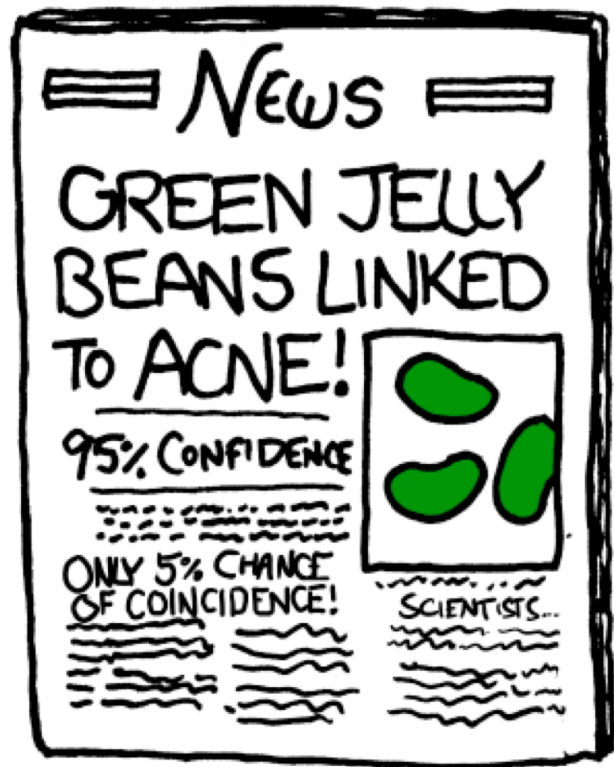
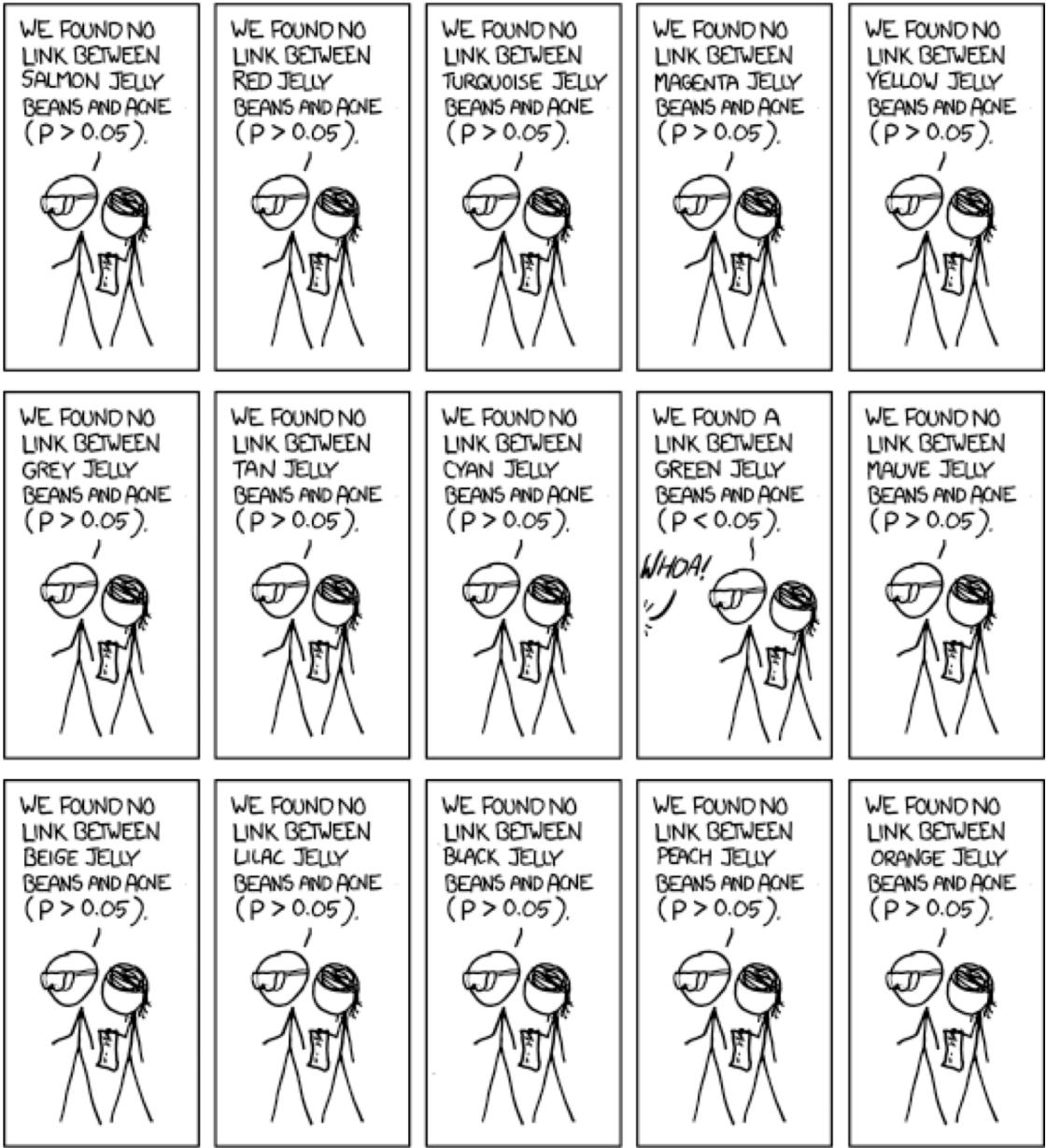


- Power increases with sample size.
- Power increases with effect size.
- Many studies are underpowered.

What happens when we test more than one hypothesis?

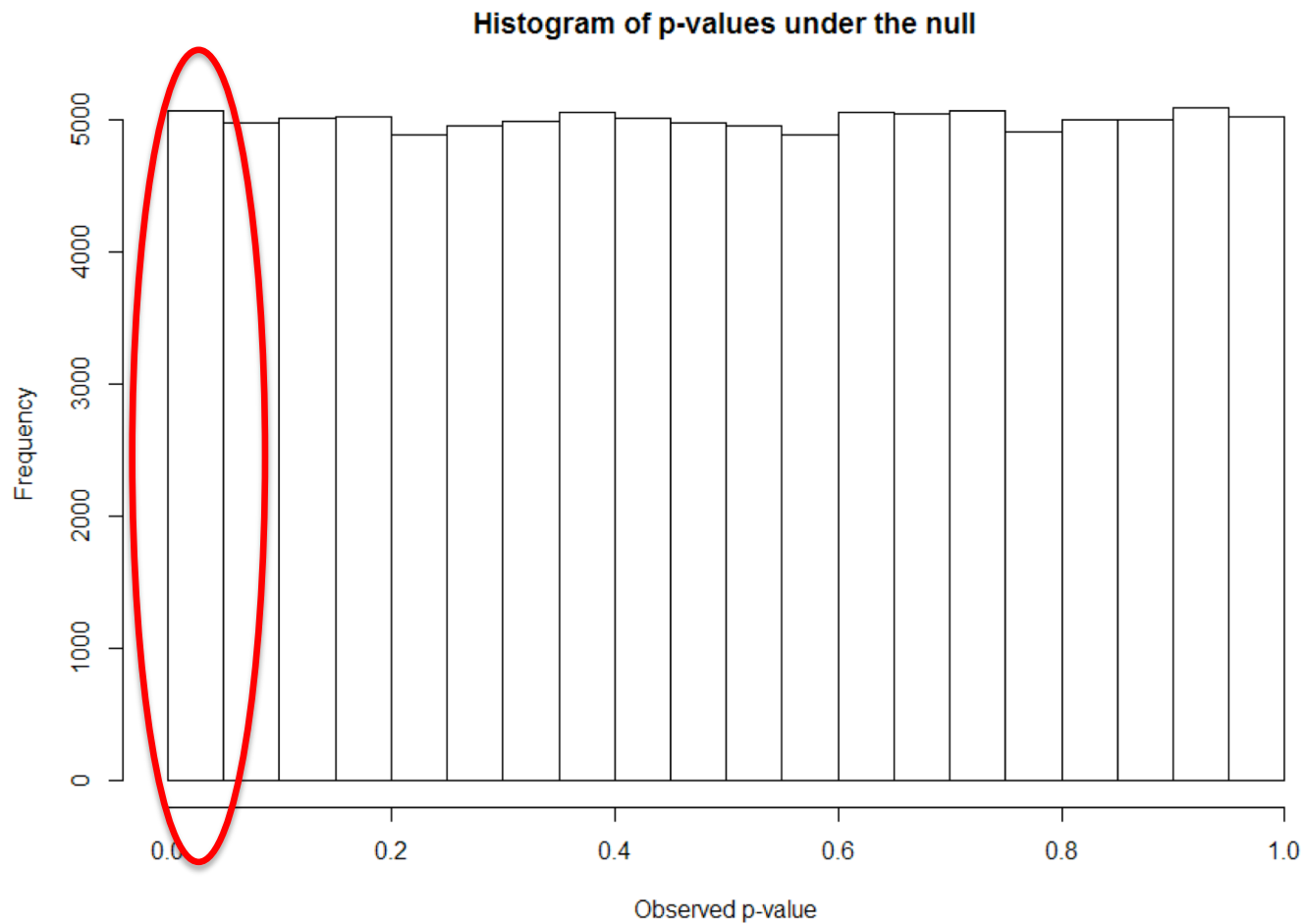
A motivational cartoon...





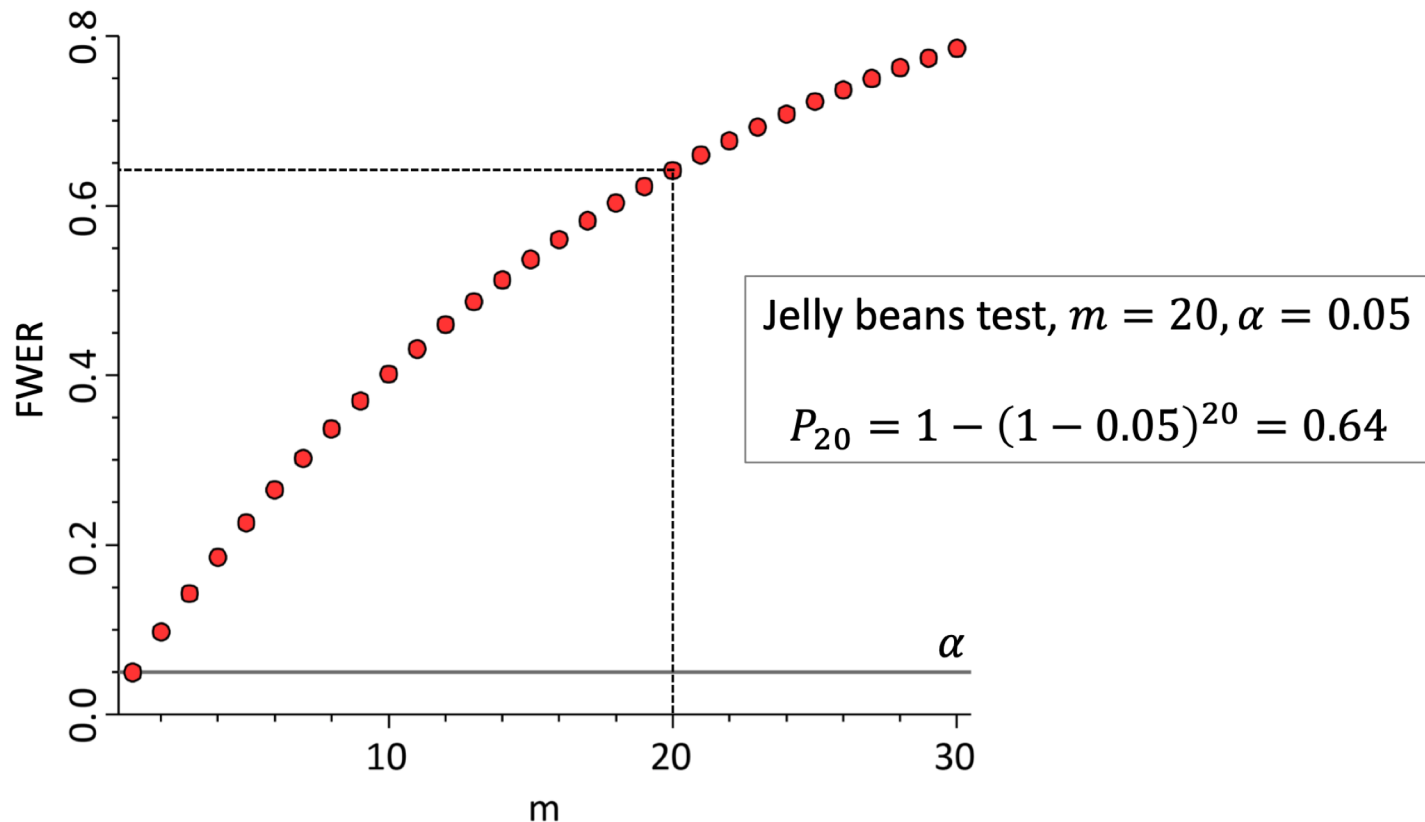
Why is testing multiple hypotheses a problem?

What is the distribution of p-values under the null?



What is the chance under the null of at least one p-value $< \alpha$ in m ind. tests?

$$1 - (1 - \alpha)^m$$



Outcomes of Many Tests

| | | Reject | Do Not Reject | |
|-------|-------|--------------------------|------------------------|--------------------------------|
| H_0 | True | # False Positives (FP) | # True Negatives (TN) | $m_0 = \# \text{ true nulls}$ |
| | False | # True Positives (TP) | # False Negatives (FN) | $m_1 = \# \text{ false nulls}$ |
| | | $r = \#$ reject nulls | $m-r$ | $m = \# \text{ tests}$ |

What can we do?

- In one test, α controls the family-wise error rate (**FWER**), the probability of at least one false positive :

$$P(\mathbf{FP} > 0) \leq \alpha$$

- Over all m tests, this is:

$$P(\mathbf{\#FP} > 0) \leq \alpha$$

Bonferroni Correction

To control FWER over m tests, adjust the p-value threshold (α) we use:

$$\alpha_{\text{Bonferroni}} = \alpha / m$$

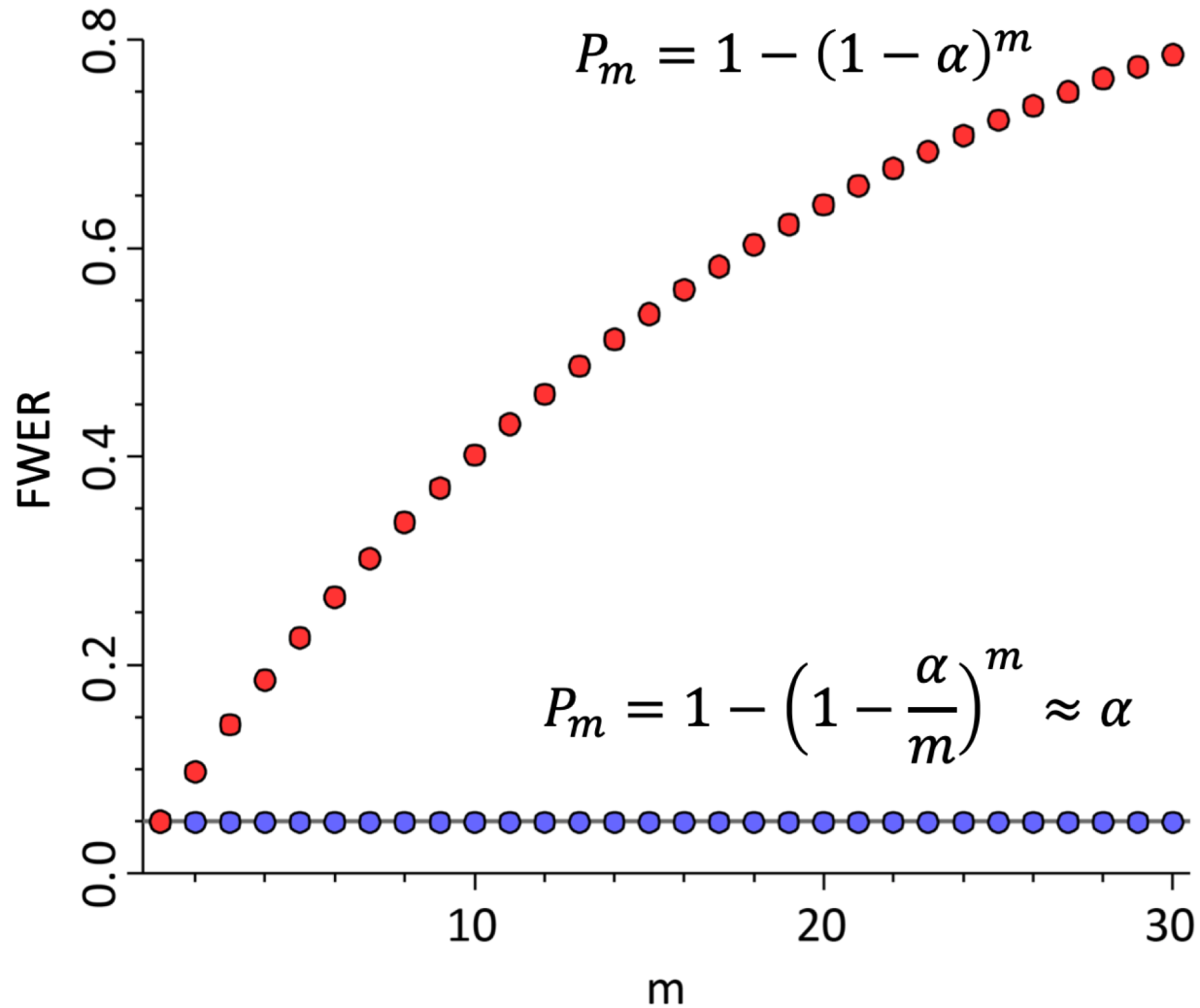
If $\alpha = .05$ and 20 tests:

$$\alpha_{\text{Bonferroni}} = 0.05 / 20 = 0.0025$$

Or, equivalently, correct the p-values:

$$p_{\text{Bonferroni}} = p * 20$$

Bonferroni Correction Graph



Proof

- Let p_1, \dots, p_m be the p-values for all tests
- Let I_0 be the set of all m_0 true null hypotheses
- We are interested in:

$$P(p_i \leq \frac{\alpha}{m}) \text{ for at least one } i \text{ in } I_0$$

- By Boole's inequality, this is \leq :

$$\sum_{i \in I_0} P(p_i \leq \frac{\alpha}{m}) = \sum_{i \in I_0} \frac{\alpha}{m} = \frac{m_0 \alpha}{m} \leq \alpha$$

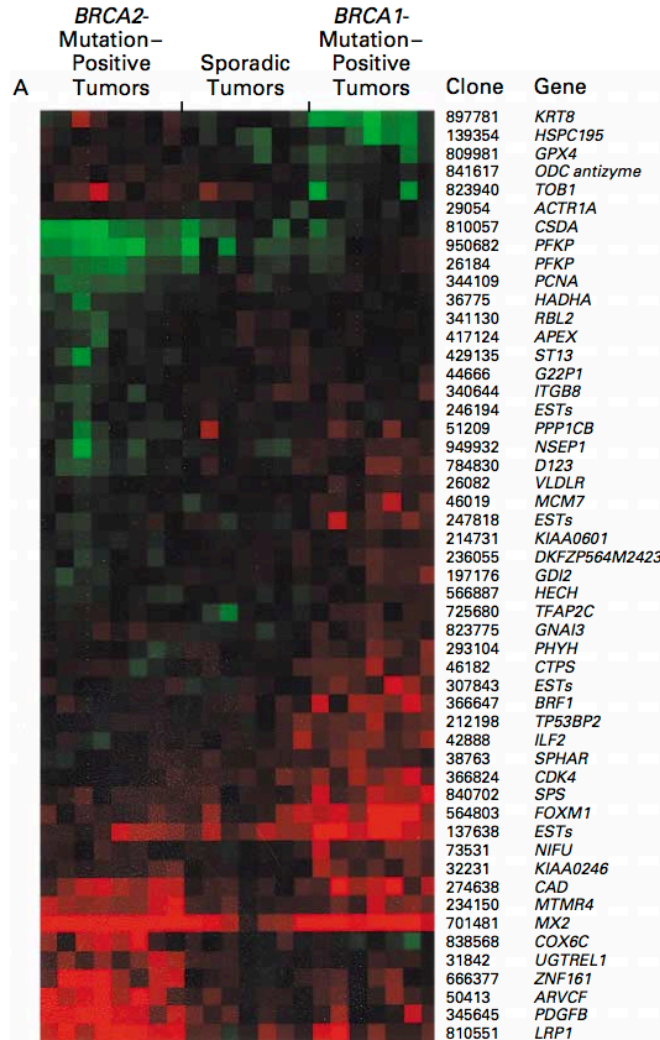
Problems with Bonferroni

- Bonferroni correction is conservative
 - Can use Holm-Bonferroni instead:

$$P_k < \frac{\alpha}{m + 1 - k}$$

- Bonferroni says little about the mix of **TPs** and **FPs** in the set of hypotheses called significant.
- If we expect that many tests should reject H_0 , we may be fine with more than one **FP**.

Genome-wide Analyses



Many genes are likely to be differentially expressed between conditions.

Why not control # FPs in tests called significant?

False Discovery Rate (FDR)

$$FP / (FP + TP)$$

vs.

False Positive Rate (FPR)

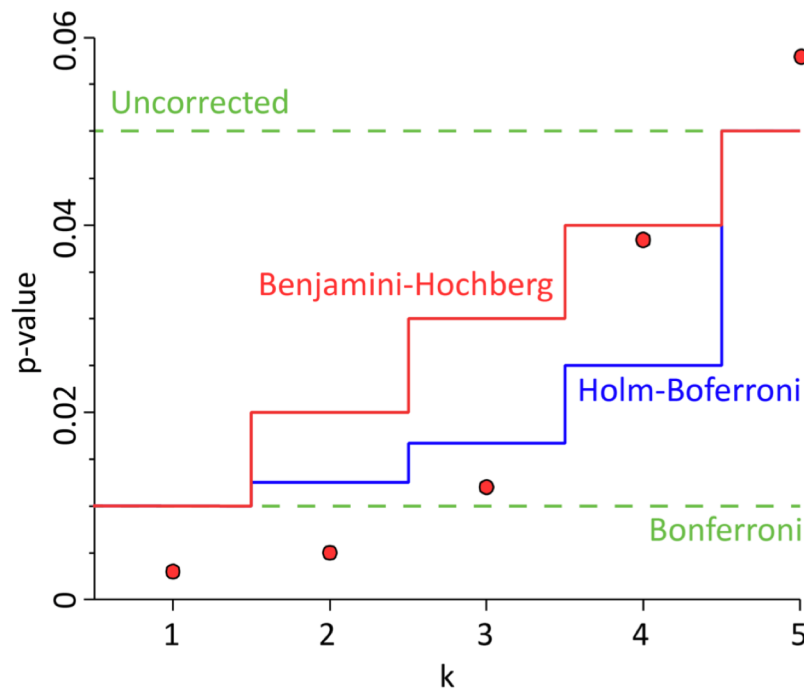
$$FP / (FP + TN)$$

| | | Reject | Do Not Reject |
|----------------|-------|------------------------|------------------------|
| H ₀ | True | # False Positives (FP) | # True Negatives (TN) |
| | False | # True Positives (TP) | # False Negatives (FN) |

q-value: the FDR analog of the p-value

Benjamini-Hochberg Procedure

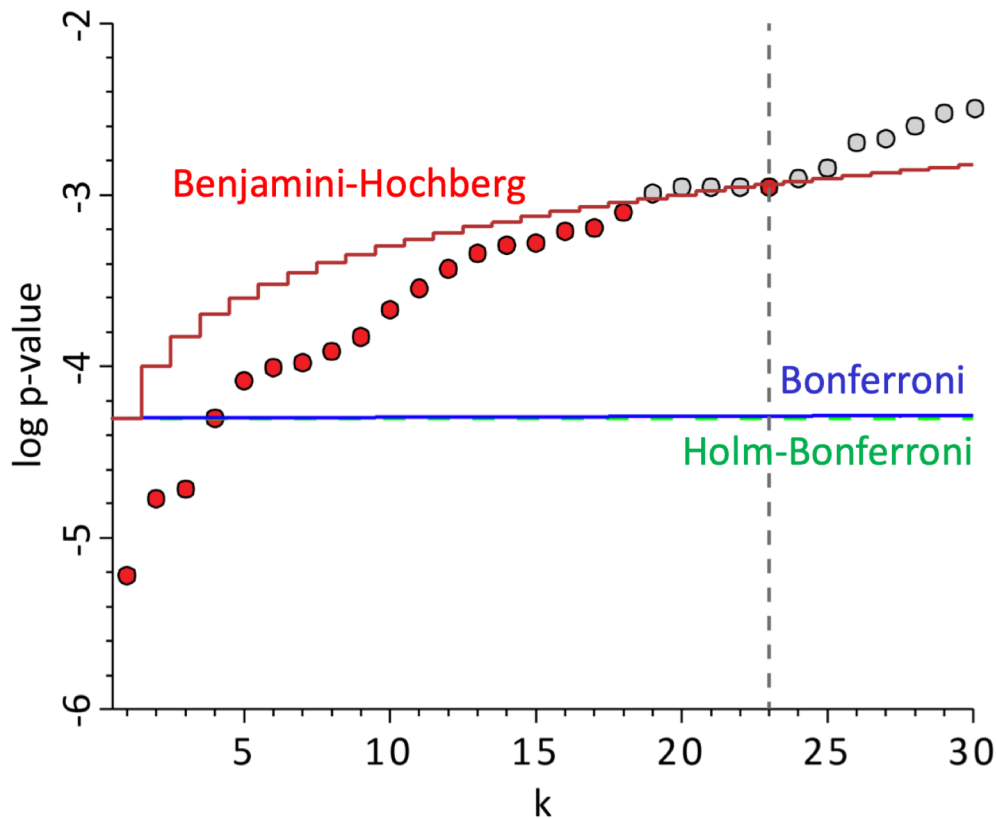
1. Rank p-values in ascending order: $P_{(1)} \dots P_{(m)}$.
2. For a given α , find largest k such that $P_{(k)} \leq \frac{k}{m} \alpha$.
3. Reject the null for all $H_{(i)}$ for $i = 1, \dots, k$.



Benjamini-Hochberg Procedure

1. Rank p-values in ascending order: $P_{(1)} \dots P_{(m)}$.
 2. For a given α , find largest k such that $P_{(k)} \leq \frac{k}{m} \alpha$.
 3. Reject the null for all $H_{(i)}$ for $i = 1, \dots, k$.
- BH procedure is less conservative than Bonferroni correction.
 - In genomics, we often expect many rejections of the null and can tolerate a few false positives.

BH Graphical Example



| Benjamini-Hochberg | | | |
|--------------------|---------------------|----------------------|-------|
| | H ₀ true | H ₀ false | Total |
| Significant | 2 | 21 | 23 |
| Not significant | 968 | 9 | 977 |
| Total | 970 | 30 | 1000 |

Other useful metrics

Sensitivity, Recall, True Positive Rate

$$TP / (TP + FN)$$

Specificity, True Negative Rate

$$TN / (TN + FP)$$

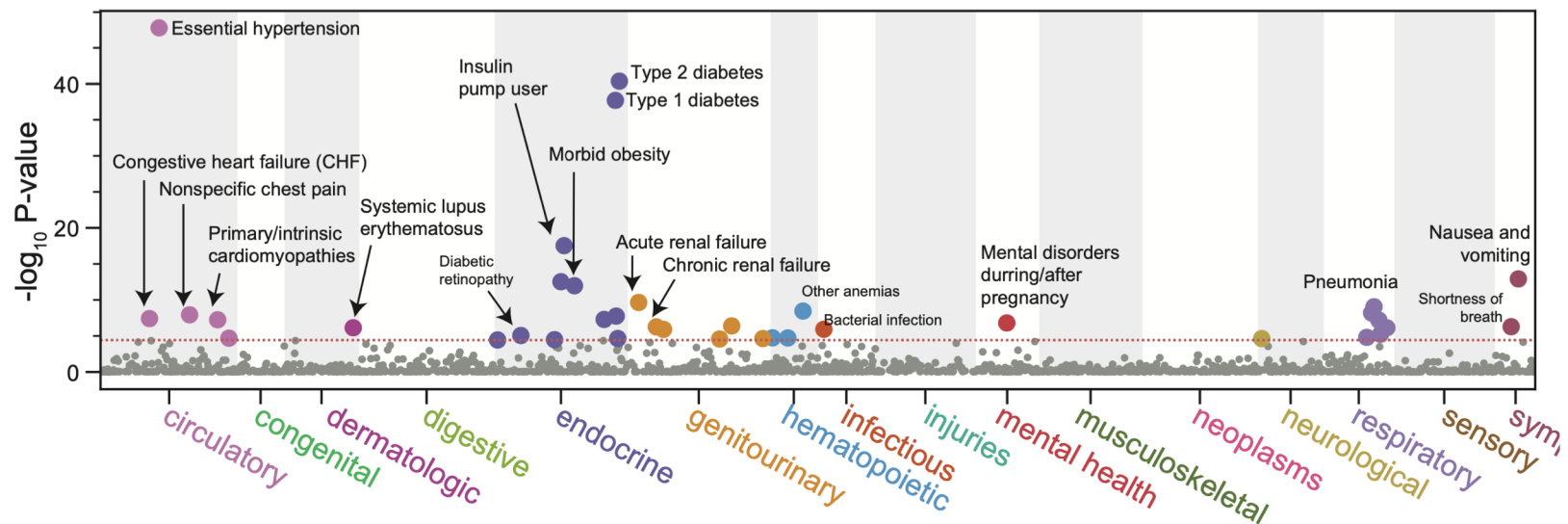
Precision, Positive Predictive Value

$$TP / (TP + FP)$$

| | | Reject | Do Not Reject |
|----------------|-------|-----------------------|------------------------|
| | | True | # False Positives (FP) |
| H ₀ | False | # True Positives (TP) | # False Negatives (FN) |

Discussion

1. How can we account for correlation structure among the results of our multiple tests?



2. Should you perform multiple testing correction for all the hypotheses you test in your life?

Tomorrow

- Examples of permutation and bootstrap methods for jointly adjusting for multiple testing

[Home](#) » [Bioconductor 3.14](#) » [Software Packages](#) » [multtest](#)

multtest

platforms [all](#) rank [44 / 2083](#) support [0 / 0](#) in Bioc > [16.5 years](#)
build [warnings](#) updated [before release](#) dependencies [14](#)

DOI: [10.18129/B9.bioc.multtest](#) [f](#) [t](#)

Resampling-based multiple hypothesis testing

Bioconductor version: Release (3.14)

Non-parametric bootstrap and permutation resampling-based multiple testing procedures (including empirical Bayes methods) for controlling the family-wise error rate (FWER), generalized family-wise error rate (gFWER), tail probability of the proportion of false positives (TPFP), and false discovery rate (FDR). Several choices of bootstrap-based null distribution are implemented (centered, centered and scaled, quantile-transformed). Single-step and step-wise methods are available. Tests based on a variety of t- and F-statistics (including t-statistics based on regression parameters from linear and survival models as well as those based on correlation parameters) are included. When probing hypotheses with t-statistics, users may also select a potentially faster null distribution which is multivariate normal with mean zero and variance covariance matrix derived from the vector influence function. Results are reported in terms of adjusted p-values, confidence regions and test statistic cutoffs. The procedures are directly applicable to identifying differentially expressed genes in DNA microarray experiments.

Author: Katherine S. Pollard, Houston N. Gilbert, Yongchao Ge, Sandra Taylor, Sandrine Dudoit

Maintainer: Katherine S. Pollard <katherine.pollard at gladstone.ucsf.edu>

Citation (from within R, enter `citation("multtest")`):

Pollard KS, Dudoit S, van der Laan MJ (2005). *Multiple Testing Procedures: R multtest Package and Applications to Genomics, in Bioinformatics and Computational Biology Solutions Using R and Bioconductor*. Springer.