How Not to be Too Wrong? An Introduction to the Multiple Testing Problem

Zhang Ruiyang

University College London

ruiyang.zhang.20@ucl.ac.uk

29 Nov 2021

This talk will be on the Multiple Testing Problem, which occurs when we are performing simultaneous hypothesis testing on the same set of data. The issue gets increasingly severe as the number of tests increases. The problem is hard to neglect, and is common in real life with examples including genetic research like GWAS and AB Testing for web designs. This talk will go over some basic remedies to the problem - the FWER controls like Bonferroni correction and the FDR controls like the BH procedure. It is self-contained and requires no prerequisite, but it would be helpful to have some knowledge of basic Statistics such as hypothesis testing and significance level, although they will be covered too. Setting: 8 cups - 4 milk then tea; 4 tea then milk

Experiment: Taste each cup, indicate whether it is milk then tea or tea then milk

Hypothesis: No such ability.

Significance Level: 5%

Prob. of 4 Correct Guesses at random: $1 \div {\binom{8}{4}} \approx 0.014 < 5\%$ Prob. of 3+ Correct Guesses at random: $(1 + {\binom{4}{3}} {\binom{4}{1}}) \div {\binom{8}{4}} \approx 0.24 > 5\%$ Hypothesis: One (null) hypothesis in mind.

p-value: Probability p of having the given data assuming the null hypothesis is true.

Significance Level: A pre-determined, reasonable yet arbitrary threshold α to reject the null hypothesis.

Outcome: (1) $p < \alpha$, reject null (2) Otherwise, not reject null

イロト イヨト イヨト イヨ

Hypothesis: Null Hypothesis H_0 , the claim that we aim to reject, and Alternative Hypothesis H_1 , some different claim.

p-value: Probability p of having the given data assuming the null hypothesis is true.

Significance Level: Predetermined. Probability of rejecting the null hypothesis when the null hypothesis is actually true.

Outcome: (1) $p < \alpha$, reject null (2) Otherwise, not reject null

	H_0 is true	H_0 is false
	Type I Error	Power
Reject H ₀	False Positive	True Positive
	α	$1 - \beta$
Not Reject H ₀	True Negative	Type II Error
	$1 - \alpha$	False Negative
	$1 - \alpha$	β

Under true null, the random variable **p-value** follows Uniform(0, 1).

If $\alpha = 0.05$, we will make a false positive on average for every 20 true null tests.

20 true nulls, 1 false positive on average. 1000 true nulls, 50 false positives on average.

Given *m* null hypotheses denoted by H_1, H_2, \cdots, H_m :

	True Null Hypothesis	False Null Hypothesis	Total
Reject	V	S	R
Not Reject	U	Т	m - R
Total	m_0	$m - m_0$	т

V, S, U, T are random variables, R, m, m_0 are known.

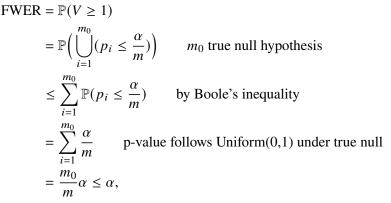
Family-wise Error Rate: FWER = $\mathbb{P}(V \ge 1) = 1 - \mathbb{P}(V = 0)$.

• • • • • • • • • • • • • •

Bonferroni Correction

Boole's Inequality: For a countable set of events A_1, A_2, A_3, \dots , $\mathbb{P}(\bigcup_i A_i) \leq \sum_i \mathbb{P}(A_i)$, due to sub-additivity property of probability measure.

With this inequality, since we are rejecting each null hypothesis H_i when its p-value $p_i \leq \frac{\alpha}{m}$, we would have



which indicates that the FWER is controlled under level α ?

Zhang Ruiyang (UCL)

An Introduction to the Multiple Testing Problem

(Holm, 1979)

For the *m* null hypotheses H_1, H_2, \dots, H_m , we compute their respective p-values P_1, P_2, \dots, P_m and we rank them such that $P_{(k)}$ denotes the *k*-th smallest p-value. So, $P_{(1)} \leq P_{(2)} \leq \dots \leq P_{(m)}$, and we denote the corresponding null hypotheses as $H_{(1)}, H_{(2)}, \dots, H_{(m)}$. We would want to control the FWER at α .

• Is $P_{(k)} \leq \frac{\alpha}{m-k+1}$? If so, reject $H_{(k)}$ and continue. Otherwise, EXIT.

To show that this method does in fact keep the FWER at α , we let I_0 be the set of indices corresponding to the true null hypothesis. This set is unknown to us and has size m_0 . Let us assume that we make the first false positive decision at $H_{(h)}$. Based on the procedures stated above, all the decisions for null hypothesis $H_{(1)}, H_{(2)}, \dots, H_{(h-1)}$ are true positives. Also, we know for a fact that $h - 1 \le m - m_0$ due to the definition of $m - m_0$. This implies that $m - h + 1 \ge m_0$, and then $\frac{1}{m-h+1} \le \frac{1}{m_0}$. Now, since $H_{(h)}$ is rejected, we would have $P_{(h)} \le \frac{\alpha}{m-h+1}$ by definition, so we will then have $P_{(h)} \le \frac{\alpha}{m-h+1} \le \frac{\alpha}{m_0}$. This means, if there is any false positive, we have at least one true null hypothesis with p-value less than $\frac{\alpha}{m_0}$.

FWER =
$$\mathbb{P}(V \ge 1)$$

= $\mathbb{P}\left(\bigcup_{i \in I_0} (P_i \le \frac{\alpha}{m_0})\right)$
 $\le \sum_{i \in I_0} \mathbb{P}(P_i \le \frac{\alpha}{m_0})$ by Boole's inequality
= $m_0 \frac{\alpha}{m_0}$ p-value follows Uniform(0,1) under true null
= α ,

which indicates that the FWER is controlled at level α .

Holm-Bonferroni Method is uniformly more powerful than Bonferroni.

	True Null Hypothesis	False Null Hypothesis	Total
Reject	V	S	R
Not Reject	U	Т	m-R
Total	m_0	$m - m_0$	т

(Benjamini & Hochberg, 1995)

Define a new random variable Q = V/(V + S) where Q = 0 when V + S = 0. This is the proportion of the false positives over all the rejected null hypotheses. This is unobservable since we do not know V, or S, or their realisations v or s. We will define the False Discovery Rate, or **FDR**, as the expectation of Q,

$$FDR = E[Q] = E[V/(V+S)] = E[V/R].$$

To avoid the division by zero issue, we would have the alternative formula for FDR as

$$FDR = E[V/R|R > 0]\mathbb{P}(R > 0).$$

(1) If all the null hypotheses are true, FDR = FWER. When s = 0 and v = r, Q = 0 if v = 0 and Q = 1 if v > 0, which means $\mathbb{P}(V \ge 1) = \mathbb{E}[Q]$. This means a control of FDR is a control of FWER in the weak sense.

(2) When $m_0 < m$, FDR is no bigger than FWER. Given $m_0 < m, v > 0$ implies $v/(v + s) \le 1$, which means $V \ge 1 \implies V \ge v/(v + s)$ and $\mathbb{P}(V \ge 1) \ge \mathbb{E}[Q]$. This means a control of FWER will control FDR.

For the *m* null hypotheses H_1, H_2, \dots, H_m , we compute their respective p-values P_1, P_2, \dots, P_m and we rank them such that $P_{(k)}$ denotes the *k*-th smallest p-value. So, $P_{(1)} \leq P_{(2)} \leq \dots \leq P_{(m)}$, and we denote the corresponding null hypotheses as $H_{(1)}, H_{(2)}, \dots, H_{(m)}$. We let the level of FDR that we would want to control at as α . The procedure works as the following:

- Let k be the largest i for which $P_{(i)} \leq \frac{i}{m} \alpha$.
- Reject all $H_{(i)}$ where $i = 1, 2, \cdots, k$.

Theorem

For **independent** test statistics and for any configuration of false null hypotheses, the above precedure controls the FDR at α .

• • • • • • • • • • • • • •

- Weaken Independence Criteria for BH Procedure
- FDR Control using Empirical Bayes methods
- Online Multiple Testing
- Knockoff for Variable Selection

- Benjamini Y., Hochberg Y. (1995) Controlling the False Discovery Rate: a Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society, Series B*, **57** No. 1, 289-300.
- Benjamini Y. (2010) Discovering the False Discovery Rate. *Journal of the Royal Statistical Society, Series B*, **72**, Part 4, 405-416.
- Holm S. (1979) A Simple Sequentially Rejective Multiple Test Procedure. *Scandinavian Journal of Statistics*. **6**, 65-70.