## Benjamini-Hochberg method (1995) for "false discovery rate",
http://www.jstor.org/stable/2346101 ;
implemented in Statistical Analysis of Microarrays (SAM) and BRBTOOLS from the National Cancer Institute's Biometric Research Branch.

*Number of errors committed when testing m null hypotheses*

|  | Declared non-significant | Declared significant | Total |
|---|---|---|---|
| True null hypotheses | U | V | $m_0$ |
| Non-true null hypotheses | T | S | $m - m_0$ |
|  | $m - R$ | R | $m$ |

B-H is a classical frequentist technique. False discovery rate could be defined as
$$\text{FDR} = Q_e = E(Q) = E\{V/(V + S)\} = E(V/R).$$
Here $R$ = # declared significant.
What if $R$ can equal 0? The $Q_e$ will be infinite.
And when all null hypotheses are true, then all discoveries are false and FDR=1.
Instead B&H define

$$\text{FDR} = P(R > 0)\, E(V/R \,|\, R > 0).$$

The BH procedure is:

$$\text{let } k \text{ be the largest } i \text{ for which } P_{(i)} \leqslant \frac{i}{m} q^*;$$

$$\text{then reject all } H_{(i)} \; i = 1, 2, \ldots, k.$$

Then FDR is no bigger than $q^*$.

Example: An randomized clinical trial of treatments rt-PA versus APSAC when myocardial infarction occurs. There are 15 different clinical endpoints.
The 15 P values are ranked. Here are the first few with the critical cutoff values with $q^*$=0.05.

|  | *1* | *2* | *3* | *4* |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|
| **Pvalue** | **0.0001** | **0.0004** | **0.0019** | **0.0095** | 0.0201 | 0.0278 | 0.0298 | ... | 1 |
| **cutoff** | **0.0033** | **0.0067** | **0.0100** | **0.0133** | 0.0167 | 0.0200 | 0.0233 | ... | 0.05 |

The first 4 hypotheses are "significant" with this rule. (In bold face)

Suppose we tweak the results. Decrease the 4th P value just a little:

|  | *1* | *2* | *3* | *4* | *5* |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|
| **Pvalue** | **0.0001** | **0.0004** | **0.0019** | **0.0095** | **0.0151** | 0.0278 | 0.0298 | ... | 1 |
| **cutoff** | **0.0033** | **0.0067** | **0.0100** | **0.0133** | **0.0167** | 0.0200 | 0.0233 | ... | 0.05 |

Great, now the 5th is significant.

Now increase the 4th P value slightly:

|  | *1* | *2* | *3* | *4* |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|
| **Pvalue** | **0.0001** | **0.0004** | **0.0019** | 0.0135 | 0.0151 | 0.0278 | 0.0298 | ... | 1 |
| **cutoff** | **0.0033** | **0.0067** | **0.0100** | 0.0133 | 0.0167 | 0.0200 | 0.0233 | ... | 0.05 |

Suddenly the 5th is ALSO not significant. Why not? Does this make sense? The data for the 5th test has not changed at all. That's typical that weird things happen in frequentist approaches to multiple comparisons. It's unsatisfying to many people.