

Correcting false discovery rates for their bias toward false positives

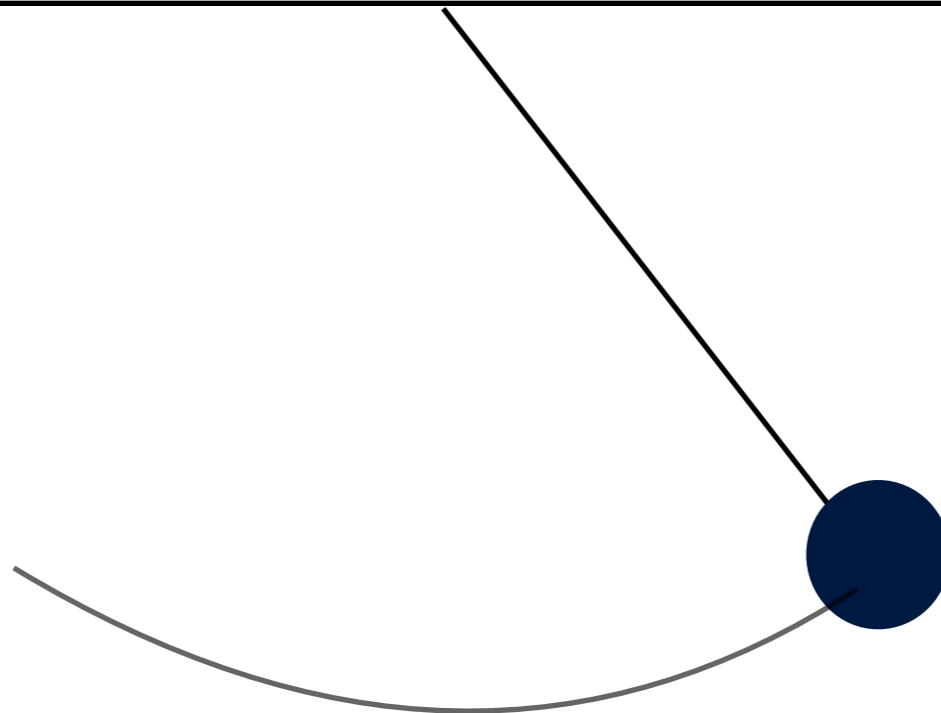
2017 Annual Meeting of the Statistical Society of Canada
University of Manitoba, Winnipeg

13 June 2017

David Bickel
University of Ottawa



Swung too far?



Family-wise error rates

Too many false negatives

False discovery rates

Too many false positives

Local false discovery rates

Large variance

The rise of false discovery rates

FDR software used in genomics

- Desktop software:
 - GSEA
 - Cyber-T
 - MACS

Why the FDR became popular

- Web software:
 - DAVID
 - Topppfunn
 - GREAT

Approach	FWER control	FDR control
Significance measure	Adjusted p-value	q-value
Interpretation	<i>(Achieved FWER)</i>	<i>(Achieved FDR)</i>
Decisions	<i>(Rejection only)</i>	<i>(Rejection only)</i>
Prior distribution	None	None
Many tests	<i>(Highly conservative)</i>	Adequate
Few tests	Adequate	Adequate

Multiple comparisons

Table 1 Summary of strengths and weakness of the four major approaches to multiple hypothesis testing.

Approach	Error-rate control approaches		Posterior probability approaches	
	FWER control	FDR control	Classical Bayes	Empirical Bayes
Significance measure	Adjusted p-value	q-value	LFDR	Estimated LFDR
Interpretation	<i>(Achieved FWER)</i>	<i>(Achieved FDR)</i>	Level of belief	Estimated prob.
Decisions	<i>(Rejection only)</i>	<i>(Rejection only)</i>	Optimal, flexible	<u>Flexible</u>
Prior distribution	None	None	<i>(Specified)</i>	<u>Estimated</u>
Many tests	<i>(Highly conservative)</i>	Adequate	Adequate	Adequate
Few tests	Adequate	Adequate	Adequate	<i>(Estimation error)</i>

Each of the last five rows has practical advantages and disadvantages of each approach according to the consideration given in the first column. A **bold entry** means the approach of a column is among the best for the consideration, an underlined entry means it is advantageous but notably less so, and an *(italicized entry in parentheses)* means it is relatively disadvantageous.

FDRs & local FDRs

$$\begin{aligned}
 \text{FDR}(0.01) &= \frac{\text{average number of false discoveries at 0.01 significance}}{\text{average number of discoveries at 0.01 significance}} \\
 &= \frac{\text{average number of p-values} < 0.01 \text{ for equivalently expressed genes}}{\text{average number of p-values} < 0.01}
 \end{aligned}$$

$$\begin{aligned}
 \widehat{\text{FDR}}(\alpha) &= \frac{\text{estimated average number of false discoveries}}{\text{estimated average number of discoveries}} \\
 &= \frac{\text{estimated average number of false discoveries}}{\text{number of discoveries}} \\
 &= \begin{cases} \frac{\alpha d}{\#(p(x_i) \leq \alpha)} & \text{if } \frac{\alpha d}{\#(p(x_i) \leq \alpha)} < 1 \\ 1 & \text{if } \frac{\alpha d}{\#(p(x_i) \leq \alpha)} > 1. \end{cases}
 \end{aligned}$$

$$\text{FDR}(\alpha) \approx \frac{\text{LFDR}(p_1) + \text{LFDR}(p_2) + \dots + \text{LFDR}(p_{\#(p(x_i) \leq \alpha)})}{\#(p(x_i) \leq \alpha)}$$

Interpreting the false discovery rate

- If a discovery of differential expression is made whenever the p-value is less than 0.05, then the false discovery rate is the average of all LFDRs corresponding to discoveries

$$\text{FDR}(0.05) = \text{mean}(\text{LFDR}(p(x_i)) | p(x_i) < 0.05)$$

- false discovery rate = probability that randomly selected discovery is false

$$\text{FDR}(0.05) = P(A_i = 0 | p(x_i) < 0.05)$$

Local false discovery rates

- local false discovery rate (LFDR) = posterior probability of equivalent expression: $\text{LFDR}(0.00832) = P(A_i = 0 | p(X_i) = 0.00832)$
- evidence of differential expression = likelihood ratio:

$$\frac{L_i(1)}{L_i(0)} \approx \frac{P(p(X_i) \approx 0.00832 | A_i = 1)}{P(p(X_i) \approx 0.00832 | A_i = 0)}$$

- posterior odds that gene i of p-value 0.00832 is differentially expressed:

$$\frac{1 - \text{LFDR}(0.00832)}{\text{LFDR}(0.00832)} = \frac{P(A_i = 1 | p(X_i) = 0.00832)}{P(A_i = 0 | p(X_i) = 0.00832)} = \frac{P(A_i = 1)}{P(A_i = 0)} \times \frac{L_i(1)}{L_i(0)}$$

Achieved FDR

$$\begin{aligned} \text{FDR}(0.01) &= \frac{\text{average number of false discoveries at 0.01 significance}}{\text{average number of discoveries at 0.01 significance}} \\ &= \frac{\text{average number of p-values} < 0.01 \text{ for equivalently expressed genes}}{\text{average number of p-values} < 0.01} \end{aligned}$$

$$\widehat{\text{FDR}}(\alpha) = \begin{cases} \frac{\alpha d}{\#(p(x_i) \leq \alpha)} & \text{if } \frac{\alpha d}{\#(p(x_i) \leq \alpha)} < 1 \\ 1 & \text{if } \frac{\alpha d}{\#(p(x_i) \leq \alpha)} > 1 \end{cases}$$

$$\widehat{\text{FDR}}(p(x_j)) = \begin{cases} \frac{p(x_j)d}{\#(p(x_i) \leq p(x_j))} & \text{if } \frac{p(x_j)d}{\#(p(x_i) \leq p(x_j))} < 1 \\ 1 & \text{if } \frac{p(x_j)d}{\#(p(x_i) \leq p(x_j))} > 1 \end{cases}$$

Bias in false discovery rates

$$\begin{aligned} \text{FDR}(p(x_j)) &\approx \frac{\text{LFDR}(p_1) + \text{LFDR}(p_2) + \cdots + \text{LFDR}(p_{\#(p(x_i) \leq p(x_j))})}{\#(p(x_i) \leq p(x_j))} \\ &= \frac{\text{LFDR}(p_1) + \text{LFDR}(p_2) + \cdots + \text{LFDR}(p_j)}{\#(p(x_i) \leq p(x_j))} \end{aligned}$$

$$p_1 < p_2 < \cdots < p_j$$

$$\text{LFDR}(p_1) < \text{LFDR}(p_2) < \cdots < \text{LFDR}(p_j)$$

$$\text{LFDR}(p_j) > \frac{\text{LFDR}(p_1) + \text{LFDR}(p_2) + \cdots + \text{LFDR}(p_j)}{\#(p(x_i) \leq p(x_j))} = \text{FDR}(p(x_j))$$

$$\text{FDR}(p(x_j)) < \text{LFDR}(p_j)$$

Correcting the bias

$$p_1 < p_2 < \dots < p_j$$

$$\text{LFDR}(p_1) < \text{LFDR}(p_2) < \dots < \text{LFDR}(p_j)$$

$$\widehat{\text{FDR}}(p(x_j)) = \begin{cases} \frac{p(x_j)d}{\#(p(x_i) \leq p(x_j))} & \text{if } \frac{p(x_j)d}{\#(p(x_i) \leq p(x_j))} < 1 \\ 1 & \text{if } \frac{p(x_j)d}{\#(p(x_i) \leq p(x_j))} > 1 \end{cases}$$

$$\widehat{\text{LFDR}}(p_j) = \left(\frac{1}{j-1+1} + \frac{1}{j-2+1} + \dots + \frac{1}{j-j+1} \right) \widehat{\text{FDR}}(p(x_j))$$

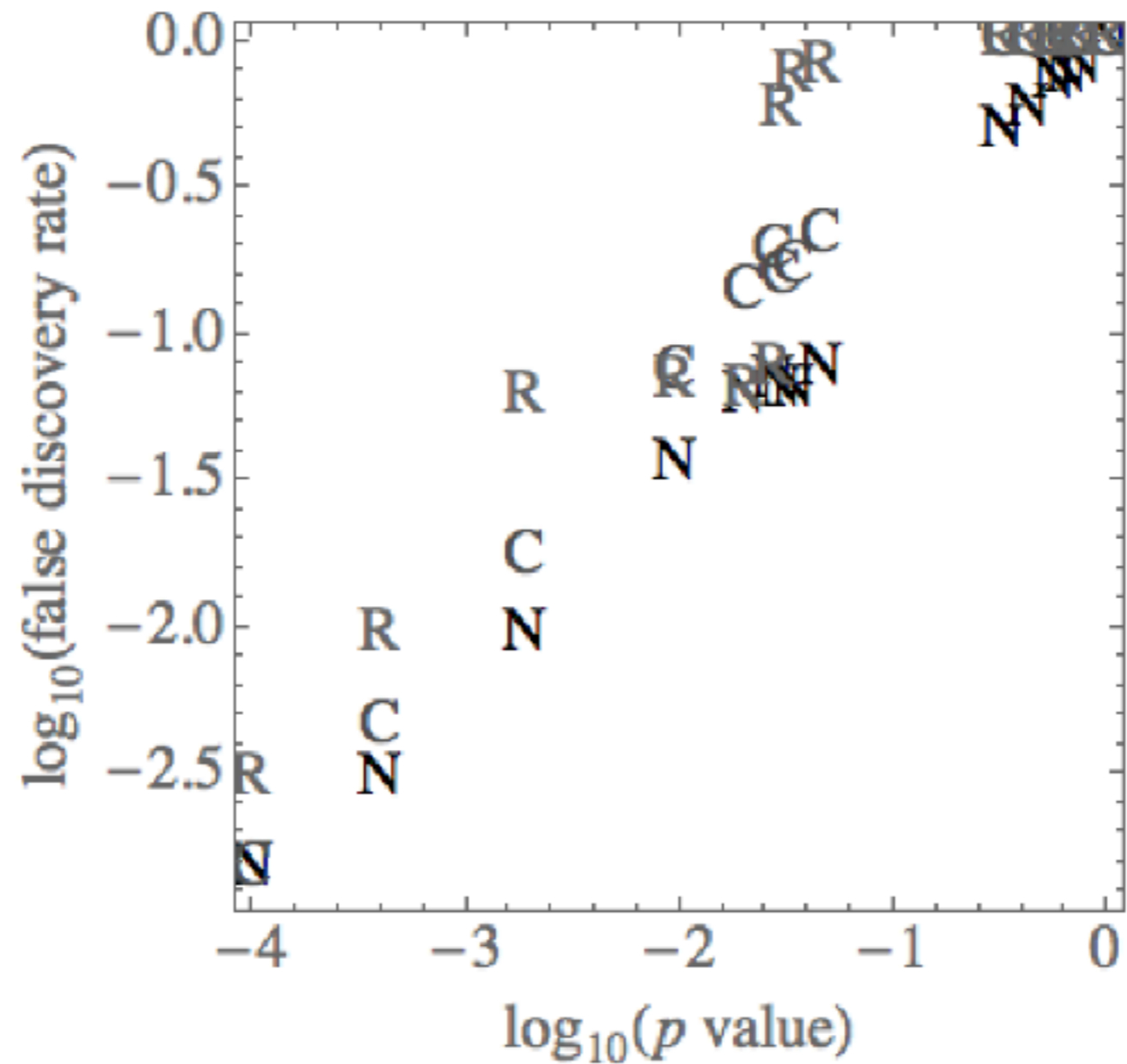
Re-ranked FDRs

Let $q(\alpha)$ denote the smallest value of q such that all hypothesis with p values in $[0, \alpha]$ are rejected according to some procedure that guarantees that the FDR, NFDR, or an estimate of either is no higher than q .

$$\text{RFDR}(x_{(i)}) = q(p(x_{([i/F^*(\text{NFDR})])})) \text{ if } [i/F^*(\text{NFDR})] \leq d \text{ else } \text{RFDR}(x_{(i)}) = 1$$

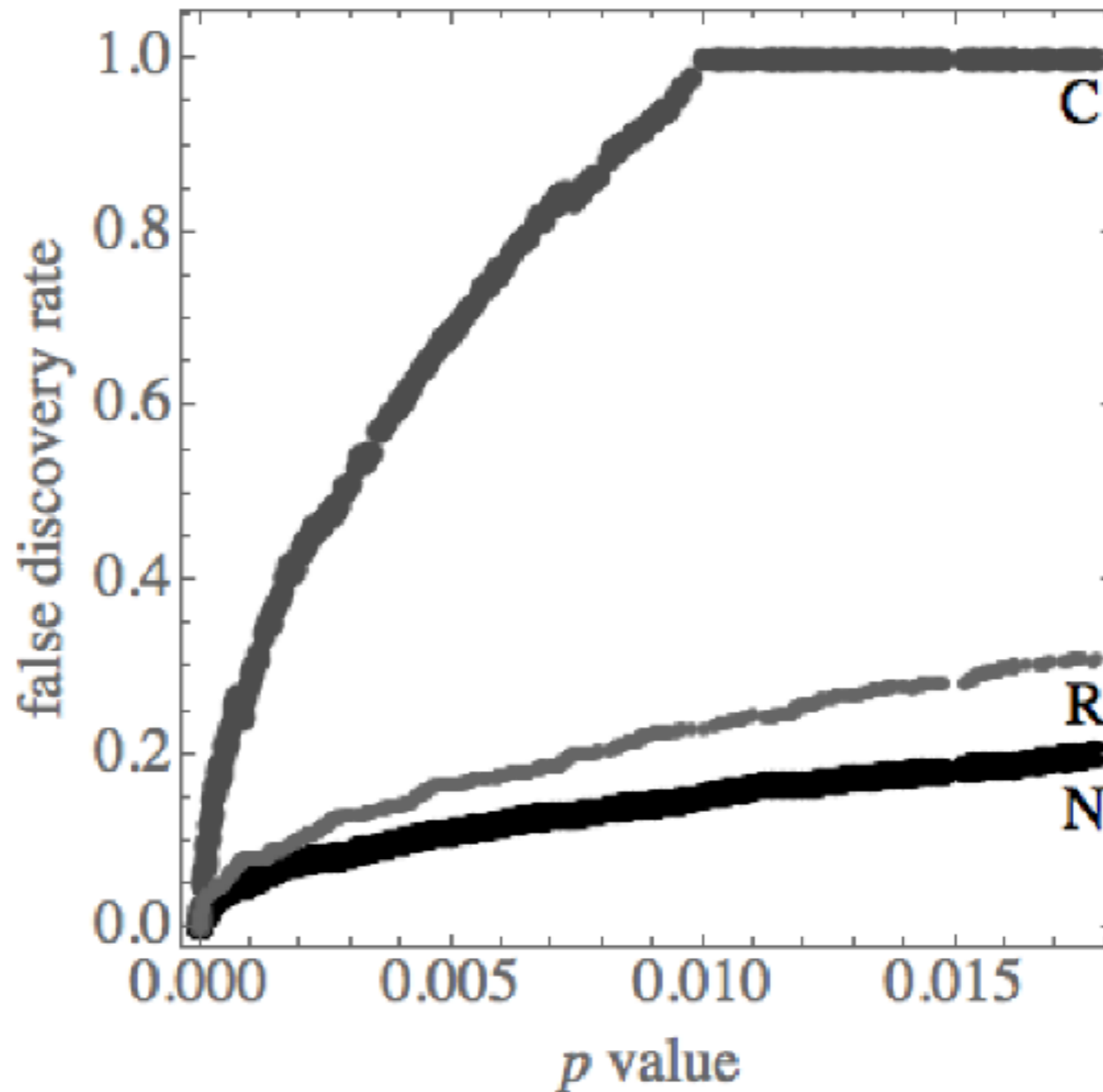
$$F^*(\text{NFDR}) = 1 - e^{-1}$$

Biomedical data



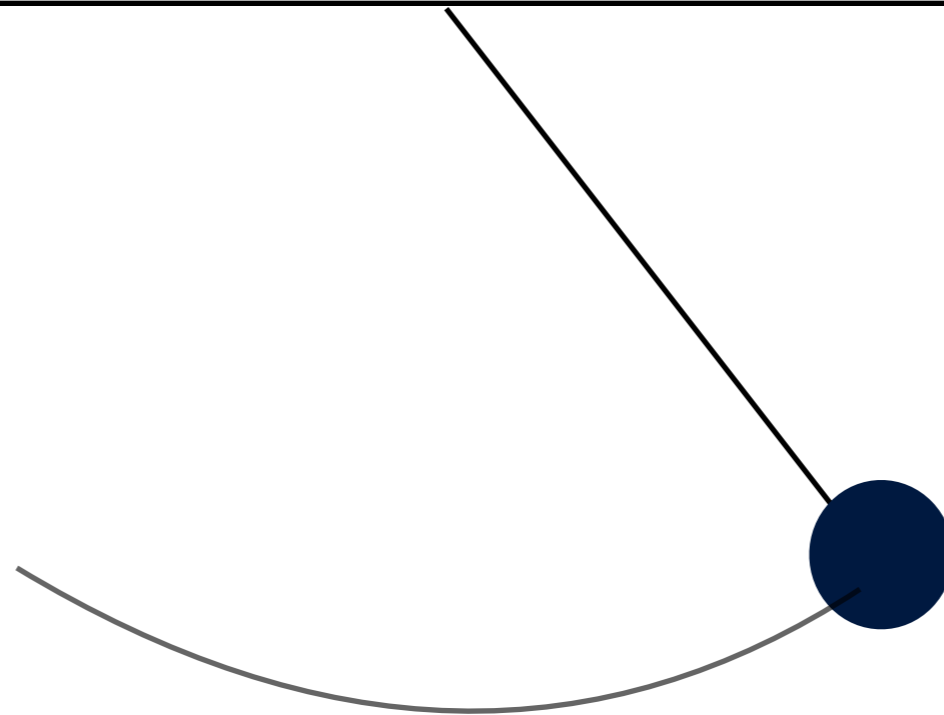
D. R. Bickel, deposited in uO Research at <https://goo.gl/GcUjJe>

Gene expression data



D. R. Bickel, deposited in uO Research at <https://goo.gl/GcUjJe>

The right balance



Family-wise error rates

Too many false negatives

CFDR

RFDR

Other local FDR estimators?

False discovery rates

Too many false positives

Slides and preprint: www.davidbickel.com

Acknowledgements

- Collaborators:

- Alexandre Blais

- Abbas Rahal



uOttawa

- Funding:

- Agriculture and Agri-Food Canada

- Faculty of Medicine of the University of Ottawa