

A SEMIPARAMETRIC APPROACH TO EMPIRICAL BAYES ESTIMATION OF DISCRETE FALSE DISCOVERY RATES USING KERNELS

D. B. Daya¹ S. M. Escalante²

¹School of Statistics
University of the Philippines, Diliman

²*Graduate Student*
Aboitiz School of Innovation, Technology, & Entrepreneurship
Asia Institute of Management

IASC-ARS Interim Conference, December 2022

AGENDA

- 1 INTRODUCTION
- 2 THE MULTIPLE HYPOTHESIS PROBLEM
- 3 USING DISCRETE KERNELS
- 4 SIMULATION RESULTS
- 5 APPLICATION TO HEDENFALK [2001]

INTRODUCTION

Contemporary statistical computing sees a fertile ground of problems in genomics, where inference can be done simultaneously on thousands of genes.

Efron et al. [2001] took the idea of False Discovery Rates and approached it from a Bayesian perspective. Also introduced the **local false discovery rate** (hereafter referred to as the lowercase 'fdr'), quantifying the probability that a given null hypothesis is true given the observed value of its test statistic sense.

Robin et al. [2007] and Guedj et al. [2009], this mixture model comes center field in a novel estimation procedure for the **local false discovery rate involving kernel functions**.

THE MULTIPLE HYPOTHESIS PROBLEM

Consider testing m null hypotheses simultaneously. In the following table, m_0 denotes the number of true null hypotheses, of which V of them are falsely rejected.

TABLE: The multiple hypothesis testing problem

		Decision		
		Do not reject H_0	Reject H_0	Total
Actual Setting	True H_0	U	V	m_0
	False H_0	T	S	$m - m_0$
	Total	$m - R$	R	m

THE MULTIPLE HYPOTHESIS PROBLEM

Benjamini and Hochberg [1995] defined FDR as

$$FDR = E \left[\frac{V}{\max(R, 1)} \right] = E \left[\frac{V}{R} \right] \text{ for } R \geq 1 \quad (1)$$

THE MULTIPLE HYPOTHESIS PROBLEM

The FDR has also been investigated on a Bayesian perspective, as posterior probability distribution holding certain assumptions on the said framework.

Consider the problem of testing m hypotheses. Let H_i be equal to 0 if the null hypothesis is true, 1 if it is false.

Note that $H_i \sim \text{Bernoulli}(1 - \pi_0)$, where π_0 is our prior knowledge of the proportion of the null hypotheses being tested that are actually true (i.e. $\pi_0 = m_0/m$).

Then π_0 also happens to be the probability a null hypothesis being true, which we write as $P(H_i = 0)$.

THE MULTIPLE HYPOTHESIS PROBLEM

Define $\pi_1 = 1 - \pi_0$, the probability of null hypothesis being false, $P(H_i = 1) = 1 - P(H_i = 0)$.

Furthermore, for the m corresponding test statistics x_1, x_2, \dots, x_m , let f_0 be the density of x_i , if the null hypothesis is true; if it is false, let f_1 be its density.

Then the distribution of X is given by

$$f(x) = \pi_0 f_0(x) + \pi_1 f_1(x) \quad (2)$$

THE MULTIPLE HYPOTHESIS PROBLEM

With the assumptions discussed above, Storey [2003] presented the form of Bayesian FDR (FDR) as follows:

$$FDR(\mathcal{R}) = P(H = 0 | x \in \mathcal{R}) \quad (3)$$

where \mathcal{R} is the set of "rejected" null hypotheses.

$FDR(\mathcal{R})$ represents the resulting false discovery rate of choosing to declare the tests in \mathcal{R} as significant.

THE MULTIPLE HYPOTHESIS PROBLEM

Now let F_0 be the cumulative distribution function (CDF) corresponding to f_0 , and F be the CDF corresponding to f .

Storey [2003] showed that by applying Bayes' rule, the FDR takes the form

$$FDR(\mathcal{R}) = \frac{\pi_0 F_0(\mathcal{R})}{F(\mathcal{R})} \quad (4)$$

THE MULTIPLE HYPOTHESIS PROBLEM

Efron [2005] further explored this Bayesian interpretation of the FDR. While FDR comprises the whole set of hypotheses, Efron introduced the local fdr , which is the probability that a hypothesis is true given the data, which is given by

$$fdr(x) = \frac{\pi_0 f_0(x)}{f(x)} \quad (5)$$

THE MULTIPLE HYPOTHESIS PROBLEM

Methods for estimating these figures have already been explored for the continuous case. Efron and Tibshirani [2002] let f_0 be known and assumed to follow the standard normal $N(0, 1)$ distribution. The CDF $F(\mathcal{R})$ is then estimated by

$$\hat{F}(R) = \frac{\#\{z_i \in \mathcal{R}\}}{m} \quad (6)$$

THE MULTIPLE HYPOTHESIS PROBLEM

As for the local fdr, a close analog is the following estimator:

$$\widehat{fdr}(x) = \frac{\pi_0 f_0(x)}{\widehat{f}(x)} \quad (7)$$

USING DISCRETE KERNELS

Kernel Density Estimation is a nonparametric approach, which attempts to improve upon standard histogram methods by reweighting the data with the use of “kernel functions.” A rough analog of the histogram method may be presented as follows:

$$\hat{f}(x) = \frac{1}{n} \sum_j I(x_j = x) \quad (8)$$

USING DISCRETE KERNELS

The kernel function will have to be replaced with one that is more suited for a discrete support. We refer to two such kernels, the first being the Dirac kernel mentioned in Li and Racine[2007]:

$$k_j(x; h) = \begin{cases} 1 - h & x = x_j \\ \frac{h}{c-1} & x \neq x_j \end{cases} \quad (9)$$

as well as the Aitken kernel defined in Li et al [2008]:

$$k_j(x; h) = \begin{cases} 1 & x = x_j \\ h & x \neq x_j \end{cases} \quad (10)$$

USING DISCRETE KERNELS

Algorithm 1 Estimating the Empirical Bayes FDR

Initialize: $\hat{\tau}_j = 1 - \pi_0$

while $err \leq \epsilon$ **do**

$$\hat{f}_1(x) \leftarrow \sum_j \tau_j K(x_j; h) / \sum_i \tau_i$$

$$\hat{f}(x) \leftarrow \pi_0 f_0(x) + (1 - \pi_0) \hat{f}_1(x)$$

$$\hat{\tau}_j \leftarrow 1 - \pi_0 f_0(x) / \hat{f}(x)$$

$$err \leftarrow \max_j \frac{|\tau_j^{(l)} - \tau_j^{(l-1)}|}{\tau_j^{(l-1)}}$$

end while

USING DISCRETE KERNELS

Algorithm 2 Cross-validation on the Kernel Smoothing Parameter

Initialize: randomly divide set of values x into V non-overlapping folds: Y_1, Y_2, \dots, Y_V . Define X_k as the collection of all $Y_{i \neq k}$, that is the data-set taking out the k th fold.

for fold $i \in \{1, 2, \dots, V\}$ **do**

Fit \hat{f}_1 on X_i (the training set) using Algorithm 1 with given h

Using the fitted \hat{f}_1 , estimate $\hat{f}_1(y)$ for every y from Y_i (the testing set)

Define the log-likelihood of the subset y_i as $L(y_i; h) = \sum_j \ln \hat{f}_1(y_j)$

Then the V -fold cross-validation log-likelihood is defined as $L_{CV}(h) = \frac{1}{V} \sum_j L(y_j; h)$

end for

We choose $h = \arg \max_h L_{CV}(h)$

USING DISCRETE KERNELS

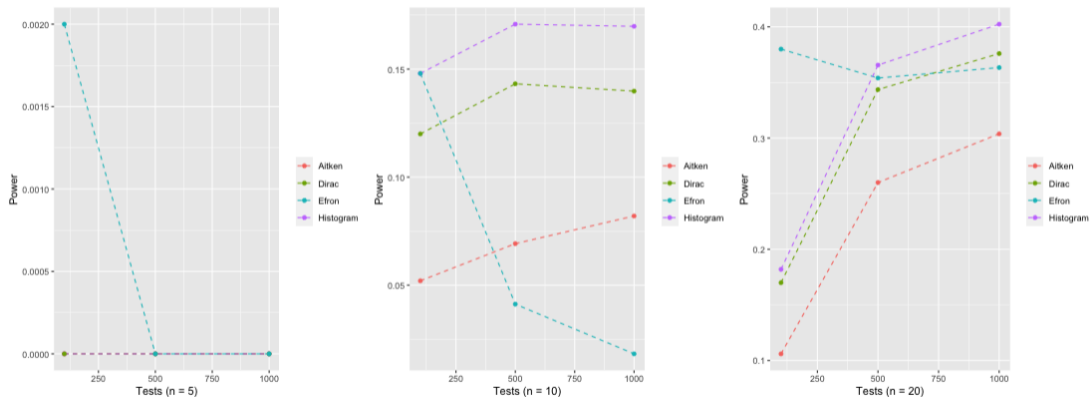


FIGURE: Simulated power, assuming true $\pi_0 = 0.95$, in the mixed case

USING DISCRETE KERNELS

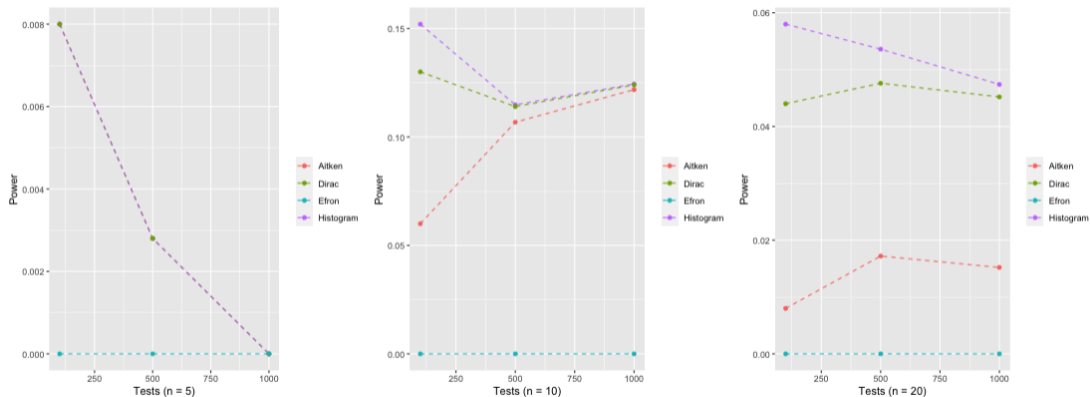


FIGURE: Simulated power, assuming true $\pi_0 = 0.95$, in the overlapping case

USING DISCRETE KERNELS

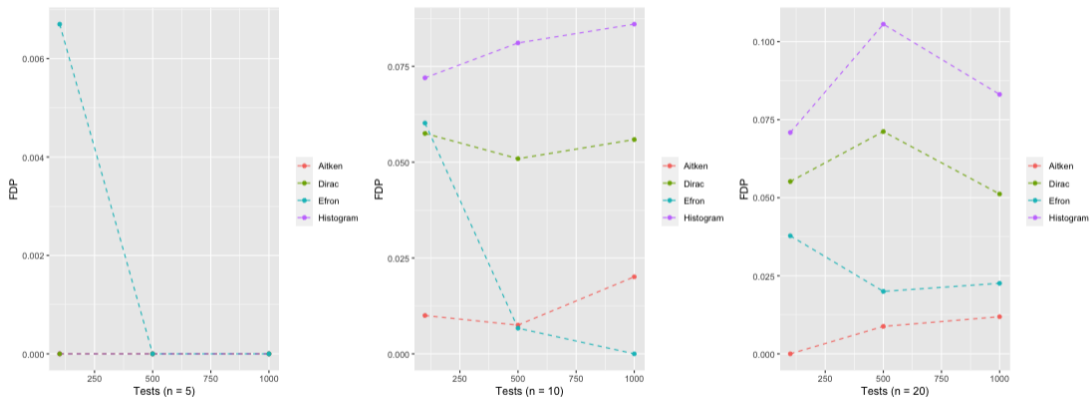


FIGURE: Simulated FDP, assuming true $\pi_0 = 0.95$, in the mixed case

USING DISCRETE KERNELS

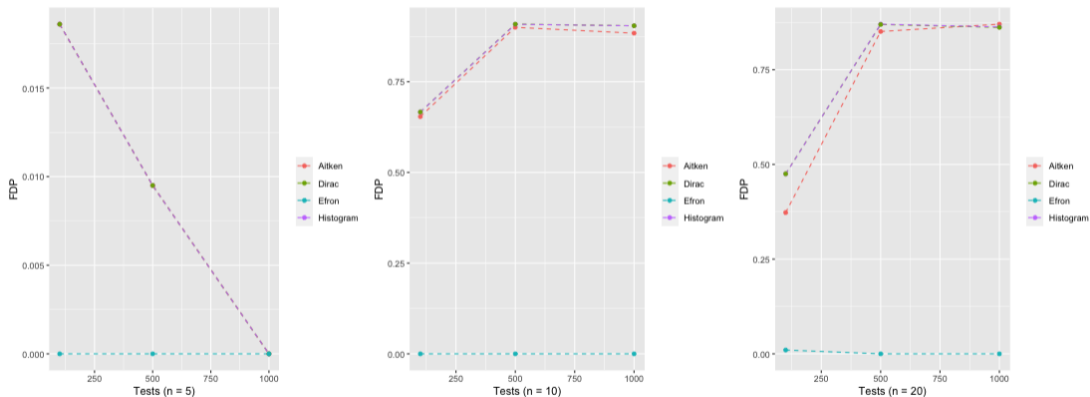


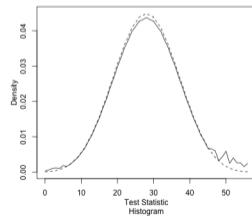
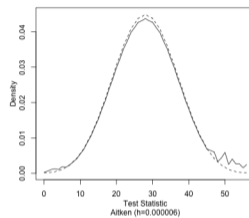
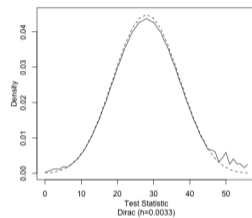
FIGURE: Simulated FDP, assuming true $\pi_0 = 0.95$, in the overlapping case

APPLICATION TO HEDENFALK [2001]

We apply the proposed procedure on data reported by Hedenfalk et al. [2001], concerning patients diagnosed with two different kinds of breast cancer.

The data used in the present paper concerns seven BRCA1 and eight BRCA2 patients (corresponding to gene mutations characterizing the disease), whose expression ratios were measured across $m = 3226$ genes.

APPLICATION TO HEDENFALK [2001]



APPLICATION TO HEDENFALK [2001]

A total of 51 genes were found to be significant at $FDR_i < 0.10$ if estimation using the Dirac kernel is applied.

We observe that Efron's procedure tends to overestimate not only the proportion of true nulls but also the corresponding local false discovery rates.

Cutoff, q	Genes with $fdr < q$	
	Dirac	Efron
0.050	0.00	0.00
0.075	26.00	0.00
0.100	51.00	0.00

APPLICATION TO HEDENFALK [2001]

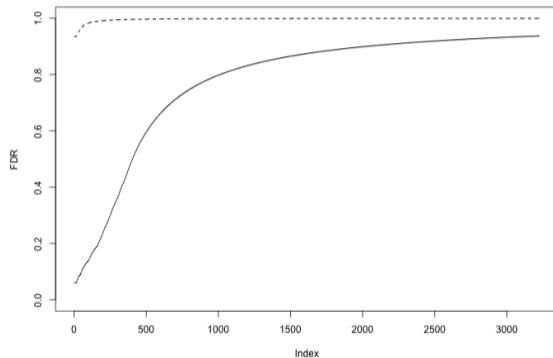


FIGURE: Estimated Bayesian FDR (ordered) for Dirac (solid) versus Efron (dashed)

SELECTED REFERENCES

- 1 Efron, B (2004). "Large-scale simultaneous hypothesis testing: the choice of a null hypothesis." Journal of the American Statistical Association. 2004, 99:96-104.
- 2 Chen, Xiongzhi and Doerge, R.W. (2015). "A weighted FDR procedure under discrete and heterogeneous null distributions." [URL](https://arxiv.org/abs/1502.00973)
- 3 Guedj, Mickael; Robin, Stephane; Celisse, Alain; and Nuel, Gregory (2009). "Kerfdr: a semi-parametric kernel-based approach to local false discovery rate estimation." BMC Bioinformatics 2009, 10:84. DOI:10.1186/1471-2105-10-84
- 4 Li, Q.; and Racine, J. (2007). "Nonparametric econometrics: Theory and practice." Princeton University Press.
- 5 Benjamini, Y; and Hochberg, Y (1995): "Controlling the false discovery rate: a practical and powerfull approach to multiple testing." JRSSB 1995, 57:289-300.
- 6 Hedenfalk, I; Duggan, D; Chen, YD; Radmacher, M; Bittner, M; Simon, R; Meltzer, P; Gusterson, B; Esteller, M; Kallioniemi, OP; Wilfond, B; Borg, A; and Trent, J (2001). "Gene expression profiles in hereditary breast cancer." New England Journal of Medicine, 344:539-548.