***Warning:*** *These notes may contain factual and/or typographic errors. They are based on Emmanuel Candès's course from 2018 and 2021, and scribe notes.*

# Outline

**Agenda**: False Discovery Rate.

1. Empirical Process viewpoint of BHq.

2. Empirical Process viewpoint of FDR control.

3. Improving on BHq.

Much of the material in this lecture is taken from Storey, Siegmund, and Taylor (2004) [1].

## 7.1 The Empirical Process Viewpoint of BHq

In previous lectures, we introduced the BH procedure by looking at the sorted $p$-values on the $x$-axis and whether they fall below a critical line. An alternative way to view BH is to flip the axes and view the sorted $p$-values on the $y$-axis. This is illustrated in the following figure.



(a) P-values on the $y$ axis, indices on $x$ (b) P-values on the $x$ axis, indices on $y$
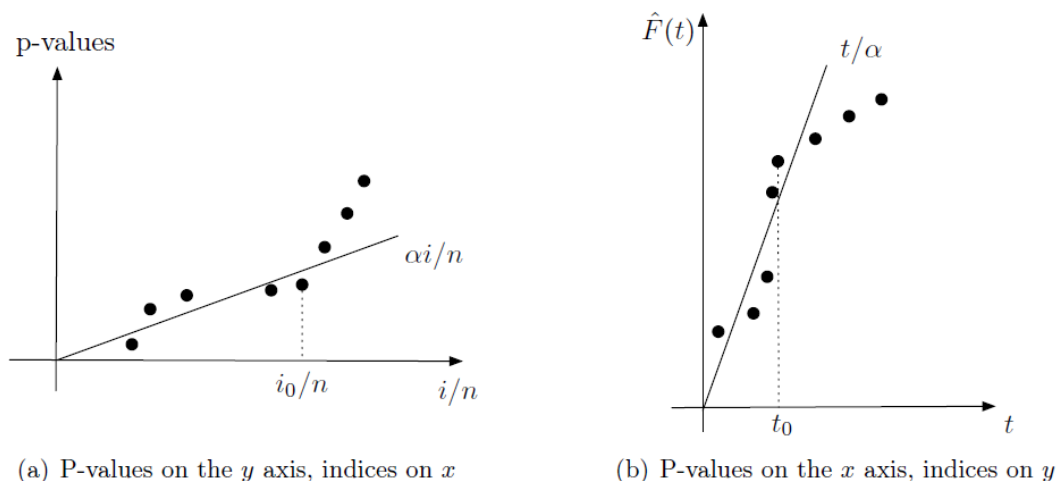
Figure 1: Sorted p-values and BH($q$) threshold line

This alternative view allows us to describe the BH procedure in terms of an empirical process. The coordinates on the $y$-axis are the values of the empirical CDF

$$\widehat{F}_n(t) = \frac{\#\{i : p_i \le t\}}{n}$$

evaluated at the $p$-values. Assume that the $p$-values and hypotheses are ordered in increasing order:

$$p(1) \le ... \le p(n), \qquad H(1) \le H(n)$$

BHq is defined to reject $H(1), ..., H(i_0)$ where

$$i_0 = \max\left\{i : p(i) \le \frac{qi}{n}\right\}$$

The critical $p$-value is $p^* = p(i_0)$ and can be written as

$$\begin{aligned}
p^* &= \max\left\{p(i) : p(i) \le \frac{qi}{n}\right\} \\
&= \max\left\{p(i) : p(i) \le q\widehat{F}_n(p(i))\right\} \\
&= \max\left\{t \in \{p_1, ..., p_n\} : t \le q\widehat{F}_n(t)\right\}
\end{aligned}$$

If the set is empty, the convention is $p^* = q/n$. Therefore, the BHq procedure is equivalent to rejecting all hypotheses with $p_i \le \tau_{BH}$:

$$\tau_{BH} = \max\left\{t : \frac{t}{\widehat{F}_n(t) \vee 1/n} \le q\right\}$$

Notice that $\tau_{BH} \ge q/n$.

This formulation has a simple interpretation. Let $t \in (0, 1)$ be fixed and consider rejecting $H_i$ iff $p_i \le t$. We can construct the rejection/acceptance table for the hypotheses whose values depend on $t$.

|            | $H_0$ accepted | $H_0$ rejected | Total |
|------------|----------------|----------------|-------|
| $H_0$ true | $U(t)$ | $V(t)$ | $n_0$ |
| $H_0$ false | $T(t)$ | $S(t)$ | $n - n_0 = n_1$ |
|            | $n - R(t)$ | $R(t)$ | $n$ |

We define

$$FDP(t) = \frac{V(t)}{R(t) \vee 1}, \qquad FDR(t) = \mathbb{E}[FDP(t)]$$

The idea is to choose the threshold $t$ as large as possible while controlling the $FDR$ at level $q$. If we had an estimate $\widehat{FDR}$ of the $FDR$, we can take the threshold $\tau$ to be

$$\tau = \sup\{t \le 1 : \widehat{FDR}(t) \le q\}$$

and define the rejection rule to reject $H_i$ iff $p_i \leq \tau$. We hope that this method indeed controls the false discovery rate and that having such a liberal test will not decrease the power too much. The first question is how to estimate $FDR(t)$.

We know that $\mathbb{E}[V(t)] = n_0 t$, but $n_0$ is not known. Therefore, a conservative estimate of $n_0 t$ is $nt$, which leads to our first estimate

$$\widehat{FDR}(t) = \frac{nt}{R(t) \vee 1} = \frac{t}{\widehat{F}_n(t) \vee 1/n}$$

This leads us to exactly the BH procedure since

$$\tau_{BH} = \sup \left\{ t \leq 1 : \frac{nt}{R(t) \vee 1} \leq q \right\}$$

**Theorem 1.** Under independence, this FDR estimate is biased upwards:

$$\mathbb{E}[\widehat{FDR}(t)] \geq FDR(t)$$

For a proof, see [1].

## 7.2   Martingale Theory and FDR Control

We can invert the estimate of FDR to prove FDR control using martingales, giving us an alternate proof of the Benjamini-Hochberg result.

**Theorem 2.** BH (1995).
The procedure rejecting all hypotheses with $p_i \leq \tau_{BH}$ controls the FDR:

$$\mathbb{E}[FDR(\tau_{BH})] = qn_0/n$$

*Proof.* We let $\tau = \tau_{BH}$. Define the filtration

$$\mathcal{F}_t = \sigma(V(s), R(s) : t \leq s \leq 1)$$

Notice this is a backwards filtration: for $t_1 < t_2$, $\mathcal{F}_{t_2} \subset \mathcal{F}_{t_1}$. Define the reverse martingale $\{V(t)/t, 0 \leq t \leq 1\}$. We prove this is indeed a martingale: Let $s \leq t$.

$$\begin{aligned}
\mathbb{E}\left[ \frac{V(s)}{s} \Big| \mathcal{F}_t \right] &= \frac{1}{s} \mathbb{E}[V(s)|\mathcal{F}_t] \\
&= \frac{1}{s} \cdot \frac{s}{t} V(t) \qquad (*) \\
&= \frac{V(t)}{t}
\end{aligned}$$

where in $(*)$ we used the fact that under $\mathcal{F}_t$, $V(t) = \#\{p_i : p_i \leq t, H_i \text{ null}\}$ and these $p_i \sim U[0, t]$ and are independent. This proves $\{V(t)/t, 0 \leq t \leq 1\}$ is a martingale.

Next, notice that $\tau_{BH}$ is a stopping time with respect to $\{\mathcal{F}_t\}$. This is because knowing $V(s), R(s) = n\widehat{F}_n(s)$ for $s \geq t$ will determine whether $\tau \leq t$. Therefore, $\{\tau \leq t\} \in \mathcal{F}_t$ and $\tau$ is a stopping time.

We are ready to apply Doob's Optional Stopping Theorem. By definition, $R(\tau) \vee 1 = n\tau/q$. Therefore,

$$
\begin{aligned}
FDR(\tau) &= \mathbb{E}\left[\frac{V(\tau)}{R(\tau) \vee 1}\right] \\
&= \frac{q}{n}\mathbb{E}\left[\frac{V(\tau)}{\tau}\right] \\
&= \frac{q}{n}\mathbb{E}\left[\frac{V(1)}{1}\right] \\
&= \frac{q}{n} \cdot n_0
\end{aligned}
$$

$\square$

## 7.3   Improving on BHq

We want to improve the simple conservative estimate of $\hat{\pi}_0 = 1$ by using the distribution of $p$-values. Fix $\lambda \in [0, 1)$ and define

$$
\hat{\pi}_0^\lambda = \frac{n - R(\lambda)}{(1 - \lambda)n}
$$

We usually will take $\lambda = 1/2$, while $\lambda = 0$ recovers the BHq procedure. The motivation for this estimation is the following:

$$
\hat{\pi}_0^\lambda = \frac{n_0 - V(\lambda) + n_1 - S(\lambda)}{(1 - \lambda)n}
$$

We would expect the non-null $p$-values to be small, so $n_1 - S(\lambda) \approx 0$, and hence

$$
\hat{\pi}_0^\lambda \approx \frac{n_0 - S(\lambda)}{(1 - \lambda)n} \approx \frac{n_0 - (n_0/2)}{n/2} = \frac{n_0}{n}
$$

Our estimate for the false discovery rate is

$$
\widehat{FDR}^\lambda(t) = \hat{\pi}_0^\lambda \cdot \frac{nt}{R(t) \vee 1}
$$

and the natural test would be to reject $H_i$ iff $p_i \leq \tau$,

$$
\tau = \sup\{t \leq 1 : \widehat{FDR}(t) \leq q\}
$$

In cases where $\hat{\pi}_0^\lambda$ is smaller than 1, say 0.8, we may get more powerful results than BHq because we have a significant proportion of non-nulls.

There are several drawbacks to this approach. One drawback is that we may have $\hat{\pi}_0^\lambda > 1$, in which we are being even more conservative in our estimation. In addition, the threshold $\tau$ may not control the FDR, so we introduce a modified version called Storey's procedure.

## 7.4   Storey's Procedure

Storey's procedure involves a simple modification of to the estimate of $\pi_0$ defined in the previous section. Define

$$\hat{\pi}_0 = \frac{1 + n - R(1/2)}{n/2}$$

The only difference between $\hat{\pi}_0$ and $\hat{\pi}_0^{1/2}$ is the added 1 in the numerator. Our test now becomes reject $H_i$ iff $p_i \leq \tau$,

$$\tau = \sup\left\{t \leq \frac{1}{2} : \widehat{FDR}(t) = \frac{1 + n - R(1/2)}{n/2} \cdot \frac{nt}{R(t) \vee 1} \leq q\right\}$$

Notice that we only take the supremum over $t \leq \frac{1}{2}$, which is necessary because the estimate of $\pi_0$ used the information of the $p$-values $> 1/2$.

**Theorem 3.** Storey's procedure controls FDR at level $q$.

*Proof.* We use martingale theory in a proof similar to the proof of Theorem 2. We know that $\widehat{FDR}(\tau) = q$. Then

$$FDR(\tau) = \mathbb{E}\left[\frac{V(\tau)}{R(\tau) \vee 1}\right]$$

$$= \mathbb{E}\left[\frac{V(\tau)}{n\tau} \cdot \frac{n\tau}{R(\tau) \vee 1} \cdot \frac{1 + n - R(1/2)}{n/2} \cdot \frac{n/2}{1 + n - R(1/2)}\right]$$

$$= \mathbb{E}\left[\widehat{FDR}(\tau) \cdot \frac{V(\tau)}{n\tau} \cdot \frac{n/2}{1 + n - R(1/2)}\right]$$

$$= q \cdot \mathbb{E}\left[\frac{V(\tau)}{\tau} \cdot \frac{1/2}{1 + n - R(1/2)}\right]$$

Applying Doob's Optional Stopping Theorem to the martingale $\{V(t)/t : t \in [0, 1/2]\}$ and stopping time $\tau$, we have

$$FDR(\tau) = q \cdot \mathbb{E}\left[\frac{V(1/2)}{1/2} \cdot \frac{1/2}{1 + n - R(1/2)}\right]$$

$$= q \cdot \mathbb{E}\left[\frac{V(1/2)}{1 + n - S(1/2) - V(1/2)}\right]$$

$$\leq q \cdot \mathbb{E}\left[\frac{V(1/2)}{1 + n_0 - V(1/2)}\right]$$

where the last inequality holds because $n_1 - S(1/2) \geq 0$.

We directly calculate $\mathbb{E}[\frac{V(1/2)}{1+n_0-V(1/2)}] \leq 1$. We know that $V(1/2) \sim Bin(n_0, 1/2)$. Then

$$
\begin{aligned}
\mathbb{E}\left[\frac{V(1/2)}{1+n_0-V(1/2)}\right] &= \sum_{i=1}^{n_0} \mathbb{P}(V(1/2) = i)) \cdot \frac{i}{1+n_0-i} \\
&= 2^{-n_0} \sum_{i=1}^{n_0} \binom{n_0}{i} \cdot \frac{i}{1+n_0-i} \\
&= 2^{-n_0} \sum_{i=1}^{n_0} \frac{i \cdot n_0!}{(n_0-i+1) \cdot (n-i)! \cdot i!} \\
&= 2^{-n_0} \sum_{i=1}^{n_0} \frac{n_0!}{(n_0-i+1)!(i-1)!} \\
&= 2^{-n_0} \sum_{j=0}^{n_0-1} \binom{n_0}{j} \\
&= 2^{-n_0}(2^{n_0} - 1) \\
&= 1 - 2^{-n_0} \\
&\leq 1
\end{aligned}
$$

Therefore, $FDR(\tau) \leq q$ and this concludes the proof. $\qquad\square$

# References

[1]  Taylor J. Storey J. and Siegmund D. "Strong Control, Conservative Point Estimation and Simultaneous Conservative Consistency of False Discovery Rates: A Unified Approach". In: *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 66.1 (2004), pp. 187–205.