# Chapter 4: False Discovery Rate Control

Benjamini & Hochberg FDR Control

Under the null

$$H_{0i}: P_i \sim \text{Unif}(0,1)$$

which we can order

$$P_{(1)} \leq P_{(2)} \leq \cdots \leq P_{(N)}$$

Let $R_D$ be the number of cases rejected then the false discovery proportion with respect to $D$ is

$$Fdp_D = \frac{a_p}{R_D}$$

where

Decision

| | Null | Non Null | |
| --- | --- | --- | --- |
| Null | N — a | | |

Actual

| | Null | | $N_0$ |
|---|---|---|---|
| Nonnull | $N_1 - b$ | $b$ | $N_1$ |
| | $N - R$ | $R$ | $N$ |

So we have $a$ Type I errors

and $N - R$ Type II errors.

Here we have $R = a + b$ hypothesis

Task: How many of $R$ are true

discoveries.

Def: Family wise error rates (FWER)

$$\mathbb{P}(a > 0)$$

Check: When $N = 1$

$$\mathbb{P}(a = 1 \mid N_0 = 1) = \alpha$$

$$\mathbb{P}(b = 1 \mid N_1 = 1) = \beta$$

Under independence with rejection region $R$

$a | R \sim Binom(R, \phi(z))$

counts number of type II

$\phi(z) = \mathbb{P}(\text{reject} | \text{null is true})$

Suppose $f(z) = \pi_0 f_0(z) + \pi_1 f_1(z)$

is the marginal dist. of the test stat.

$$\phi(z) = \frac{\mathbb{P}_0(z \in Z) \pi_0}{\mathbb{P}(z \in Z)}$$

$$= \frac{\pi_0 \int_Z f_0(z) dz}{\int_Z f(z) dz}$$

$$= FDR(z)$$

What is the dist. of a p-value?

$$\mathbb{P}(P_i \leq u) = \mathbb{P}(F(z_i) \leq u)$$

$$= \mathbb{P}(z_i \leq F^{-1}(u))$$

$$= F_z(F_z(u)) = u$$

So the false discovery prop.

$$FdP_D = \frac{a_D}{R_D} \longleftarrow$$ # false discovery for region / decision rule

D.

## Benjamini Hochberg

For a fixed value $z \in (0,1)$ let

$$i_{max} = \text{argmax} \quad r_{(i)} \leq \frac{i}{N} q$$

Decision: reject $H_{0(i)}$ if $i \leq i_{max}$

and the adjusted p-value is given by

$$\frac{N}{i} P_{(i)} < q$$

Thrm: Under independence the BH algorithm controls the expected false discovery prop. at $q$

$$\mathbb{E}\left(Fdp_{BH(q)}\right) = \pi_0 q \leq q \qquad \pi_0 = \frac{N_0}{N}$$

Pf: $t \in [0,1]$ $\quad R(t) = \# p_i \leq t$

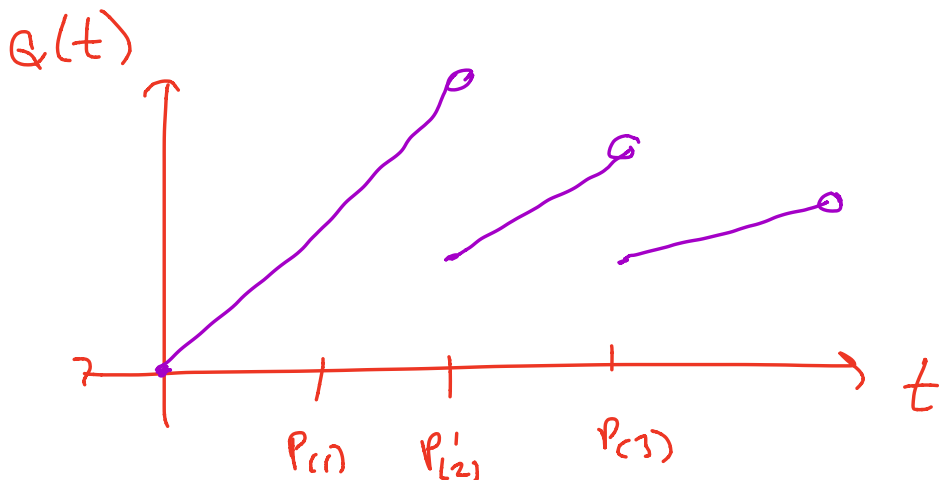$\quad a(t) = \#$ of false discoveries

$$Fdp(t) = \frac{a(t)}{max(R(t)), 1}$$

$$Q(t) = \frac{\# \text{ rejects}}{\max(R(t), 1)}$$

$$\overset{H_0}{=} \frac{Nt}{\max(R(t), 1)}$$

In the BH procedure

$$P_i \leq \frac{i}{N} z \iff \frac{N}{i} P_i \leq z$$

$$\iff Q(P_{(i)}) \leq z$$



Using a generalized inference,

$$t = \sup \{ Q(t) \leq q \}$$

$$z_t = \frac{}{t}(\quad) = t)$$

So the decision becomes

$$P_{(i)} \leq t_{\bar{q}}.$$

Now define $A(t) = \dfrac{a(t)}{t}$

We claim this is a martingale.

$$\frac{1}{s} \mathbb{E}\left(a(s) \mid a(t) = y_2\right)$$

$$\mathbb{P}\left(a(s) = y_1 \mid a(t) = y_2\right) \qquad y_1 \leq y_2$$

$$= \frac{\mathbb{P}\left(a(s) = y_1, \ a(t) = y_2\right)}{\mathbb{P}\left(a(t) = y_2\right)}$$

$$= \frac{\mathbb{P}\left(a(s) = y_1, \ a(t) - a(s) = y_2 - y_1\right)}{\mathbb{P}\left(a(t) = y_2\right)}$$

Now $\quad a(t) \sim \text{Binom}(N_0, t)$

by looking at the rows

$$= \frac{\binom{N_o}{y_1 \; y_2 - y_1} s^{y_1} (t-s)^{y_2-y_1} (1-t)^{N_o - y_2}}{\binom{N_o}{y_2} t^{y_2} (1-t)^{N_o - y_2}}$$

$$\xleftarrow{\;\; \overset{y_1}{|} \;\; \overset{y_2 - y_1}{|} \quad N - y_2 \;\;} \to$$

$$\underset{s}{\phantom{|}} \quad \underset{t}{\phantom{|}}$$

$$= \frac{\binom{N_o}{y_1, y_2 - y_1}}{\binom{N_o}{y_2}} \left(\frac{s}{t}\right)^{y_1} \left(\frac{(t-s)}{t}\right)^{y_2} \cdot$$

$$= \frac{\dfrac{\cancel{N_o!}}{y_1! \; y_2 - y_1! \; \cancel{N_o - y_2!}}}{\dfrac{\cancel{N_o!}}{y_2! \; \cancel{N_o - y_2!}}} \; \left(s/t\right)^{y_1} \left(1 - s/t\right)^{y_2}$$

$$= \binom{y_2}{y_1} \left(\frac{1}{t}\right)^{\sigma_1} \left(1 - s/t\right)^{\sigma_2 - \sigma_1}$$

$$\frac{1}{s} \mathbb{E}\left(a(s) \mid a(t) = y_2\right) = \frac{1}{s} y_2 \frac{s}{t} = \frac{y_2}{t}$$

$$\mathbb{E}\left(A(s) \mid A(t) = t_2\right)$$

$$= \mathbb{E}\left(A(s) \mid a(t) = t t^2\right)$$

$$= \frac{t t_2}{t} = t_2$$

So $A(s)$ is a decreasing martingale with stopping time $t_2$. So by optimal sampling thearm

$$\mathbb{E}\left(A(t_2)\right) = \mathbb{E}(A(1)) = \mathbb{E}(a(1))$$

$$= N_0$$

$$\max\{R(t_\#),1\} = \frac{Nt_\#}{Q(t_\#)} \simeq \frac{Nt_\#}{\#}$$

So the false discovery prop.

$$\mathbb{E}(FDP_D) = \frac{a(t_\#)}{\max(1, R(t_\#))} = \frac{\# \, a(t_\#)}{Nt_\#}$$

$$\mathbb{E}(FDP(t_q)) = \frac{\#}{N} \mathbb{E}(A(t_\#))$$

$$= \frac{\#}{N} N_0$$

$$= \Pi_0 \# \leq \#$$

- On expect. we control the rate. What about variability

- How should $\#$ be chosen?

- Is the theoretical null

$$p_i \sim Unif(0,1)$$

correct?

## Empirical Bayes Interp.

$$p_{(i)} = F_0(z_{(i)})$$

$$z_{(1)} \leq \cdots \leq z_{(N)}$$

$$\overline{Fdr}(z) = \frac{\pi_0 F_0(z)}{\overline{F}(z)}$$

BH is equivalent to

$$\frac{\pi_0 F_0(z_{(i)})}{\overline{F}(z_i)} \leq \pi_0 \, q$$