

Large scale hypothesis testing and false discovery rate (FDR)

$K \geq 1$ hypotheses to be tested simult.

$H_0: H_{01} \wedge H_{02} \wedge H_{03} \wedge \dots \wedge H_{0K}$ vs $H_1: H_{11} \text{ or } H_{12} \text{ or } \dots \text{ or } H_{1K}$

$H_0: H_{0i} \quad i=1, \dots, K$ vs $H_1: H_{1i} \quad i=1, \dots, K$

For $i=1, \dots, K$ let π_i (p-value) of H_{0i} vs H_{1i}

Bonferroni strategy: if $\alpha \in (0, 1)$ reject H_{0i} if $\pi_i < \frac{\alpha}{K}$

Example

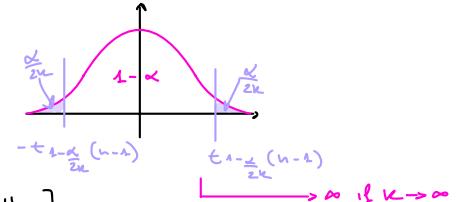
X_1, \dots, X_n iid $\sim N_p(\mu, \Sigma)$

$\alpha_1, \dots, \alpha_K \in \mathbb{R}^p$

$H_{0i}: \underline{\alpha}_i^\top \underline{\mu} = \delta_i \quad \text{vs} \quad H_{1i}: \underline{\alpha}_i^\top \underline{\mu} \neq \delta_i$

Bonf: $\alpha \in (0, 1)$ reject H_{0i} if $\frac{|(\underline{\alpha}_i^\top \underline{\mu}) - \delta_i|}{\sqrt{\underline{\alpha}_i^\top \Sigma \underline{\alpha}_i}} > t_{1-\frac{\alpha}{2K}}(n-1)$

Bonf strategy: $\mathbb{P}\left[\bigcup_{i=1}^K \{\text{rejecting } H_{0i}\} \mid \bigcap_{i=1}^K H_{0i}\right] \leq \sum_{i=1}^K \mathbb{P}[\text{rejecting } H_{0i} \mid H_{0i}] \leq \sum_{i=1}^K \frac{\alpha}{K} = \alpha$



$\rightarrow \alpha \text{ if } K \rightarrow \infty$

FDR: Benjamini & Hochberg

$H_0: H_{0i} \quad i=1, \dots, K$ vs $H_1: H_{1i} \quad i=1, \dots, K$

Let Δ be any strategy for running H_0 vs H_1

		Decision following Δ	
		Do not res H_0	Reject H_0
Truth	H_0	V	V false discovery
	H_1	T Missed discovery	S true discovery
	$K-R$	R Observables	K

Assume H_{0i} is true iff $i \in I_0$ $|I_0| = k_0$

$\Delta = \text{Bonf.}$

$$\begin{aligned} \mathbb{P}[V \geq 1] &= \mathbb{P}[\text{At least one false discovery}] = \mathbb{P}\left[\bigcup_{j \in I_0} \{\text{rejecting } H_{0j}\} \mid \bigcap_{j \in I_0} H_{0j}\right] \leq \\ &\leq \sum_{j \in I_0} \mathbb{P}[\text{rejecting } H_{0j} \mid H_{0j}] = \sum_{j \in I_0} \frac{\alpha}{K} = k_0 \frac{\alpha}{K} \leq \alpha \end{aligned}$$

Familly wise Error rate: Bonf FWER $\leq \alpha$

Consider $\frac{V}{R}$ with the convention $\frac{V}{R} = 0$ if $R=0$, $Q = \begin{cases} 0 & \text{if } R=0 \\ \frac{V}{R} & \text{if } R>0 \end{cases}$

Definition (FDR)

False discovery rate: $FDR = \mathbb{E}[Q]$

Obs: 1. Assume $k_0 = K$ (no discovery to be made) $\implies S=0 \implies V=R$

$$Q = \begin{cases} 0 & \text{if } R=0 \\ 1 & \text{if } R>0 \end{cases} \iff V=0$$

$$FDR = \mathbb{E}[Q] = \mathbb{P}[V>0] = \mathbb{P}[V \geq 1] = \text{FWER}$$

2. Assume $k_0 < K$:

$$\text{if } V=0 \implies Q=0 = V$$

$$\text{if } V>0 \implies Q = \frac{V}{R} < 1$$

$$\mathbb{E}[V>0] = \mathbb{E}_{\frac{V}{R}}[V>0] \implies Q \leq \mathbb{E}[V>0]$$

$$FDR = \mathbb{E}[Q] \leq \mathbb{E}[\mathbb{E}[V>0]] = \mathbb{P}[V \geq 1] = \text{FWER}$$

}

FDR \leq FWER

Hence: controlling FDR is less restrictive than controlling FWER

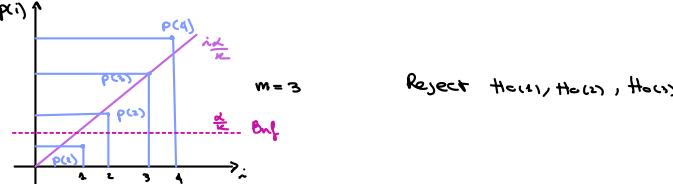


A strategy for controlling FDR

Let p_i : p-value of H_{0i} vs H_{1i}

Consider the k p-values: $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(k)}$
 $\downarrow \quad \downarrow$
 $H_{0(1)} \text{ vs } H_{1(1)} \quad H_{0(k)} \text{ vs } H_{1(k)}$

Let $\alpha \in (0, 1)$. $m = \max \{ i = 1, \dots, k : p(i) \leq i \frac{\alpha}{m} \}$



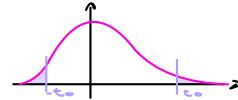
Theorem

If p_1, \dots, p_k are independent. The strategy D_α that rejects $H_{0(i)}$ if $i \leq m$ controls FDR at level α

Obs: p-value: $p = p(x_1, \dots, x_n)$

Example

$$p(x_1, \dots, x_n) = P\left[\frac{|x_1 - \bar{x}|}{\sqrt{s^2/n}} \leq |t| \right]$$



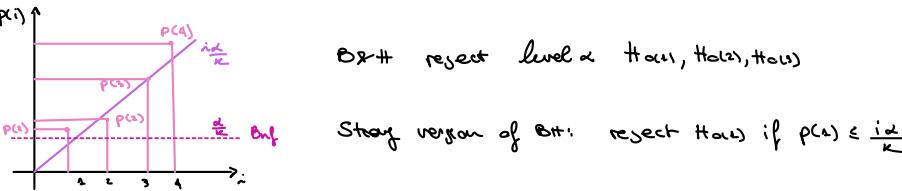
What if p_1, \dots, p_n are NOT independent? (Benjamini & Yekutieli, 2001)

1. p_1, \dots, p_n are positively correlated \implies BH (1995) controls FDR

2. p_1, \dots, p_n are negatively correlated BH? consider this strategy:

reject $H_{0(i)}$ if $i \leq m^* = \max \{ i \in \{1, \dots, k\} : p(i) \leq \frac{i}{k} C(\alpha) \}$ with $C(\alpha) = \sum_{j=1}^k \frac{1}{j}$
 \implies Controls FDR at level α

Obs: It might happen



Comparing means of multivariate Gaussian dist.

Paired data n statistical units obs. twice:

$$\mathbb{R}^p \ni \underline{x}_{1i}, \underline{x}_{2i} \in \mathbb{R}^p \quad i = 1, \dots, n$$

$$\begin{aligned} \mathbb{R}^p &\ni \begin{pmatrix} \underline{x}_{11} \\ \underline{x}_{21} \end{pmatrix}, \begin{pmatrix} \underline{x}_{12} \\ \underline{x}_{22} \end{pmatrix}, \dots, \begin{pmatrix} \underline{x}_{1n} \\ \underline{x}_{2n} \end{pmatrix} \quad \text{iid} \sim N_{2p}\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \Sigma_{2p \times 2p}\right) \\ \mathbb{R}^p &\ni \begin{pmatrix} \underline{x}_{11} \\ \underline{x}_{21} \end{pmatrix}, \dots, \begin{pmatrix} \underline{x}_{1n} \\ \underline{x}_{2n} \end{pmatrix} \quad \text{iid} \sim N_{2p}\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \Sigma_{2p \times 2p}\right) \end{aligned}$$

Attention: x_{1i} and x_{2i} might be dependent!

Typical inf. questions: $H_0: \mu_1 - \mu_2 = \delta \text{ vs } H_1: \mu_1 - \mu_2 \neq \delta \quad (\text{e.g. } \delta = 0)$
 $\text{CR}_{1-\alpha}(\mu_1 - \mu_2)$

Consider $D_i = \underline{x}_{1i} - \underline{x}_{2i}$ iid $N_p(\delta, \Sigma_0) \implies \bar{D} = \frac{1}{n} \sum_{i=1}^n D_i$
 $S_D = \frac{1}{n-1} \sum_{i=1}^n (D_i - \bar{D})(D_i - \bar{D})'$

Note: $n(\bar{D} - \delta)' S_D^{-1} (\bar{D} - \delta) \sim \frac{(n-1)p}{(n-p)} F(p, n-p) \text{ pivotal}$



$$\alpha \in (0, 1) \quad C\bar{\Sigma}_{1-2}(\mu_1 - \mu_2) = CR_{1-2}(\Delta) = \{ \delta \in \mathbb{R}^p : n(\Delta - \delta)^T S\delta^*(\Delta - \delta) \leq \frac{(n-1)p}{n-p} F_{1-2}(p, n-p) \}$$

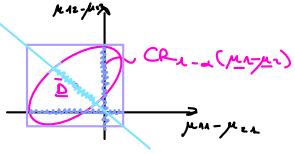
$$H_0: \mu_1 - \mu_2 = \Delta \quad \text{vs} \quad H_A: \mu_1 - \mu_2 \neq \Delta$$

$\alpha \in (0, 1)$ reject H_0 if $n(\Delta - \delta_0)^T S\delta^*(\Delta - \delta_0) > \frac{(n-1)p}{n-p} F_{1-2}(p, n-p)$

$$\delta = \mu_1 - \mu_2$$

$$\text{Sim } C\bar{\Sigma}_{1-2}(\mu_1 - \mu_2) = \left[\bar{\mu}_1^T \bar{\Delta} \pm \sqrt{\frac{(n-1)p}{n-p} F_{1-2}(p, n-p)} \sqrt{\frac{\alpha' S \alpha}{n}} \right]$$

Example



$$\mu_1 - \mu_2 = \Delta \in \mathbb{R}^p$$

$$H_0: \Delta = 0 \quad \text{vs} \quad H_A: \Delta \neq 0$$

Reject H_0 at level α : $\mu_1 \neq \mu_2$

$$\Delta' \mu_1 = \Delta' \mu_2 \quad \Delta \neq 0$$

Repeated Measurements

Univariate case

$$x_1, \dots, x_n \text{ iid } N(\mu, \Sigma)$$

$$\underline{x}_i = (x_{i1}, \dots, x_{iq})^T \in \mathbb{R}^q \quad q \text{ measurements repeated on the same unit } i=1, \dots, n$$

$$\mu = (\mu_1, \dots, \mu_q) \in \mathbb{R}^q$$

$$H_0: \mu_1 = \mu_2 = \dots = \mu_q \quad \text{vs} \quad H_A: \exists i, j \in \{1, \dots, q\} \text{ st. } \mu_i \neq \mu_j$$

$$H_0: \mu \in L(\Delta) \quad \text{vs} \quad H_A: \mu \notin L(\Delta)$$

Contrast matrix

$C \in (q-1) \times q$ is contrast matrix

$$\text{if } C = \begin{bmatrix} c_1 \\ \vdots \\ c_{q-1} \end{bmatrix} \quad c_i \in \mathbb{R}^q$$

st. 1. c_1, \dots, c_{q-1} are linearly indep. $\iff L^1(\Delta) = \text{Span}(c_1, \dots, c_{q-1})$

2. $c_i^T \Delta = 0$ for $i=1, \dots, q-1$

Then we can rewrite the hypothesis: $H_0: C\mu = 0 \quad \text{vs} \quad H_A: C\mu \neq 0$

• $\bar{\Sigma} = \frac{1}{n} \sum_{i=1}^n \Sigma_i$ is unbiased for Σ

$C\bar{\Sigma}$ is unbiased for $C\mu$

• $\bar{\Sigma} \sim N_q(\mu, \frac{1}{n}\Sigma) \Rightarrow C\bar{\Sigma} \sim N_{q-1}(C\mu, \frac{1}{n}C\Sigma C^T)$

• S estimate for Σ : $(n-1)S \sim \text{Wish}(\Sigma, n-1)$

$(n-1)CSC^T \sim \text{Wish}(CC^T, n-1) \perp\!\!\!\perp C\bar{\Sigma}$

Hotelling's T₀ $\Rightarrow n(C\bar{\Sigma} - C\mu)'(CSC^T)^{-1}(C\bar{\Sigma} - C\mu) \sim \frac{(n-1)(q-1)}{n-q+1} F(q-1, n-q+1)$

private hole

$$\alpha \in (0, 1) \quad T_0^2 = n(C\bar{\Sigma})'(CSC^T)^{-1}(C\bar{\Sigma})$$

reject at level α if $T_0^2 > \frac{(n-1)(q-1)}{n-q+1} F_{q-1, n-q+1}$

Obs: C and \tilde{C} two different contrast matrices:

$$\Rightarrow B \in (q-1) \times (q-1), \text{ Det}(B) \neq 0 \quad \text{and} \quad C = B\tilde{C}$$

$$\begin{aligned} T_0^2 &= n(C\bar{\Sigma})'(CSC^T)^{-1}(C\bar{\Sigma}) = n(B\tilde{C}\bar{\Sigma})'(B\tilde{C}S\tilde{C}^T)^{-1}(B\tilde{C}\bar{\Sigma}) = \\ &= n(\tilde{C}\bar{\Sigma})' \tilde{B}'(B^T)^{-1} \tilde{B}^T \tilde{B}(\tilde{C}\bar{\Sigma}) = \tilde{T}_0^2 \end{aligned}$$

Ex

$$C = \begin{bmatrix} 1 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & \dots & -1 \end{bmatrix} \quad \tilde{C} = \begin{bmatrix} 1 & -1 & 0 & \dots & 0 \\ 1 & 0 & -1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 & -1 \end{bmatrix}$$



Obs: $\underline{x}_1, \dots, \underline{x}_n$ iid $\sim N_p(\mu, \Sigma)$

$H_0: \mu \in \mathcal{L} \subseteq \mathbb{R}^p$ $H_1: \mu \notin \mathcal{L}$

\mathcal{L} = linear subspace of \mathbb{R}^p , $\dim(\mathcal{L}) = k \geq 1$

Assume: $\mathcal{L}^\perp = \text{span}(c_1, \dots, c_{p-k})$

$$\text{Let } C = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_{p-k} \end{bmatrix} \quad (p-k) \times p$$

$H_0: C\mu = 0$ vs $H_1: C\mu \neq 0$

C is unbiased for $C\mu$

Hotelling's theorem: if H_0 is true $T_0^2 = n(C\bar{\underline{x}})'(CSC')^{-1}(C\bar{\underline{x}}) \sim \frac{(n-1)(p-k)}{n-p+k} \mathcal{F}_{p-k, n-p+k}$

Multivariate cases

unit i : $\begin{pmatrix} x_{i1}(h) \\ x_{i2}(h) \\ \vdots \\ x_{ik}(h) \end{pmatrix}, \begin{pmatrix} x_{i1}(w) \\ x_{i2}(w) \\ \vdots \\ x_{ik}(w) \end{pmatrix}, \dots, \begin{pmatrix} x_{i1}(g) \\ x_{i2}(g) \\ \vdots \\ x_{ig}(g) \end{pmatrix} \quad i = 1, \dots, n$ units are independent

$\begin{cases} H_0: \mu_1(h) = \mu_1(w) = \dots = \mu_1(g) \quad \text{and} \quad \mu_2(h) = \mu_2(w) = \dots = \mu_2(g) \\ H_1: H_0^c \end{cases}$

change of rep: $\underline{x}_i = (x_{i1}(h), x_{i2}(h), \dots, x_{ik}(h), x_{i1}(w), \dots, x_{ig}(w))' \in \mathbb{R}^{2g}$

