Homework #2 Stat 212A, Fall 2015: Topics in Selective Inference Instructor: Will Fithian

Assigned Oct. 19, 2015. Due 11:59pm Nov. 5, 2015

You are welcome to work with each other or consult articles or textbooks online, but you should then go away and write up the problem by yourself. If you collaborate or use other resources, please list your collaborators and cite the resources you used. Please show your work and include code where appropriate.

You can turn in the problem set in class Nov. 5 or under my door (Evans 301) Thursday night.

1. FDR Control for Arbitrary Dependence Consider the usual setup with m p-values p_1, \ldots, p_m and null hypotheses $H_i : p_i \sim U[0, 1]$, true for $i \in I_0$.

In class we stated, but did not prove, the theorem from Benjamini and Yekutieli (2001) that under any dependendency structure among the *p*-values, the BH procedure at level α controls FDR at

$$\operatorname{FDR} \le \alpha \sum_{k=1}^m \frac{1}{k}$$

In this problem we will prove the theorem. The key to the proof is the quantity $\pi_{i,j,r}$, defined as

$$\pi_{i,j,r} = \mathbb{P}\left(R^i = r, \ \frac{(j-1)\alpha}{m} \le p_i \le \frac{j\alpha}{m}\right).$$

Recall the definition of R^i as the number of hypotheses BH would reject if we replaced p_i with 0, so that $R^i \perp y_i$. Also recall our lemma that

$$R = r, i \in S \iff R^i = r, p_i \le \frac{\alpha r}{m}$$

(a) Show that

$$FDR = \sum_{i \in I_0} \sum_{r=1}^m r^{-1} \sum_{j=1}^r \pi_{i,j,r}$$

(b) Continuing, show that we can exchange the indices j and r to obtain

FDR
$$\leq \sum_{i \in I_0} \sum_{j=1}^m j^{-1} \sum_{r=j}^m \pi_{i,j,r}$$

(Note: be sure to justify replacing r^{-1} with j^{-1} .)

(c) Show that $\sum_{r=1}^{m} \pi_{i,j,r} \leq \alpha/m$ and use this to complete the proof.

2. BY Intervals for Marginal Screening For the construction of the Benjamini-Yekutieli (BY) intervals in class, we defined

$$R_{\min}^{i} = \min\left\{ |S(y_{1}, \dots, y_{i-1}, x, y_{i+1}, \dots, y_{m})| : x \in \mathbb{R}, i \in S(y_{1}, \dots, y_{i-1}, x, y_{i+1}, \dots, y_{m}) \right\}.$$

We showed in particular that for the Benjamini-Hochberg selection procedure, we always have $R_{\min}^i = R$ for every $i \in S(y)$.

A simpler selection algorithm we could use is marginal screening: selecting all indices for which $|y_i|$ surpasses some fixed threshold t. That is,

$$S(y_1, \dots, y_m) = \{ i \in \{1, \dots, m\} : |y_i| > t \}.$$

- (a) Show that for this selection procedure, we also have $R_{\min}^i = R$ for every $i \in S(y)$.
- (b) What is $R_{\min}^i R$ for the unselected indices $i \notin S(y)$?

3. Convergence of Conditional to Nominal Intervals In class we defined the equal-tailed conditional intervals for the truncated $N(\mu, 1)$ distribution conditional on |y| > t for a fixed value of t > 0. If $k_1(\mu)$ and $k_2(\mu)$ are the lower and upper $\alpha/2$ quantiles of the conditional distribution with density

$$f^t_{\mu}(y) = \frac{e^{-(y-\mu)^2/2}}{\sqrt{2\pi}} \cdot \frac{1\{|y| > t\}}{\mathbb{P}_{\mu}(|y| > t)},$$

we defined the conditional intervals as

$$C(y;t) = \left[k_2^{-1}(y), k_1^{-1}(y)\right].$$

Consider taking a limit with fixed t > 0 and $y \to \infty$. Show that the conditional interval tends to the nominal interval $y \pm z_{\alpha/2}$. That is, prove that

$$\lim_{y \to \infty} \left| k_2^{-1}(y) - (y - z_{\alpha/2}) \right| = \lim_{y \to \infty} \left| k_1^{-1}(y) - (y + z_{\alpha/2}) \right| = 0.$$

4. Computing Conditional Intervals For the threshold t = 4, compute the equal-tailed conditional interval for y = 4.001, 4.1, 5, and 10 for $y \sim N(\mu, 1)$ given |y| > t. Show your code.

Hint: if you can implement $k_1(\mu)$ and $k_2(\mu)$ as functions then you can invert them via numerical root-finding.

5. A Bit of Philosophy (Note: Graded for completion only; write as little or as much as you want, but write something.)

In class we have discussed the contrast between the "full model" viewpoint and the "submodel viewpoint" of PoSI. Come up with a concrete applied example (other than the "coefficient of IQ on salary" example I gave in class) where you think the submodel viewpoint seems more appropriate, and another concrete applied example where you think the full model viewpoint seems more appropriate. Here, "appropriate" really means "aligned with the actual scientific goals of the analysis."

References

Yoav Benjamini and Daniel Yekutieli. The control of the false discovery rate in multiple testing under dependency. *Annals of statistics*, pages 1165–1188, 2001.