

## Lecture 4: The post-hoc inference problem

May 8, 2019

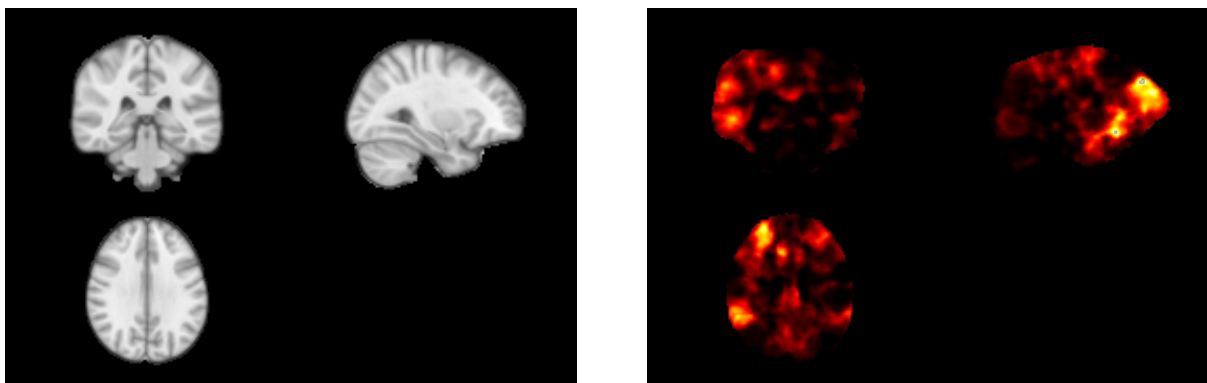
*Lecturer: Aldo Solari*

*Closed testing* (Marcus et al., 1976) is a fundamental principle of FWER control. [3] showed that closed testing can be used to obtain simultaneous confidence bounds for the false discovery proportion (FDP). Used in this way, closed testing allows a form of *post-selection inference*. Theoretical results for the special case of Simes local tests are discussed in [1]. Applications to genomics and fMRI data are discussed in [2] and [4], respectively.

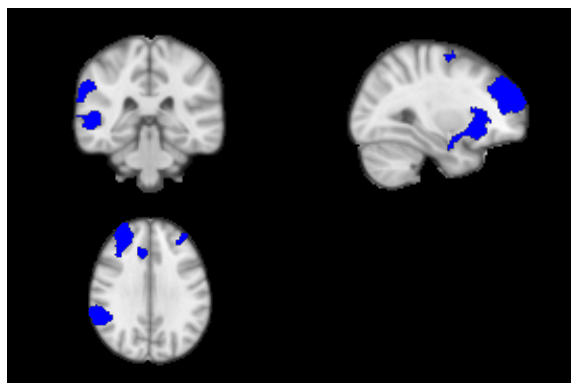
### 1 fMRI data

Analysis of fMRI data supplies an *activation map*: a  $p$ -value testing the null hypothesis of ‘no activation’ (e.g.  $H_i : \mu_i = 0$ ) at each location (*voxel*).

The goal is to find brain regions of activations. Note that the data is at higher resolution than units of interest: we have  $p$ -values at the voxel level, but we wish to perform inference at the region level.



We can select interesting regions by looking at the data (activation map):



Now the problem is the following: *How to assess the significance of selected regions?* Regions are both selected and tested with the same data. We need to correct for the overoptimism in inference due to data-driven selection.

## 1.1 Classical multiple testing

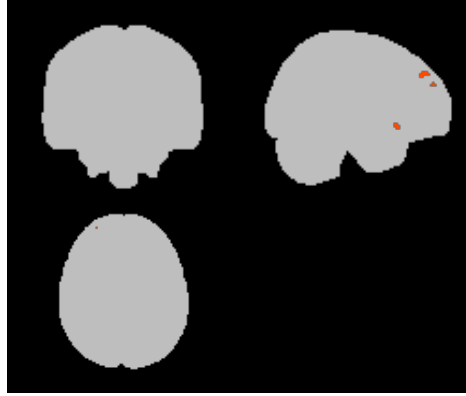
Suppose we have  $m$  hypotheses  $H_1, \dots, H_m$  with corresponding  $p$ -values  $p_1, \dots, p_m$ . Denote by  $T \subseteq \{1, \dots, m\}$  and  $F = \{1, \dots, m\} \setminus T$  the index sets of true and false hypotheses, respectively. For a selection  $S \subseteq \{1, \dots, m\}$ , let

$$\pi_1(S) = \frac{|F \cap S|}{|S|}$$

be the proportion of false hypotheses (true activations) in  $S$ , where  $|\cdot|$  denotes the size of a set.

If we use Bonferroni method at level  $\alpha$ , we reject all the hypotheses with indices in

$$R_\alpha = \{i : p_i \leq \alpha/m\}$$



Then the proportion of true activations in a selected region  $S$  can be estimated by

$$\hat{\pi}_1(S) = \frac{|R_\alpha \cap S|}{|S|}$$

with

$$P(\pi_1(S) \leq \hat{\pi}_1(S) \quad \forall S) \geq 1 - \alpha$$

FWER has the *subsetting property* that if a set  $R$  of hypotheses is rejected by an FWER-controlling procedure, then FWER control is also guaranteed for any subset  $S \subset R$ , i.e.

$$P(R \cap T = \emptyset) \geq 1 - \alpha \Rightarrow P(S \cap T = \emptyset) \geq 1 - \alpha \quad \forall S \subseteq R$$

This property does not hold for FDR control. If FDR control holds for the full set  $R$  only, this does not translate to any subset  $S$ , i.e.

$$E\left(\frac{|R \cap T|}{|R| \vee 1}\right) \leq \alpha \not\Rightarrow E\left(\frac{|S \cap T|}{|S| \vee 1}\right) \leq \alpha$$

## 2 Closed testing

Let  $\theta \in \Theta$  be a parameter of interest. Suppose we have  $m$  hypotheses  $H_1, \dots, H_m$  with  $H_i \subset \Theta$  true if and only if  $\theta \in H_i$ . Denote by  $T \subseteq \{1, \dots, m\}$  the index set of true hypotheses.

In the closed testing procedure, the collection of hypotheses is augmented with all possible intersection hypotheses

$$H_I = \bigcap_{i \in I} H_i$$

with  $I \subseteq \{1, \dots, m\}$ . An intersection hypothesis  $H_I$  is true if and only if  $H_i$  is true for all  $i \in I$ . Note that  $H_i = H_{\{i\}}$ , so all original hypotheses, known as *elementary hypotheses*, are also intersection hypotheses.

The closed testing procedure starts by testing all intersection hypotheses with a *local test*, i.e. an  $\alpha$ -level test for  $H_I$ :

$$\sup_{\theta \in H_I} P_\theta(\phi_I = 1) \leq \alpha$$

Let the collection of all subsets of  $\{1, \dots, m\}$  be denoted by

$$\mathcal{I} = 2^{\{1, \dots, m\}}$$

and the collection of index sets corresponding to true intersection hypotheses by

$$\mathcal{T} = \{I \in \mathcal{I} : I \subseteq T\}$$

We define  $\mathcal{U}_\alpha$  as the collection of  $I \in \mathcal{I}$  such that  $H_I$  is rejected by a local test at level  $\alpha$

$$\mathcal{U}_\alpha = \{I \in \mathcal{I} : \phi_I = 1\}$$

For each  $I$ ,  $H_I$  is rejected by the closed testing procedure if and only if  $I \in \mathcal{X}_\alpha$  where

$$\mathcal{X}_\alpha = \{I \in \mathcal{I} : J \in \mathcal{U}_\alpha \ \forall J \supseteq I\}$$

**Theorem 2.1.** *The closed testing procedure controls the FWER for all hypotheses  $H_I$  at level  $\alpha$ , i.e.*

$$P(\mathcal{T} \cap \mathcal{X}_\alpha = \emptyset) \geq 1 - \alpha$$

*Proof.* We have

$$\{T \notin \mathcal{U}_\alpha\} \subseteq \{I \notin \mathcal{X}_\alpha \ \forall I \subseteq T\} \subseteq \{\mathcal{T} \cap \mathcal{X}_\alpha = \emptyset\}$$

and because  $H_T$  is tested by an  $\alpha$  level test

$$1 - \alpha \leq P(T \notin \mathcal{U}_\alpha) \leq P(\mathcal{T} \cap \mathcal{X}_\alpha = \emptyset)$$

□

## 2.1 Closed testing for FWER control

Let

$$R_\alpha = \{i : \{i\} \in \mathcal{X}_\alpha\}$$

be the index set of elementary hypotheses rejected by the closed testing procedure. Then  $P(\mathcal{T} \cap \mathcal{X}_\alpha = \emptyset) \geq 1 - \alpha$  implies

$$P(T \cap R_\alpha = \emptyset) \geq 1 - \alpha$$

Local tests tend to be easy to specify in most models, as each local test is a test of a single hypothesis, so that standard statistical test theory may be used. Below some examples:

- Bonferroni local test: reject  $H_I$  if and only if

$$\bigcup_{i \in S} \left\{ p_i \leq \frac{\alpha}{|I|} \right\}$$

The CT procedure based on Bonferroni local tests gives Holm's method

- Simes local test: reject  $H_I$  if and only if

$$\bigcup_{i \in S} \left\{ p_{(i:I)} \leq \frac{i\alpha}{|I|} \right\}$$

where  $p_{(i:I)}$  is the  $i$ th smallest  $p$ -value among the multiset  $\{p_i : i \in I\}$ . The Simes local test is level  $\alpha$  if we assume the Simes' inequality, i.e.

$$P\left(\bigcup_{i \in T} \left\{ p_{(i:T)} \leq \frac{i\alpha}{|T|} \right\}\right) \geq 1 - \alpha$$

The CT procedure based on Simes local tests gives the Hommel method: compute

$$h_\alpha = \max\{i \in \{0, \dots, m\} : ip_{(m-i+j)} > j\alpha \text{ for } j = 1, \dots, i\}$$

and

$$R_\alpha = \{i : p_i \leq \alpha/h_\alpha\}$$

- Fisher local test: reject  $H_I$  if and only if

$$-2 \sum_{i \in I} \log(p_i) \geq c_{|I|}$$

where  $c_r$  is the  $1 - \alpha$  quantile of a  $\chi^2$  distribution with  $2r$  degrees of freedom. Fisher's combination test is level  $\alpha$  if we assume that the null  $p$ -values are independent.

### 3 Simultaneous confidence bounds for the FDP

Let  $S \subseteq \{1, \dots, m\}$  be the index set of selected hypotheses, the discoveries. The number of false discoveries for  $S$  is

$$\tau(S) = |T \cap S|$$

[3] show how to calculate upper confidence bounds  $\hat{\tau}_\alpha$  for  $\tau(S)$  from the closed testing procedure:

$$\hat{\tau}_\alpha(S) = \max\{|I| : I \subseteq S, I \notin \mathcal{R}_\alpha\}$$

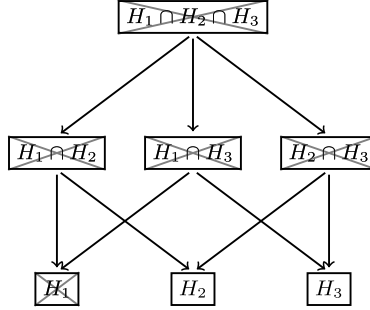
[3] prove that

$$P(\tau(S) \leq \hat{\tau}_\alpha(S) \ \forall S \in \mathcal{I}) \geq 1 - \alpha$$

meaning that the set

$$\{0, \dots, \hat{\tau}_\alpha(S)\}$$

is a  $(1 - \alpha)$ -confidence set of the parameter  $\tau(S)$ . The quantity  $\hat{\tau}_\alpha(S)$  is the size of the largest subset of  $S$  for which the corresponding intersection hypothesis is not rejected by the closed testing procedure. For, example, if the rejected hypotheses by the closed testing procedure are those marked with a cross



we obtain

$S$	Confidence set for $\tau(S)$
$\{1\}$	$\{0\}$
$\{2\}$	$\{0, 1\}$
$\{3\}$	$\{0, 1\}$
$\{1, 2\}$	$\{0, 1\}$
$\{1, 3\}$	$\{0, 1\}$
$\{2, 3\}$	$\{0, 1\}$
$\{1, 2, 3\}$	$\{0, 1\}$

To prove the coverage, remember that if the event  $E = \{T \notin \mathcal{U}_\alpha\}$  has happened, then all rejections that the closed testing procedure has made are correct. Given that  $E$  has happened, the value of  $\tau(S)$  cannot be greater than the value of  $\hat{\tau}_\alpha(S)$ , because otherwise a true intersection hypothesis would have been rejected, which is inconsistent with the definition of  $E$ .

The important thing to note about confidence sets is that they are simultaneous confidence sets, which all depend on exactly the same event  $E$  for their coverage. Because these confidence sets are simultaneous, the user can review all these confidence sets, and select the one that he or she likes best, while still keeping correct  $1 - \alpha$  coverage of the selected confidence set: under the event  $E$ , all confidence sets cover the true parameter simultaneously, and therefore, under the same event  $E$ , the selected confidence set covers the true parameter. Consequently, the selected confidence set has coverage  $P(E) \geq 1 - \alpha$ . The simultaneity of the sets makes their coverage robust against post hoc selection.

The bounds can be equivalently formulated in terms of the false discovery proportion  $\pi_0(S) = |T \cap S|/|S|$  as

$$\pi_0(S) = \frac{\hat{\tau}_\alpha(S)}{|S|}$$

or in terms of the true discovery proportion  $\pi_1(S) = |F \cap S|/|S|$  as

$$\pi_1(S) = \frac{|S| - \hat{\tau}_\alpha(S)}{|S|}$$

## References

- [1] J. Goeman, R. Meijer, T. Krebs, and A. Solari. Simultaneous control of all false discovery proportions in large-scale multiple hypothesis testing. *Biometrika (to appear)*, 2019.
- [2] J. J. Goeman and A. Solari. Multiple hypothesis testing in genomics. *Statistics in medicine*, 33(11):1946–1978, 2014.
- [3] J. J. Goeman, A. Solari, et al. Multiple testing for exploratory research. *Statistical Science*, 26(4):584–597, 2011.
- [4] J. D. Rosenblatt, L. Finos, W. D. Weeda, A. Solari, and J. J. Goeman. All-resolutions inference for brain imaging. *NeuroImage*, 181:786–796, 2018.