

# Large-scale inference for the detection of sources in MUSE hyperspectral data. Towards robust error control.



**gipsa-lab**

F. Chatelain

*Joint work with R. Bacher, C. Meillier, O. Michel*

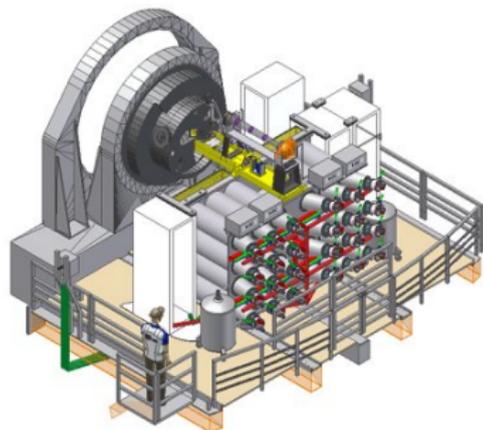


JATIA, Strasbourg, January 25, 2019

# MUSE instrument

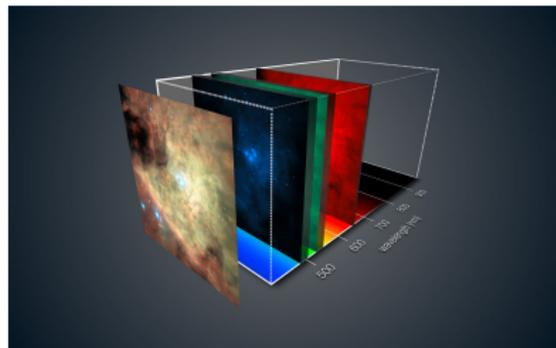
Instrument for ESO, VLT, Chili (First light in 2014) :

- ▶ [2D +  $\lambda$   $\equiv$  3D] imager = Integral field spectrograph



Observation of distant galaxies (thus, very young), and their possible halos

- ▶ dramatically faint except on a few characteristic lines
- 🔭 understanding of universe, galaxy formation...



## Data Cube

Stack of  $\sim 3600$  monochromatic images covering  $60 \times 60$  arcsec

- ▶ spatial resolution :  $.2 \times .2$  arcsec ( $300 \times 300$  pixels)
- ▶ spectral resolution :  $.14\text{nm}$  (spectral range :  $465 - 930$  nm)

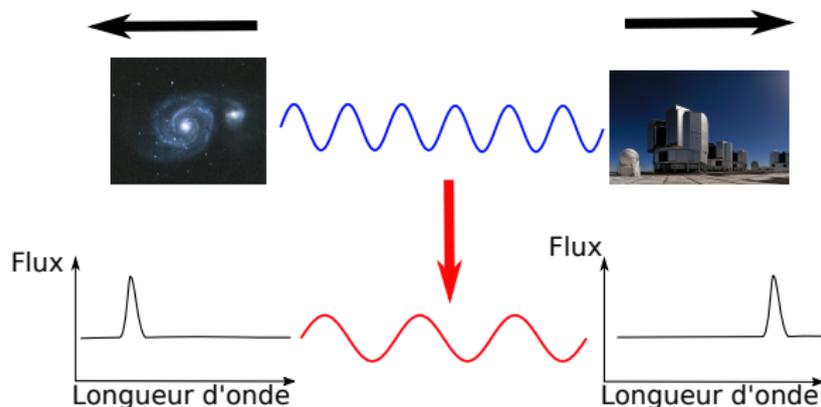
Data Cube  $300 \times 300 \times 3600$

## Redshift and detection

We want to detect 1) faint galaxies, or 2) galactic halos (hydrogen gas surrounding galaxies)

- ▶ Emission limited to a few wavelengths : Lyman- $\alpha$  emission line
- ▶ ... of unknown spectral position because of redshift

**Redshift** : during its trip to Earth, light emitted by a galaxy moving away from us (Universe expansion...) is shifted to the red (remember the ambulance!).

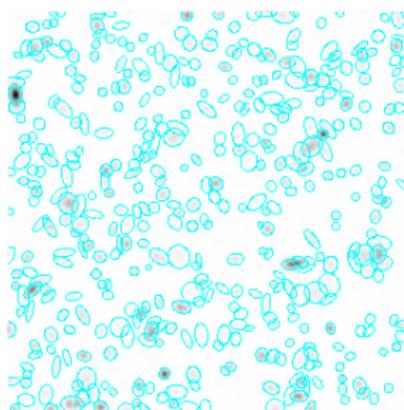


☞ Calls for detection methods adapted to these large datasets.

# 1) Detection of faint galaxies : C. Meillier's PhD

## Problem

Detect faint galaxies whose position, shape, spectrum, power, number... are unknown.



## Bayesian Nonparametric approach : (marked) point process<sup>1</sup>

- ▶ Object (galaxy) = a point (position) + marks (geometric  $\approx$  elliptical object, and spectral parameters)
- ▶ Object configuration = realization of a marked point process
- ↳ naturally sparse representation of massive data fields : configuration of marked points + noise

---

1. Meillier *et al*, *IEEE TSP* (2015)

## 1) Detection of faint galaxies : C. Meillier's PhD (Cont'd)

### SELFIE Results on HDFS data<sup>2</sup> : comparison with MUSE and Hubble (HST) catalog

Total number of detected objects	298
Number of detected objects belonging to MUSE catalog	166 / 189
Number of detected objects belonging to HST catalog	(166+78)
Number of detected objects not belonging to any catalog	54
... including potential galaxies	6

### How to assess the significance of the detection list ?

- ▶ no ground truth to assert the detection performance !
- ☞ need for a robust error control for these multiple inferences

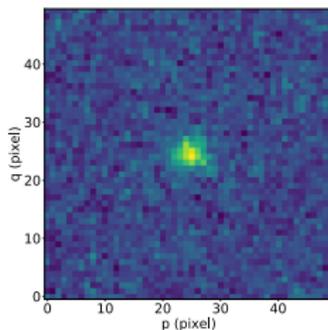
---

2. Meillier *et al*, A&A (2016)

## 2) Detection of galactic halo : R. Bacher's PhD

We have  $n$  pixels (e.g.  $n = 2500$  for a  $50 \times 50$  neighborhood) to test for :

- ▶ Which pixels have signal ? / Which pixels belong to the galactic halo ?



How to have guarantees on the detection results ?

- ▶ no ground truth to assert the detection performance !
- 👁️ need for a robust control, e.g. to guaranty the proportion of pixels among the detected set that are really part of the target (“purity” of the detection)

# Outline

## Introduction and Motivations

- MUSE instrument

- Two detection problems

## Multiple inference and Global error control

- Multiple comparisons

- False Discovery Rate FDR

- BH Procedure

## Detection of galactic sources : CGM

- CGM multiple testing

- COMET procedure

- COMET Results

## Conclusion and perspectives

# Multiplicity problem and chance correlation

## Lottery



- ▶ Winning probability for a given ticket is very low...
- ▶ But among the huge number of tickets, the probability that there is *at least one* winning ticket is quite high !

☞ **Large-scale experiments** : multiplying the comparisons dramatically increases the probability to obtain a good match by **pure chance**

## Paul the octopus



- ▶ Paul predicts eight of the 2010 FIFA World Cup matches with a perfect score !
- ▶ Does it really mean that Paul is an Oracle ?

# Multiplicity problem for statistical testing

- ▶  $T$  is the test statistics,
- ▶  $\mathcal{R}_\alpha$  is the region of rejection at level  $\alpha$  : if  $H_0$  is true,  $\Pr(T \in \mathcal{R}_\alpha) = \alpha$

## Multiple testing issue

- ▶  $N$  independent statistics  $T_1, \dots, T_N$  obtained under the null  $H_0$
- ▶ Probability to reject *at least one* of the  $N$  null hypotheses :

$$\begin{aligned}\Pr(\exists T_i \in \mathcal{R}_\alpha) &= 1 - \Pr(T_1, \dots, T_N \notin \mathcal{R}_\alpha) = 1 - \prod_{i=1}^N \Pr(T_i \notin \mathcal{R}_\alpha), \\ &= 1 - \prod_{i=1}^N (1 - \alpha) = 1 - (1 - \alpha)^N\end{aligned}$$

- ▶ for a usual significant level  $\alpha = 0.05$ , performing  $N = 20$  tests gives a probability 0.64 to find a 'significant' discovery by pure chance...
- ☞  $\Pr(\text{at least one false positive}) \gg \Pr(\text{the } i\text{-th is a false positive})$

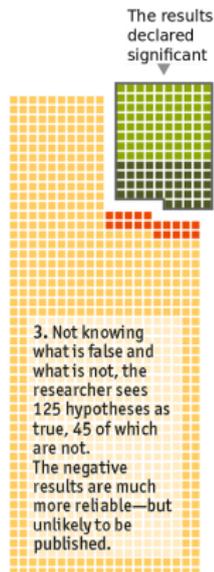
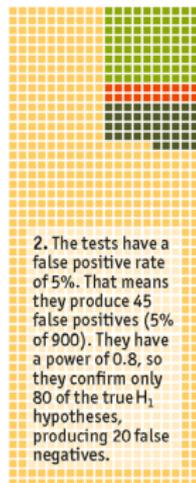
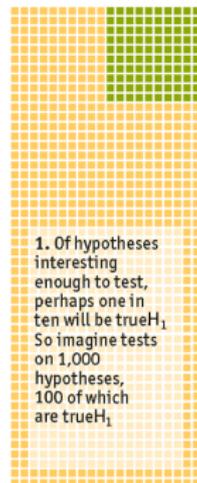
# Multiplicity problem in science

## *The Economist*, 2013, “Unreliable research”

### Unlikely results

How a small proportion of false positives can prove very misleading

False  $H_1$  True  $H_1$  False negatives False positives



Many published research findings in top-ranked journals are not, or poorly, reproducible [Ioannidis, 2005]

Source: *The Economist*

- ▶ if the test power is only 0.4, 40 true positives in average for 45 false positives. Is this significant ?

# Large-Scale Hypothesis Testing [Efron, 2010]

## Era of Massive Data Production

- ▶ “omics” revolution, e.g. microarrays measures expression levels of tens of thousands of genes for hundreds of subjects
- ▶ astrophysics, e.g. MUSE spectro-imager delivers cubes of  $300 \times 300$  images for 3600 wavelengths : detecting faint sources leads to  $N \approx 3 \times 10^8$  tests in a pixelwise approach

## Large-Scale methodology

- ▶ statistical inference and hypothesis testing theory developed in the early 20th century (Pearson, Fisher, Neyman, . . . ) for small-data sets collected by individual scientist
- ✚ corrections are needed to assess significance in large-scale experiments

# P-values : an universal language for hypothesis testing

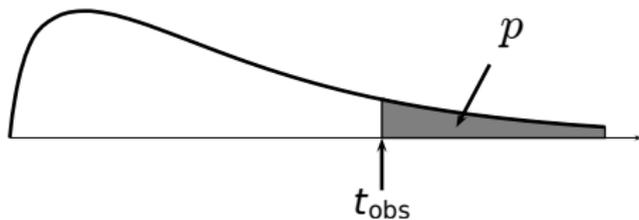
## Intuitive definition

*p-value*  $\equiv$  probability of obtaining a result as extreme or “more extreme” than the observed statistics, under  $H_0$

## One-sided test example

- ▶  $T$  is the test statistic,  $t_{\text{obs}}$  an observed realization of  $T$
- ▶  $H_0$  rejected when  $t_{\text{obs}}$  is too large :  $\mathcal{R}_\alpha = \{t : t \geq \eta_\alpha\}$

$$p(t_{\text{obs}}) = \Pr_{H_0}(T \geq t_{\text{obs}})$$



## Mathematical definition

Smallest value of  $\alpha$  such that  $t_{\text{obs}} \in \mathcal{R}_\alpha$

$$p(t_{\text{obs}}) = \inf_{\alpha} \{t_{\text{obs}} \in \mathcal{R}_\alpha\}$$

## Property of $p$ -values

Let  $P = p(T)$  be the random variable. If  $H_0$  is true

$$\Pr_{H_0}(P \leq u) = \Pr_{H_0}(T \in \mathcal{R}_u) = u,$$

- ↳  $p$ -value  $\equiv$  transformation of the test statistics to be uniformly distributed under the null (whatever the distribution of  $T$ )

### Statistical hypothesis test based on $p$ -value

$H_0$  :  $p$ -value has a **uniform distribution** on  $[0, 1]$  :  $P \sim \mathcal{U}([0, 1])$

$H_1$  :  $p$ -value is *stochastically lower* than  $\mathcal{U}([0, 1])$  :  $\Pr_{H_1}(P \leq u) = \Pr_{H_1}(T \in \mathcal{R}_u) > u$ ,

- ↳ the smaller is  $p \equiv p(t_{\text{obs}})$ , the more decisively is  $H_0$  rejected
- ↳ for a given  $\alpha$ ,  $H_0$  is rejected at level  $\alpha$  if  $p \leq \alpha$

## Counting the errors in multiple testing

- ▶  $N$  hypothesis tests with a common procedure

		Decision		Total
		$H_0$ retained	$H_0$ rejected	
Actual	$H_0$ true	$V$	$U$	$N_0$
	$H_0$ false	$S$	$T$	$N_1$
Total		$N - R$	$R$	$N$

- ▶  $N_0 = \#$  true nulls,  $N_1 = \#$  true alternatives
- ▶  $U = \#$  False Positives  $\leftarrow$  Type I Errors
- ▶  $T = \#$  True Positives,
- ▶  $R = \#$  Rejections

How to define, and control, a global Type I Error rate/criterion ?

# False Discovery Rate FDR [Benjamini and Hochberg, 1995]

## “Discovery” terminology

- ▶  $R \equiv \#$  Discoveries (Detections or Positives)
- ▶  $U \equiv \#$  False Discoveries (False Positives) ← Type I errors,
- ▶  $T \equiv \#$  True Discoveries (True Positives),

		Decision		Total
		$H_0$ retained	$H_0$ rejected	
Actual	$H_0$ true	$V$	$U$	$N_0$
	$H_0$ false	$S$	$T$	$N_1$
Total		$N - R$	$R$	$N$

## Definition

$FDP \equiv \frac{U}{R \vee 1}$ , where  $R \vee 1 \equiv \max(R, 1)$  ← False Discovery Proportion

$FDR \equiv E[FDP] = E\left[\frac{U}{R \vee 1}\right]$  ← False Discovery Rate

- ⚠ single test errors (e.g. PFA controls in average the  $U/N_0$  ratio), or power, are calculated horizontally in the table
- ⚠ False Discovery Rate is calculated vertically (Bayesian flavor)

# False Discovery Rate FDR [Benjamini and Hochberg, 1995]

## “Discovery” terminology

- ▶  $R \equiv \#$  Discoveries (Detections or Positives)
- ▶  $U \equiv \#$  False Discoveries (False Positives)  $\leftarrow$  Type I errors,
- ▶  $T \equiv \#$  True Discoveries (True Positives),

		Decision		Total
		$H_0$ retained	$H_0$ rejected	
Actual	$H_0$ true	$V$	$U$	$N_0$
	$H_0$ false	$S$	$T$	$N_1$
Total		$N - R$	$R$	$N$

## Definition

$FDP \equiv \frac{U}{R \vee 1}$ , where  $R \vee 1 \equiv \max(R, 1) \leftarrow$  False Discovery Proportion

$FDR \equiv E[FDP] = E\left[\frac{U}{R \vee 1}\right] \leftarrow$  False Discovery Rate

- ⓘ single test errors (e.g. PFA controls in average the  $U/N_0$  ratio), or power, are calculated horizontally in the table
- ⓘ False Discovery Rate is calculated vertically (Bayesian flavor)

# False Discovery Rate FDR [Benjamini and Hochberg, 1995]

## “Discovery” terminology

- ▶  $R \equiv \#$  Discoveries (Detections or Positives)
- ▶  $U \equiv \#$  False Discoveries (False Positives)  $\leftarrow$  Type I errors,
- ▶  $T \equiv \#$  True Discoveries (True Positives),

		Decision		Total
		$H_0$ retained	$H_0$ rejected	
Actual	$H_0$ true	$V$	$U$	$N_0$
	$H_0$ false	$S$	$T$	$N_1$
Total		$N - R$	$R$	$N$

## Definition

$FDP \equiv \frac{U}{R \vee 1}$ , where  $R \vee 1 \equiv \max(R, 1) \leftarrow$  False Discovery Proportion

$FDR \equiv E[FDP] = E\left[\frac{U}{R \vee 1}\right] \leftarrow$  False Discovery Rate

- ☞ single test errors (e.g. PFA controls in average the  $U/N_0$  ratio), or power, are calculated horizontally in the table
- ☞ False Discovery Rate is calculated vertically (Bayesian flavor)

# Source detection example

## Multiple testing problem

Statistical linear model (source + noise) for each  $i = 1, \dots, N$

$$X_i = \mu r_i + \epsilon_i$$

with  $\mu > 0$ ,  $r_i \in \{0, 1\}$ ,  $\epsilon_i \sim \mathcal{N}(0, 1)$

- ▶  $H_0$  : null hypothesis  $\equiv$  absence of signal, i.e.  $r_i = 0$
- ▶  $H_1$  : alternative hypothesis  $\equiv$  presence of signal, i.e.  $r_i = 1$

## Test statistics

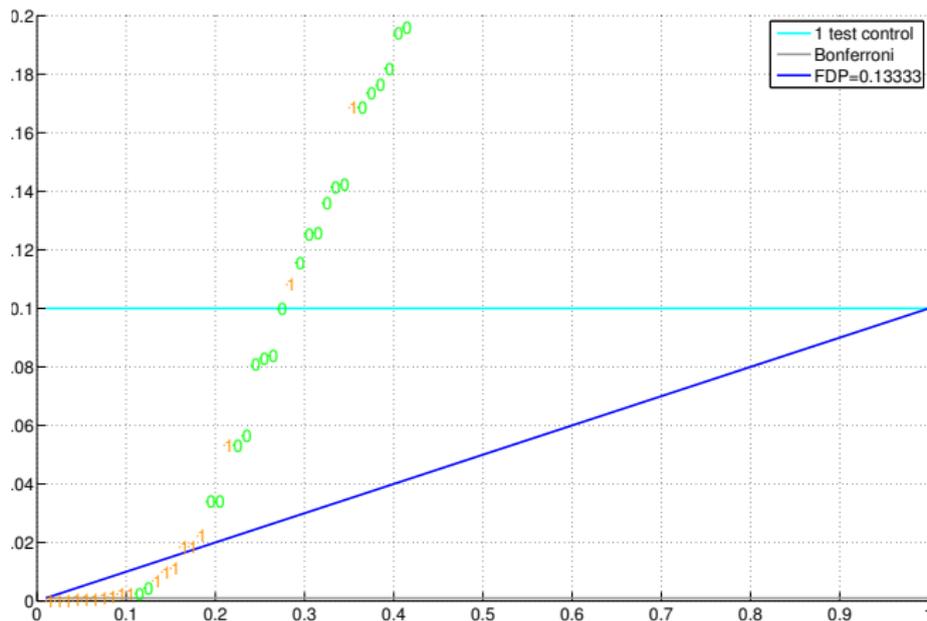
for each  $i$

- ▶  $X_i$  is the test statistics
- ▶  $p_i = 1 - \Phi(X_i)$ , where  $\Phi$  is the standard normal cdf, is the associated p-value

How to choose a good threshold  $t$  to reject the tests s.t.  $p_i \leq t$ ?

## Ordered p-values plot for $N = 100$ , $N_0 = 80$ , $\mu = 3$ , $\alpha = 0.1$

Try something between Bonferroni and single test control : choose  $t_i = q \frac{i}{N}$  (here  $q = \alpha = 0.1$ )



Ordered p-values  $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(N)}$  vs theoretical quantiles  $1/N, 2/N, \dots, 1$  under the null

# Benjamini-Hochberg (BH) procedure

## BH procedure<sup>3</sup>

- ▶ Ordered p-values  $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(N)}$ , let  $p_{(0)} = 0$  by convention
- ▶ For a given FDR control level  $0 \leq q \leq 1$  :
  - ▶ find the largest  $\hat{k}$  s.t.  $p_{(k)} \leq q \frac{k}{N}$
  - ▶ reject  $H_0$  for all  $p_{(i)}$ ,  $i = 1, \dots, \hat{k}$

## Theorem

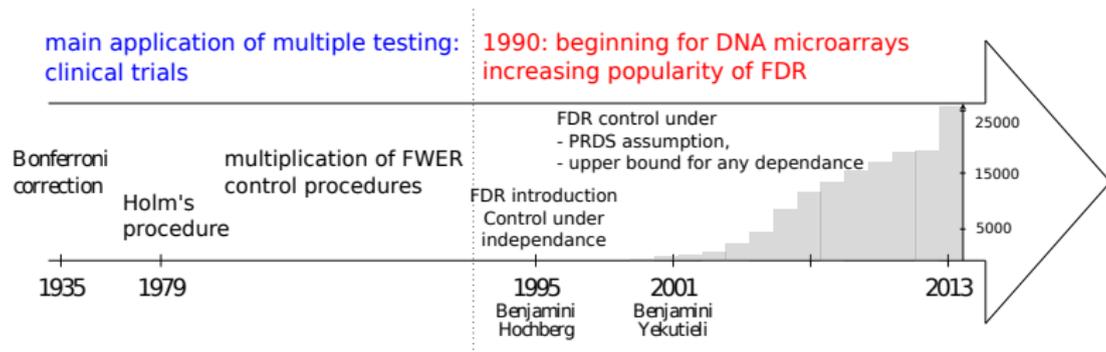
Under the independence assumption (or specific positive dependence) among the tests, BH procedure controls the FDR at level  $q$ .

- 📖 *learning from the other experiments idea*
- 📖 “testimation problem” : blurs the line between testing and estimation

---

3. Benjamini and Hochberg, *JRSS, Series B* (1995)

# Popularity of FDR and BH procedure



Historical context and citations<sup>4</sup> of the seminal paper [Benjamini and Hochberg, 1995]

## FDR for Big Data

Large-scale hypothesis testing in many fields

- ▶ DNA microarray, genomics, fMRI data, . . . . .
- ▶ Several works with astronomical imaging applications since the early 2000s

4. thanks to Marine Roux for the picture

# Outline

## Introduction and Motivations

- MUSE instrument

- Two detection problems

## Multiple inference and Global error control

- Multiple comparisons

- False Discovery Rate FDR

- BH Procedure

## Detection of galactic sources : CGM

- CGM multiple testing

- COMET procedure

- COMET Results

## Conclusion and perspectives

# Detection of galactic halo : CGM

We want to explore the gas halo surrounding a galaxy : **Circum galactic medium** or CGM).

## Galaxy properties

- ▶ Spatially limited (quasi-punctual)
- ▶ Lyman emission line + spectral continuum (+ other lines)
- ▶ Known spatial and spectral (redshift) positions

## Halo properties

- ▶ Hydrogen gas
- ▶ Emission only in Lyman line
- ▶ Spatial extension around the galaxy
- ▶ Lyman emission similar (in first approx.) to the galaxy one

🔍 **multiple testing** : need to explore a great number of pixels around the galaxy in search of the Lyman signature.

## CGM detection problem

*Goal* : Detect a quasi-connected multipixel target, while ensuring global control of errors

On each pixel  $i$ , detection of a positive signal using a one-sided test :

$$\begin{cases} H_0 : \mathbf{y}_i = \boldsymbol{\epsilon}, \\ H_1 : \mathbf{y}_i = \alpha_i \mathbf{d} + \boldsymbol{\epsilon}, \quad \text{with } \alpha_i > 0, \end{cases}$$

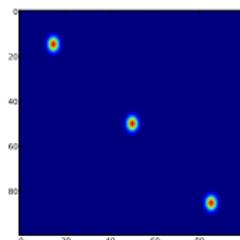
- ▶  $\boldsymbol{\epsilon} \in \mathbb{R}^l$  : noise vector of unknown distribution but assumed **symmetrical**
- ▶  $\mathbf{d} \in \mathbb{R}^l$  : known reference (Lyman signature)
- ▶  $\mathbf{y}_i \in \mathbb{R}^l$  : spectrum vector
- 📖 extension to sparse representation with multiple atoms  $\mathbf{d}_k, k = 1, \dots, K$ .

# Application of BH procedure to our case

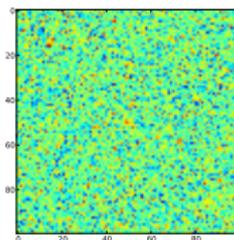
Simple but generalizable approach : matched filter test statistics  $\{w_i \equiv \mathbf{d}^T \mathbf{y}_i\}_{1 \leq i \leq n}$

- ▶ spatial, or spectral, or 3D (spatial+spectral) templates  $\mathbf{d}$

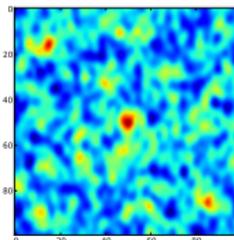
## Example : spatial matched filter



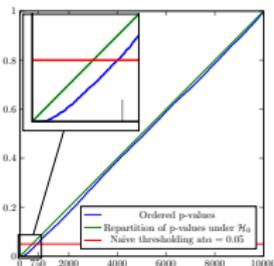
(a) Without noise



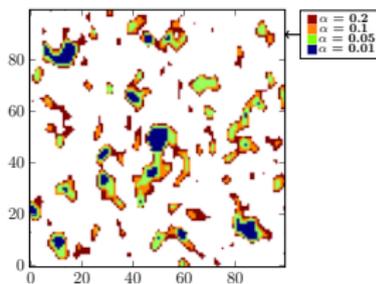
(b) Noisy image (SNR = -5dB)



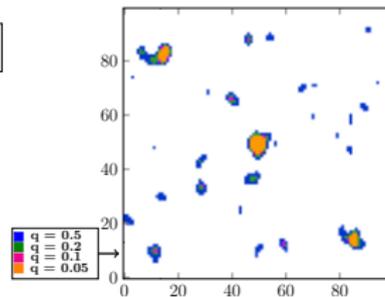
(c) Matched filter output



(d) Ordered p-values



(e) PFA thresholding



(f) FDR thresholding

# FDR control for matched filter test statistics

## Testing problem

For each pixel  $i$ ,

$$\begin{cases} H_0 : \mathbf{y}_i = \boldsymbol{\epsilon}, \\ H_1 : \mathbf{y}_i = \alpha_i \mathbf{d} + \boldsymbol{\epsilon}, \end{cases} \quad \text{with } \alpha_i > 0,$$

## Exact FDR control

Assuming

- ▶ a Gaussian noise  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ , with positive (component) covariance  $\boldsymbol{\Sigma} \succeq \mathbf{0}$
- ▶ a positive template  $\mathbf{d} \succeq \mathbf{0}$ ,

Then BH procedure applied to matched filter statistics ensures a FDR control at specified level  $q$

# Issue : Misspecification of the null distribution

## Deviation from the theoretical null

BH procedure requires so little : only the choice of the test statistics and its specification when the null hypothesis is true

- ▶ theoretical null hypothesis usually derived in an idealized framework, (e.g. does not account for complex spatial/spectral correlations, spatial inhomogeneities, standardization...)
- ☞ unlikely to be correctly specified in large-scale testing !
- ☞ large-scale testing : possibility to detect and to correct possible miss-specification of the null hypothesis

## Empirical null distribution<sup>5</sup>

Estimation of the  $H_0$  distribution : based on the observations that are the most likely under theoretical  $H_0$

---

5. Efron, B., *Cambridge University Press*, 2010

# Learning the null distribution

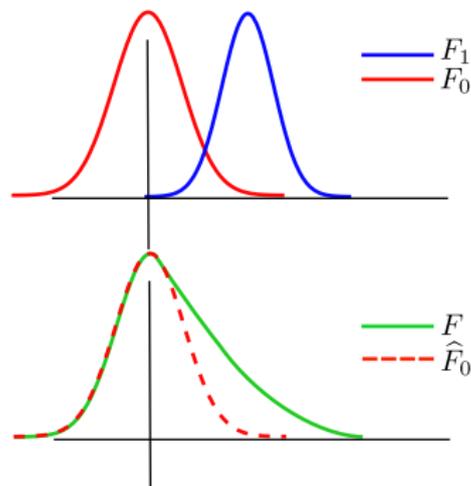
## Key assumptions

- ▶ noise distribution is symmetrical ;
- ▶ source contribution is positive.

## Empirical $p$ -values

Big picture : to have  $\hat{F}_0$ , the empirical distribution law of the  $w_i$  under  $\mathcal{H}_0$ , it is sufficient to symmetrize the negative part of the empirical distribution of the data

- ▶  $p$ -value associated to the pixel  $i$  :  $p_i = 1 - \hat{F}_0(w_i)$ .
- ☞ We can then apply the BH procedure to the empirical  $p$ -values  $\rightarrow$  empirical BH (EBH)<sup>6</sup>



6. Bacher, R. *et al.* in IEEE TSP (2017)

## Barber and Candès procedure (BC)

*BH procedure : the most well known but not always the most relevant.*

### A recent alternative : the BC procedure<sup>7</sup>

Build control statistics  $w_i$  that are

- ▶ symmetrical under  $\mathcal{H}_0$ , i.e.  $\mathbb{P}(w_i > t | i \in \mathcal{H}_0) = \mathbb{P}(w_i < -t | i \in \mathcal{H}_0)$ ,
- ▶ stochastically greater under  $\mathcal{H}_1$ , i.e.  $\mathbb{P}(w_i > t | i \in \mathcal{H}_1) > \mathbb{P}(w_i > t | i \in \mathcal{H}_0)$ .

We sort  $w_i$  by absolute decreasing order :  $|w_{(1)}| \geq |w_{(i)}| \geq |w_{(n)}|$ .

We control at level  $q$  by rejecting  $|w_{(1)}|, \dots, |w_{(\hat{k})}|$  where :

$$\hat{k} = \max \left\{ k : \frac{1 + \#\{w_{(i)}, i \leq k < 0\}}{1 \vee \#\{w_{(i)}, i \leq k > 0\}} < q \right\}$$

---

7. Barber and Candès, *Ann. Stat.* (2015)

# Control statistics

## Knockoff issue

In (BC 2015) construction of the control statistics using knockoffs

- ☞ low power and high computational cost in high dimension.

Here we already have the following hypothesis :

- ▶ noise distribution is symmetrical.
- ▶ sources have a positive contribution
- ☞ Easy build of the control statistics  $\{w_i \equiv \mathbf{d}^T \mathbf{y}_i\}_{1 \leq i \leq n}$ .

Estimate of the False Discovery Proportion FDP (among the  $w_i > 0$  discovered in a set  $\mathcal{A}$ ) :

$$\widehat{FDP} = \frac{1 + \#\{i \in \mathcal{A}, w_i < 0\}}{1 \vee \#\{i \in \mathcal{A}, w_i > 0\}}$$

With these control statistics, and BC procedures can be shown to be equivalents. How can we now gain in power ?

The BC procedure sort statistics by absolute value before looking at the signs. BUT ... we can sort the statistics in another way, for ex. here to promote connectivity.

## Algorithm

*Region growth :*

- ▶ start from an already detected region (galactic core)
  - ▶ add to the area the new pixels of interest from the neighborhood (cf next slide)
  - ▶ estimate the  $\widehat{FDP}$  on the pixels of this area
  - ▶ iterate onto the new neighborhood of the extended area
- ▶ stop onto the largest set of selection with  $\widehat{FDP}$  inferior to the given FDR
- ▶ in this set of exploration  $\mathcal{A}$  we only keep as detection pixels  $i$  with  $w_i \geq 0$

# Selection procedure

To control FDR, the selection procedure must follow [P1] :

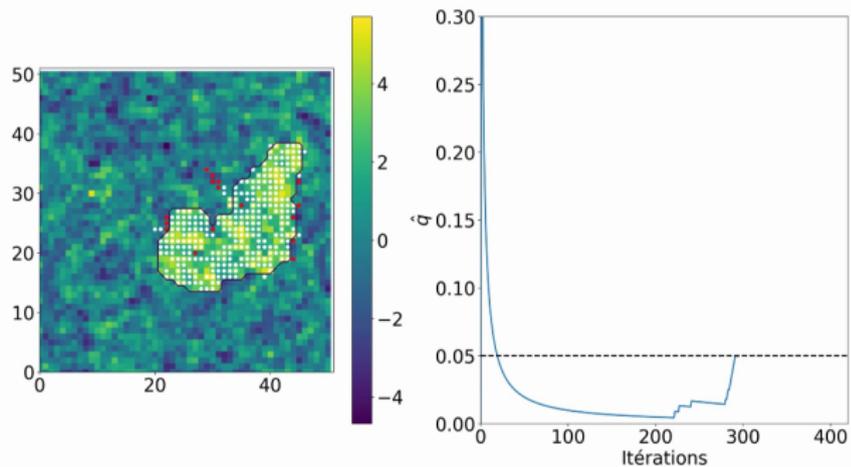
## Post-selection symmetry [P1]

For any *selected* pixel  $j$  corresponding to a true null hypothesis, the control statistic  $w_j$  is symmetrically distributed.

## Greedy approach proposed

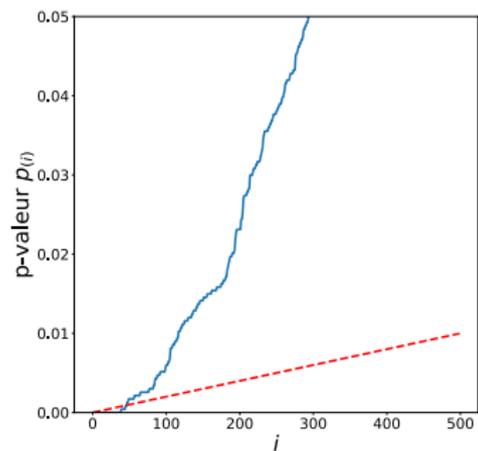
At each step, the greatest statistic (in absolute value) among the neighbors is selected.

# COMET in action

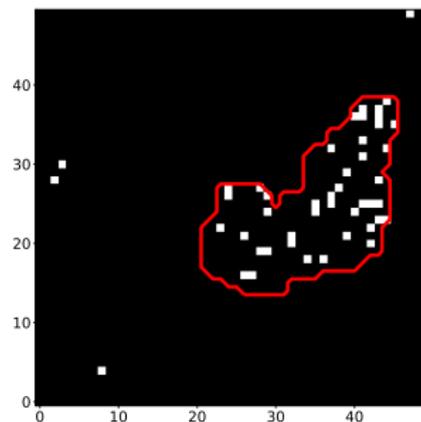


## To summarize : EBH

FDR control using EBH : thresholding of  $p$ -values pondered by the number of tests



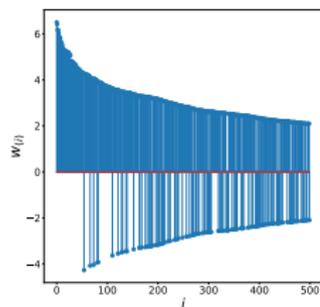
Sorted control statistics



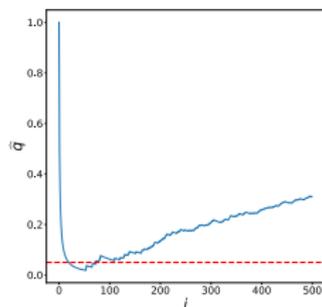
Evolution  $\widehat{FDP}$

## To summarize : BC

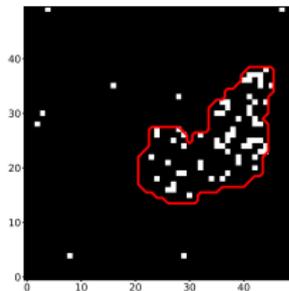
FDR control using BC : sort statistics control by absolute value



Sorted control statistics



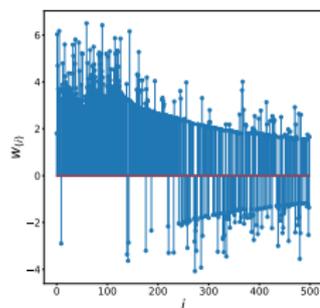
Evolution  $\widehat{FDP}$



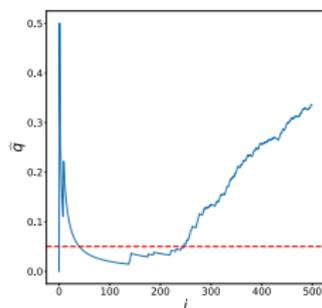
Detection map

# To summarize : COMET

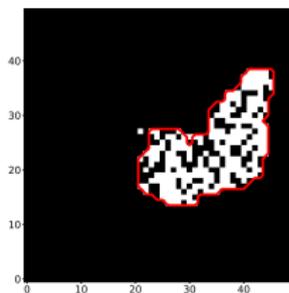
FDR control by COMET : sort statistics control using region growth



Sorted control statistics



Evolution  $\widehat{FDP}$



Detection map

## Exact control of FDR if independence of the noise

Let the noise vectors  $\epsilon_1, \dots, \epsilon_n$  be symmetrically distributed and independent. Then the COMET procedure ensures an exact control of FDR :

$$\mathbb{E} \left[ \frac{U}{R \vee 1} \right] \leq q$$

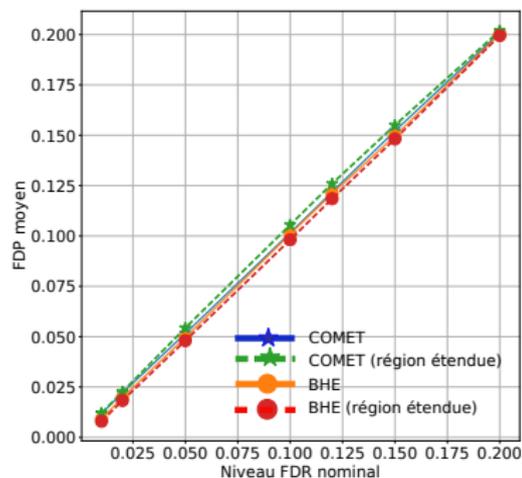
## Asymptotic control if correlated noise

Under the assumption of weakly dependent noise, if the control statistics are symmetrically distributed under  $\mathcal{H}_0$ , COMET ensures an asymptotic control of FDR

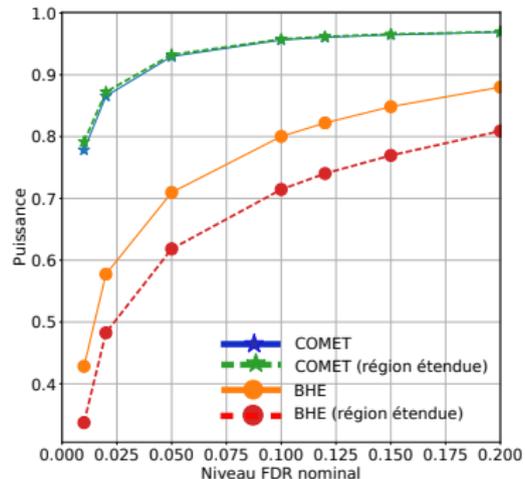
*Rk 1 : MUSE data follow this second case (short-range correlations).*

*Rk 2 : Note that these results do not require any stationarity of the noise.*

## Results (simulation)



Comparison of COMET and EBH control function of nominal FDR



Comparison of COMET and EBH power function of nominal FDR .

- ▶ Same error control
- ▶ Increase in detection power
- ▶ Power independent of total number of tests (i.e. size of explored region)

# Application to real data

## Preprocessing

- ▶ continuum subtraction
- ▶ SNR enhancement

## Spectral variability

Detection over a dictionary of spectral signatures

- ▶  $\mathbf{d}_0$  : spatial mean of galaxy core spectra.
- ▶  $\mathbf{D}$  : dictionary of shifts of  $\mathbf{d}_0$ .
- ▶ Spectral Angular Distance (SAD) .

$$SAD(\mathbf{d}_0, \mathbf{x}) = \frac{\langle \mathbf{d}_0, \mathbf{x} \rangle}{\|\mathbf{d}_0\| \|\mathbf{x}\|}$$

- ▶ Test statistics : max over all atoms  $\mathbf{w}_i = \max_k SAD(\mathbf{d}_k, \mathbf{x}_i)$

# Demo time !

Choix objet (ID UDF-10)

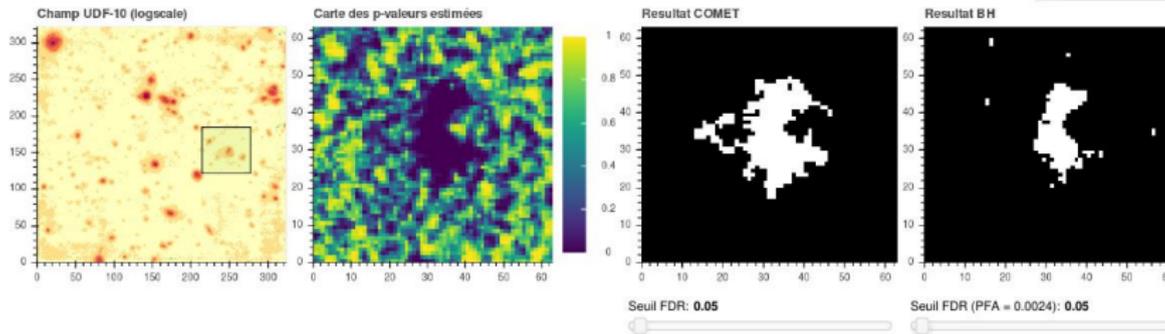
00118

Nombres d'atomes (translatées): 3

Rayon spatial: 31

Rayon spectral: 12

Go



→ <https://phd.rbacher.fr/these-app/realDetect>

# Outline

## Introduction and Motivations

- MUSE instrument

- Two detection problems

## Multiple inference and Global error control

- Multiple comparisons

- False Discovery Rate FDR

- BH Procedure

## Detection of galactic sources : CGM

- CGM multiple testing

- COMET procedure

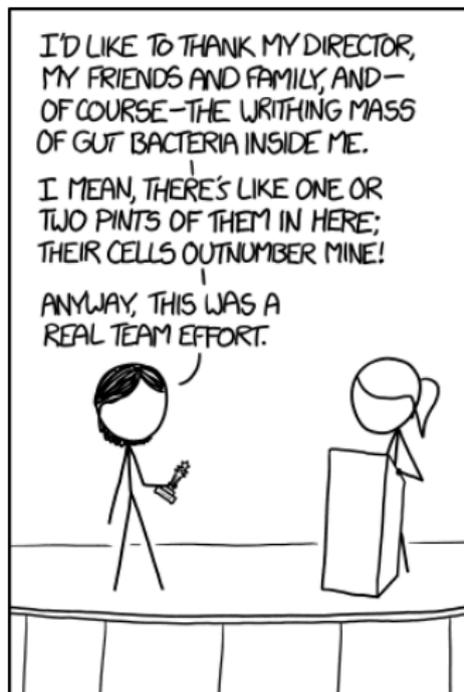
- COMET Results

## Conclusion and perspectives

# Conclusion

- ▶ Empirical approaches : no need to specify the law under  $H_0$  ( $\sim$  non parametric learning)
- ▶ Robust control of errors using this learning
- ▶ Simple hypotheses : noise symmetry and positivity of the source
- ▶ Take into account a spatial connectivity prior
- ▶ Generic detection method under FDR control with constraints
- ▶ Meaningful notion of “connected FDR” (“purity” of the detection)

Thank you !



## References

- ▶ Benjamini, Y. and Hochberg, Y. (1995), "Controlling the false discovery rate : a practical and powerful approach to multiple testing," *Journal of the Royal Statistical Society, Series B (Methodological)*, 289-300
- ▶ Barber, R. F. and Candès, E. (2015). "Controlling the False Discovery Rate via Knockoffs," *Ann. Statist.* 43 (2015), no. 5, 2055–2085.
- ▶ Efron, B. (2010), "Large-scale inference : empirical Bayes methods for estimation, testing, and prediction," (Vol. 1), *Cambridge University Press*
- ▶ Meillier, C. *et al.*, "Nonparametric Bayesian extraction of object configurations in massive data", in *IEEE TSP*, 2015, vol. 63(8).
- ▶ Meillier, C. *et al.*, "SELFIE : an object-based, Bayesian method for faint emission line source detection in MUSE deep field data cubes", in *A&A* 588, A140 (2016)
- ▶ Meillier, C. *et al.*, "Error control for the detection of rare and weak signatures in massive data", in *Proc. of EUSIPCO*, 2015, Nice, France, pp. 1974-1978
- ▶ Bacher, R. *et al.*, "Robust Control of Varying Weak Hyperspectral Target Detection With Sparse Nonnegative Representation," in *IEEE TSP*, 65 (13), pp. 3538-3550 (2017)
- ▶ Bacher, R. *et al.*, "Global error control procedure for spatially structured targets," in *Proc. of EUSIPCO*, 2017, Kos, Greece, pp. 206-210