The background features several large, overlapping, colorful swirls in shades of purple, green, and blue. Interspersed among these swirls are numerous small, yellow, triangular shapes that resemble sun rays or starbursts. The overall aesthetic is bright and dynamic.

Multiple hypothesis Testing with applications to genomics

B. Tallur

**Séminaire Symbiose
26 septembre 2008**



Plan de l'exposé

- Test statistique simple
- Hypothèses nulle et alternative, région critique, erreurs de type 1 et type 2, p-value, puissance de test
- Qu'est-ce que les tests multiples?
- Problèmes posés en génomique en termes de tests multiples

Un test statistique classique

- Hypothèse nulle: hypothèse qui sera rejetée ou non en appliquant une règle de décision
- Hypothèse alternative: c'est elle qu'on accepte en cas de rejet de l'hypothèse nulle
- Région critique: définit la règle de décision en fonction de la valeur observée d'une statistique
- Risque d'erreur de type 1 α : probabilité de rejeter une hypothèse nulle qui, en réalité, est vraie
- Risque d'erreur de type 2 β : probabilité de ne pas rejeter une hypothèse nulle qui, en réalité, est fausse
- Puissance du test $(1 - \beta)$: probabilité de rejeter une hypothèse nulle qui, en réalité, est fausse

Un Exemple: Lancer de pièce

On veut tester si une pièce de monnaie est
« équilibrée »

Hypothèse nulle à tester: la probabilité d'obtenir
« pile » = 0,5

Méthode classique consiste à fixer un seuil pour α
et définir la région critique en fonction de ce
risque

Si on fixe $\alpha=0,05$ on rejette l'hypothèse
injustement avec une probabilité de 5%

En répétant le même test 100 fois, on sait que
théoriquement on va conclure 5 fois que la pièce
n'est pas équilibrée sachant qu'en vérité elle l'est



Lancer de pièce (suite)

- On lance la pièce 10 fois et on décide de conclure que la pièce est truquée si on obtient 9 fois ou plus de pile.
- on a $\text{Prob}[9 \text{ pile ou plus}] = 0,0107$
(c'est la p-value associée à ce test)
- Avec un seuil de 0,05 on rejette notre hypothèse nulle



On lance 100 pièces

- Pour une pièce donnée, le risque d'erreur de type 1 est bien 0,0107
- Mais la proba que une des 100 pièces (peu importe laquelle) soit rejetée sera plus élevée
- Supposant toutes les pièces équilibrées, $\Pr[\text{toutes les pièces soient acceptées}] = (1 - 0,0107)^{100} = 0,34$
- Dans ce cas, la proba de conclure injustement qu'une pièce est truquée est donc plus élevée que dans le cas d'un test classique
- $\text{Prob}[\text{rejeter au moins une pièce injustement}] = 1 - 0,34 = 0,66$



Genomic Context

- Microarray experiments
- Expression levels of thousands of genes are measured simultaneously
- An important question: identification of differentially expressed genes, i.e., genes whose expression levels are associated with a response or a covariate of interest



Covariates and responses

Covariates may be:

- polytomous (nominal) e.g., cell type, treatment/control status, drug type
- Continuous e.g., dose of a drug, time

Responses could be, for example, censored survival times or other clinical outcomes

Differential expression and hypothesis testing

- The biological question of differential expression can be restated statistically as a problem of *multiple hypothesis testing*: the simultaneous test for each gene, of the null *hypothesis of no association* between the expression levels and the responses or covariates



Two types of errors

- *False positive*, or Type I error: declaring that a gene is differentially expressed when it is not
- *False negative*, or Type II error: test fails to identify a truly differentially expressed gene



Problems

- When many hypotheses are tested and each test has a specified Type I error probability, the chances of committing some Type I errors increases with the number of hypotheses



Requirements

- Need to define an appropriate Type I error rate and devising powerful multiple testing procedures that control this error rate and account for the joint distribution of the test statistics



Microarray Experiments

- m genes (variables or features)
- n samples (observations)
- For each sample a response or covariate is recorded
- The data for sample i consist of a response y_i and a gene expression profile

$$X_i = (x_{1i}, \dots, x_{mi})$$

The data

- The expression data are stored in a $m \times n$ matrix $\mathbf{X} = (x_{ji})$ with rows corresponding to genes and columns to samples

The pairs (X_i, y_i) are considered as a random sample from a population of interest

- The goal is to infer about the population, specifically, to test hypotheses concerning the joint distribution of $\mathbf{X} = (x_1, \dots, x_m)$ and Y



Differential expression

- *Simultaneous* test for each gene j
- *Null hypothesis*: H_j no association between X_j and Y (in some cases, more specific null hypothesis may be of equal mean expression levels in two populations of cells as opposed to identical distributions)



A standard approach

Consists of two steps

- Computing a test statistic T_j for each gene j
- Applying a *multiple testing procedure* to determine which hypotheses to reject while *controlling a suitably defined Type I error*
- H_j is rejected based on a statistic T_j

Set -up

- m null hypotheses $H_j, j=1, \dots, m$
- Let R be the number of rejected hypotheses

| | Not rejected | Rejected | |
|--------------------------|--------------|----------|-------|
| True null hypotheses | U | V | m_0 |
| Non-true null hypotheses | T | S | m_1 |
| Total | $m-R$ | R | m |



Set-ut (2)

- R is an observable random variable
- S , T , U and V are unobservable r.v.'s
- Ideally, one would like to minimize V (the number of false positives or type I errors) and T (the number of false negatives or type II errors)

A univariate case

- While testing a single hypothesis H_j , the procedure is: Specify level α for the Type I error rate and seek tests which minimize Type II error rate i.e. maximize power, within the class of tests with Type I error rate $\leq \alpha$.

$$\Pr\left(|T_j| \geq c_\alpha / H_j\right) \leq \alpha$$

Where c_α is the critical value

Generalization of Type I error rate to multiple testing situation

- PCER (Per-Comparison Error Rate):

$$\text{PCER} = E(V)/m$$

- PFER (Per-Family Error Rate):

$$\text{PFER} = E(V)$$

- FWER (Family-wise Error Rate):

$$\text{FWER} = \Pr(V \geq 1)$$

- FDR (False Discovery Rate):

$$\text{FDR} = E(Q) \text{ with } Q = V/R \text{ if } R > 0 \text{ and } Q = 0 \text{ if } R = 0.$$

i.e. FDR is the expected proportion of type I errors among rejected hypotheses



Type I Error rates (2)

- These error rates are defined under the true and unknown data generating distribution of $\mathbf{X} = (x_1, \dots, x_m)$ and Y
- They depend upon which specific subset

$$\Lambda_0 \subseteq \{1, 2, \dots, m\}$$

Of null hypotheses is true for this distribution

Strong and weak control of Type I error rates

- *Strong control* refers to the control of Type I error rate for *any subset* $\Lambda_0 \subseteq \{1, 2, \dots, m\}$ of true null hypotheses
- *Weak control* refers to the control of Type I error rate only *when all the null hypotheses are true* i.e. under the complete null hypothesis

$$H_0^C = \bigcap_{j=1}^m H_j$$

- In general the weak control is unsatisfactory because the complete null hypothesis is not realistic.



Comparison of Type I error rates

- In general, for a given multiple testing procedure,

$$\text{PCER} \leq \text{FWER} \leq \text{PFER}$$



Power of the test

- Within the class of multiple testing procedures that control a given Type I error rate at a level α , one seeks procedures that maximize power i.e. minimize Type II error rate

Generalization of power to multiple testing

The three common definitions are:

1. $\Pr(S \geq 1) = \Pr(T \leq m_1 - 1)$ i.e. probability of rejecting at least one false hypothesis;
2. $E(S)/m_1 =$ average power or average probability of rejecting false null hypotheses;
3. Probability of rejecting all false null hypotheses, $\Pr(S = m_1) = \Pr(T = 0)$

When the family of tests consists of pairwise mean comparisons, these quantities are called *any-pair* power, *per-pair* power and *all-pairs* power respectively



p-values

- consider a single hypothesis, say H_1 , and a test statistic T_1 . The Critical region S_α of size α is such that
- $\Pr(T_1 \in S_\alpha \mid H_1) = \alpha$
- the p-value p_1 associated with the test of H_1 is the level at which the hypothesis would be just rejected, given t_1
- $p_1 = \inf\{\alpha: t_1 \in S_\alpha\} = \Pr\{|T_1| \geq |t_1| \mid H_1\}$

Adjusted p-values (multiple testing situation)

- Let t_j and $p_j = \Pr\{|T_j| \geq |t_j| \mid H_j\}$ denote respectively, the test statistic and the unadjusted p-value for hypothesis H_j (gene j), $j=1,2,\dots,m$.
- Given any test procedure, the adjusted p-value corresponding to the single hypothesis H_j can be defined as the nominal level of the entire test procedure at which H_j would just be rejected, given the values of all test statistics

Types of multiple testing procedures

- Single-step procedures: each hypothesis is evaluated using a critical value that is independent of the results of the tests of other hypotheses

Bonferroni procedure rejects any hypothesis with unadjusted p-value $\leq \alpha/m$

The corresponding adjusted p-values are given by

$$\tilde{p}_j = \min(mp_j, 1)$$

(strong control of FWER)

Single-step procedures (2)

- Šidák procedure: it is exact for protecting the FWER under complete null hypothesis, when the unadjusted p-values are independently distributed as $U[0,1]$
- Adjusted p-values are given by

$$\tilde{p}_j = 1 - (1 - p_j)^m$$

- Single-step procedures are simple to implement but they tend to be conservative

Step-down procedures

- Let $p_{r_1} \leq p_{r_2} \leq \dots \leq p_{r_m}$ denote the observed ordered unadjusted p-values and let $H_{r_1} \leq H_{r_2} \leq \dots \leq H_{r_m}$ denote the corresponding null hypotheses.
- *The Holm's procedure* (strong control of FWER at level α):
- Define $j^* = \min \{j : p_{r_j} > \alpha / (m - j + 1)\}$
- And reject hypotheses H_{r_j} , for $j = 1, 2, \dots, j^* - 1$
- If no such j^* exists, reject all hypotheses



Step-down procedure

- Adjusted p-value for Holm's procedure is

$$\tilde{p}_{r_j} = \max_{k=1,2,\dots,j} \left\{ \min\left((m - k + 1)p_{r_k}, 1\right) \right\}$$

- Holm's procedure is less conservative than the standard Bonferroni procedure
- Other step-down procedures
 - Step-down Šidák , minP, maxT



Step-up procedures

- Hypotheses that correspond to the least significant test statistics are considered successively, with further tests dependent on the outcomes of the earlier ones. As soon as one hypothesis is rejected, all remaining hypotheses are rejected

Step-up procedures (2)

- Hochberg adjusted p-values are:

$$\tilde{p}_{r_j} = \min_{k=j, \dots, m} \left\{ \min\left((m - K + 1)p_{r_k}, 1\right) \right\}$$

- There are other step-up procedures such as
 - Benjamini and Hochberg (1995)
 - Benjamini and Yekutieli (2001)

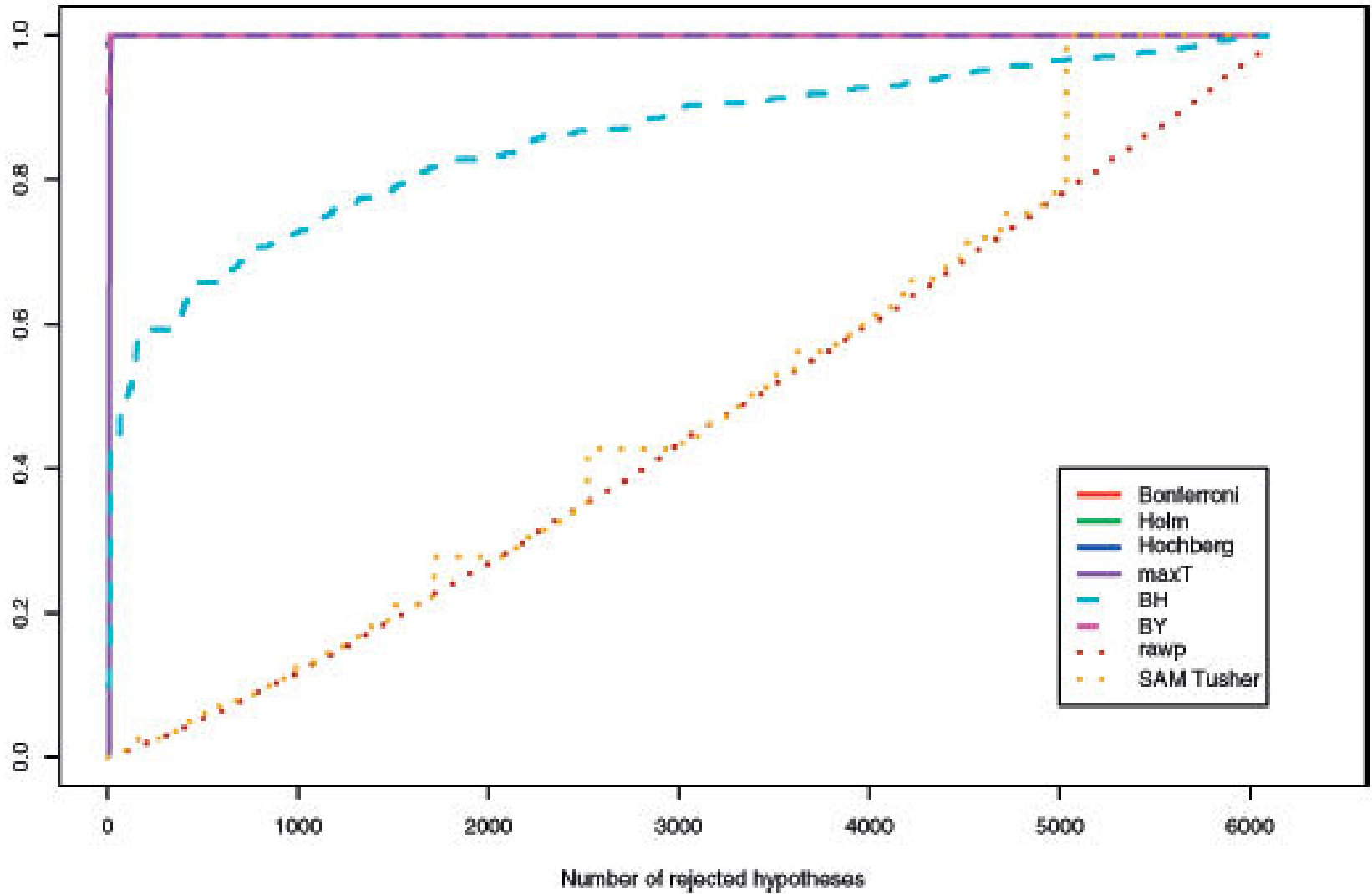
Estimation of adjusted p-values

- Permutation methods are proposed to estimate the marginal or joint distribution of the test statistics.
- (I will not talk about these algorithms today, may be next time...)



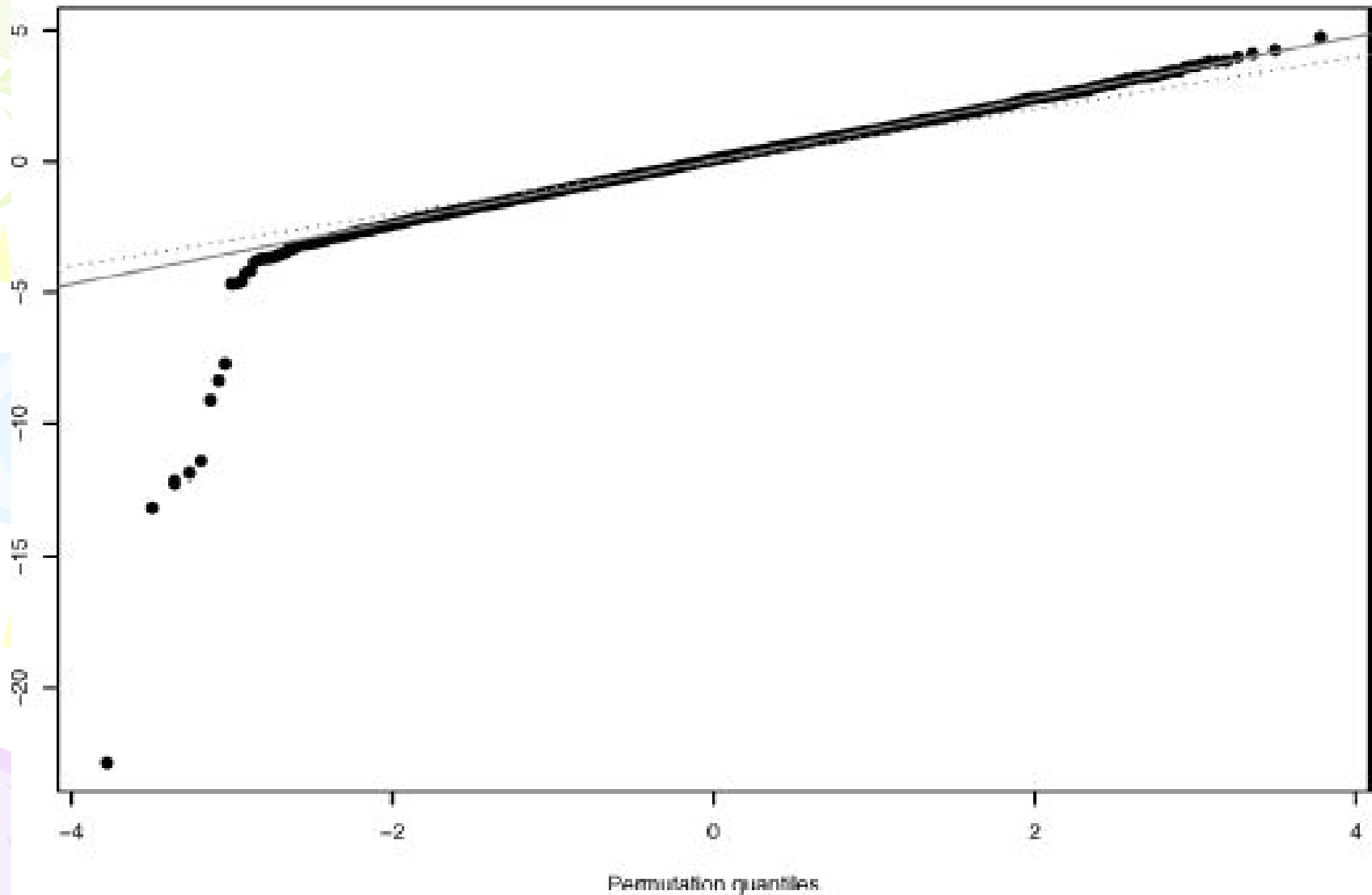
Comparison on real and simulated data

- Dudoit et al. have compared different multiple testing procedures on several real and simulated data
- Apolipoprotein AI experiment of Callow et al. (2000) as part of the study of lipid metabolism in mice



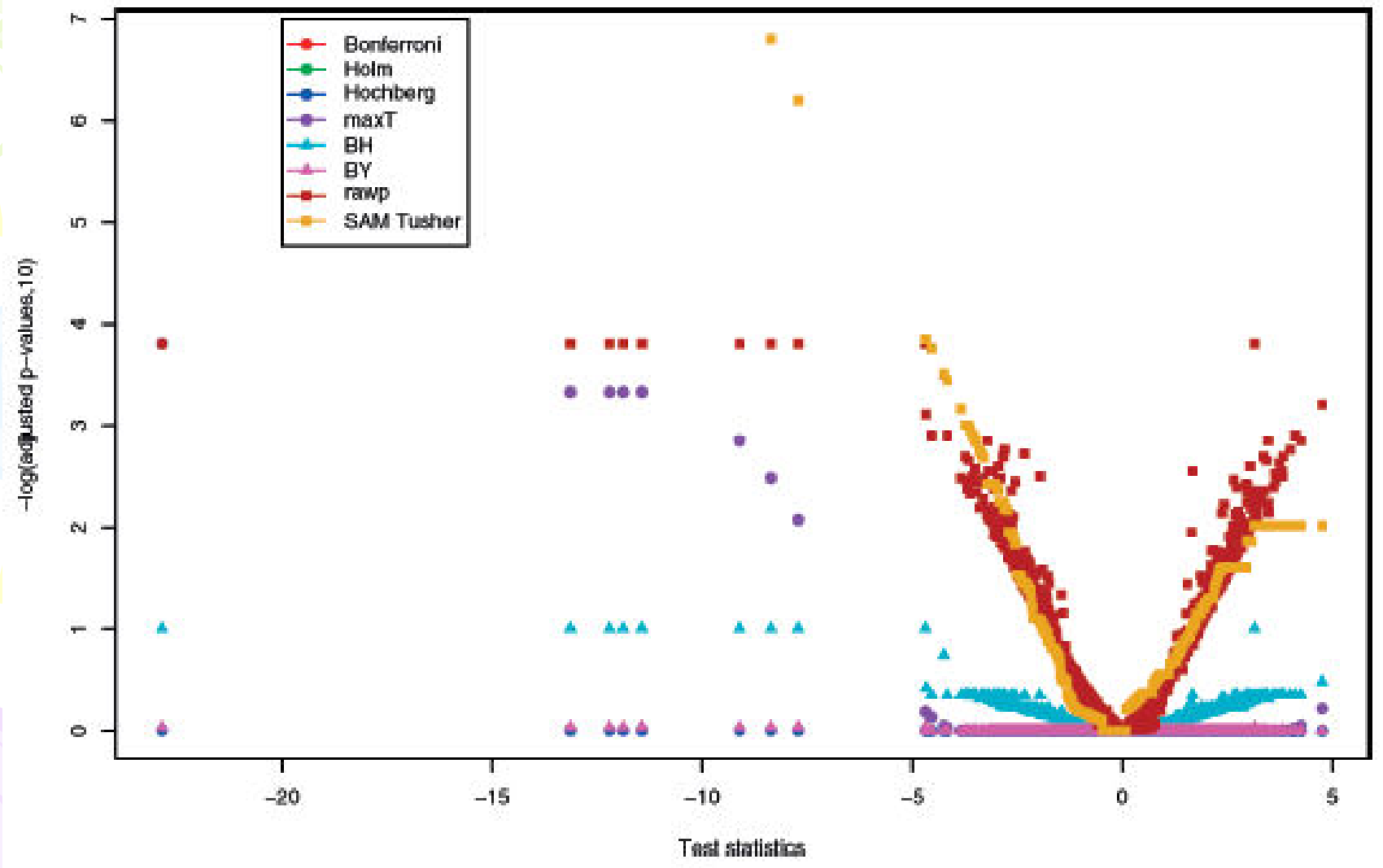
(a)

Quantile-Quantile plot

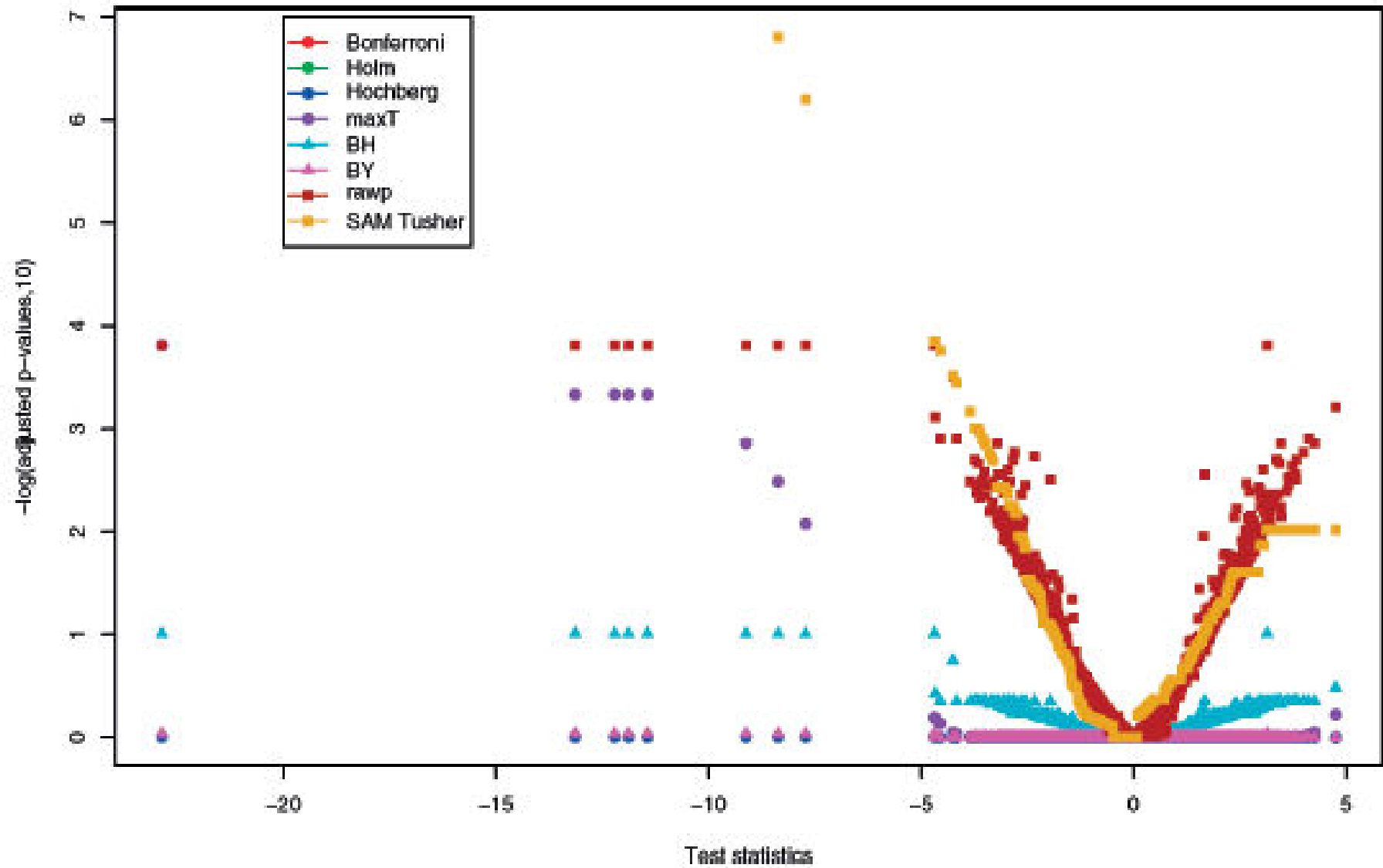


(d)

(a)



(c)



(c)



References

- Sandrine Dudoit, Juliet Popper Shaffer and Jennifer C Boldrik; « Multiple hypothesis testing in microarray experiments »; *Statistical Science*, 2003, vol. 18, No 1, 71 – 103
- Merrill D. Birkner, Katherine S. Pollard, Mark J. van der Laan, Sandrine Dudoit; « Multiple testing procedures and applications to genomics »; 2005, paper 168, university of California, Berkley Division of Biostatistics working papers series
- Shaffer, J.P.; « Multiple hypothesis testing: A review »; 1986, *Annual review of psychology*, 46, 561 – 584.



References

- Yoav. Benjamini, Yosef Hochberg
- Controlling false discovery rate: A practical and powerful approach to multiple testign, JRSSS (1995)
- J.D Storey, R. Tibshirani
Statistical significance for genomewide studies, PNAS, Aug 5, 2003