

Bayesian Maximum *a posteriori* Multiple Testing Procedure

Felix Abramovich

Tel Aviv University, Tel Aviv, Israel

Claudia Angelini

Consiglio Nazionale delle Ricerche, Napoli, Italy

Abstract

We consider a Bayesian approach to multiple hypothesis testing. A hierarchical prior model is based on imposing a prior distribution $\pi(k)$ on the number of hypotheses arising from alternatives (false nulls). We then apply the maximum *a posteriori* (MAP) rule to find the most likely configuration of null and alternative hypotheses. The resulting MAP procedure and its closely related step-up and step-down versions compare ordered Bayes factors of individual hypotheses with a sequence of critical values depending on the prior. We discuss the relations between the proposed MAP procedure and the existing frequentist and Bayesian counterparts. A more detailed analysis is given for the normal data, where we show, in particular, that by choosing a specific $\pi(k)$, the MAP procedure can mimic several known familywise error (FWE) and false discovery rate (FDR) controlling procedures. The performance of MAP procedures is illustrated on a simulated example.

AMS (2000) subject classification. Primary 62F15, 62F03.

Keywords and phrases. Bayes factor, false discovery rate, familywise error, hierarchical prior, maximum *a posteriori* rule, multiple hypothesis testing, *p*-value.

1 Introduction

Consider the standard multiple hypothesis testing set-up. Suppose that we have n independent vectors of observations \mathbf{Y}_i , $i = 1, \dots, n$ of sizes m_i , where $\mathbf{Y}_i \sim f_i(\mathbf{Y}_i|\boldsymbol{\theta}_i)$, and $\boldsymbol{\theta}_i$ is a d_i -dimensional parameter vector in $\Omega_i \subset \mathbb{R}^{d_i}$. Given this data, we want to simultaneously test n nonnested hypotheses

$$H_{0i} : \boldsymbol{\theta}_i \in \Theta_i \quad \text{vs.} \quad H_{1i} : \boldsymbol{\theta}_i \in \bar{\Theta}_i, \quad (1.1)$$

where $\Theta_i \cap \bar{\Theta}_i = \emptyset$, $\Theta_i \cup \bar{\Theta}_i = \Omega_i$, $i = 1, \dots, n$. Following a frequentist approach to test H_{0i} against H_{1i} , one chooses some test statistic $T_i(\mathbf{Y}_i)$, and the inference is typically based on the corresponding p -value p_i . Since there are n such individual tests, the resulting p -values should be adjusted for multiplicity. The multiplicity correction depends on the type of the error one wishes to control. The traditional concern for multiplicity effect has been about the familywise error (FWE) – the probability of even a single Type I error in a series of n tests. The widely known Bonferroni procedure is only an example, and more powerful FWE controlling procedures are currently available for many multiple testing problems (see Hochberg and Tamhane, 1987; Hsu, 1996 for reviews). Such a severe criterion implies however substantially reduced power properties especially when n is large. A less stringent alternative to the FWE is the false discovery rate (FDR) criterion of Benjamini and Hochberg (1995). Unlike FWE, FDR controls the expected proportion of Type I errors among hypotheses being rejected (false discoveries) rather than the probability of even a single Type I error. The resulting FWE or FDR controlling procedures are typically stepwise in nature where the ordered p -values $p_{(1)} \leq \dots \leq p_{(n)}$ are in effect compared with a series of properly chosen critical values. Step-up procedures start with testing the least significant hypothesis with the largest p -value $p_{(n)}$ and continue with decreasing p -values until the first rejection of the null hypothesis. Step-down procedures start with $p_{(1)}$ and continue with increasing p -values until the first acceptance.

A Bayesian approach to hypothesis testing in general and to the multiplicity problem in particular are conceptually different. Sarkar and Chen (2004) give an overview of the current Bayesian perspective on multiple testing. Consider first the i -th test as a single test. One assumes a prior distribution on $\boldsymbol{\theta}_i$ of the form

$$\pi_i(\boldsymbol{\theta}_i) = \begin{cases} \pi_{0i}p_{0i}(\boldsymbol{\theta}_i), & \text{if } \boldsymbol{\theta}_i \in \Theta_i \\ \pi_{1i}p_{1i}(\boldsymbol{\theta}_i), & \text{if } \boldsymbol{\theta}_i \in \bar{\Theta}_i, \end{cases}$$

where $p_{0i}(\boldsymbol{\theta}_i)$ and $p_{1i}(\boldsymbol{\theta}_i)$ are densities on Θ_i and $\bar{\Theta}_i$ respectively, and $\pi_{0i} + \pi_{1i} = 1$. The inference is then based on the posterior distribution $\pi_i(\boldsymbol{\theta}_i|\mathbf{Y}_i)$ according to the chosen loss. For the standard “0-1” loss, the null hypothesis H_{0i} is rejected if the resulting posterior odds ratio $\pi_i(\boldsymbol{\theta}_i \in \Theta_i|\mathbf{Y}_i)/\pi_i(\boldsymbol{\theta}_i \in \bar{\Theta}_i|\mathbf{Y}_i)$ is smaller than 1 or, equivalently, if the corresponding Bayes factor B_i is smaller than π_{1i}/π_{0i} . This Bayesian rule can be easily extended to a more general “0- L_k ” loss under which the null hypothesis is rejected if $B_i < (L_{0i}/L_{1i})(\pi_{1i}/\pi_{0i})$. Turning back to multiple hypothesis testing, note

that from the Bayesian perspective, placing independent priors on θ_i and using additive losses (e.g., a sum of “0-1” individual losses for each test) does not imply any multiplicity correction: the evident optimal Bayes rule in this case simply applies the corresponding individual Bayes testing rule to each one of n tests “ignoring”, therefore, the multiplicity effect. To account for multiplicity adjustments within a Bayesian framework, one should consider θ_i to be *a priori* dependent to cause the posterior distribution of θ_i to depend on all $\mathbf{Y}_1, \dots, \mathbf{Y}_n$. One possible way to introduce dependency among θ_i is via *hierarchical* prior models. In situations where all likelihoods $f_i(\mathbf{Y}_i|\theta_i)$ have the same parametric form, like in multiple comparisons, $\theta_1, \dots, \theta_n$ are usually considered as an independent sample from a population distribution $\mathcal{F}(\theta)$ that may be parametric with possibly unknown parameters or even nonparametric (e.g., Waller and Duncan, 1969; Gopalan and Berry, 1998; Berry and Hochberg, 1999; Sarkar and Chen, 2004; Scott and Berger, 2006).

Generally, however, the likelihoods $f_i(\mathbf{Y}_i|\theta_i)$ may be of different forms, θ_i may be of different dimensionalities, etc. and such type of hierarchical priors might be inappropriate. In this paper we propose a simple hierarchical prior model for a multiple hypothesis testing set-up (1.1) by imposing a prior distribution on the *number* of hypotheses arising from the alternatives (false nulls). We then apply the maximum *a posteriori* (MAP) rule and find the most plausible configuration of nulls and alternatives with the maximal posterior probability. Such a Bayes rule essentially corresponds to the following (non-additive) “nothing or everything” multiple “0-1” loss : the loss is zero if all n hypotheses are inferred upon correctly and one if there is at least one wrong decision of any type regardless of their actual number. In the case of independent likelihoods, the resulting MAP procedure is based on the sequence of ordered individual Bayes factors somewhat similar in spirit to frequentist procedures operating with ordered p -values. We also present closely related step-up and step-down versions of the MAP multiple testing procedure and establish interesting parallels between them and several known frequentist and Bayesian multiple testing procedures.

The paper is organized as follows. The main results are presented in Section 2 where we propose the Bayesian MAP testing procedure and its stepwise versions. The relations between them and several existing Bayesian and frequentist multiple testing procedures are discussed in Section 3 while a more detailed study for the normal data is given in Section 4. In Section 5, we illustrate the performance of MAP procedures on a simulated example. Some concluding remarks are made in Section 6.

2 Main Results

2.1. Hierarchical prior model. Consider again the data set of n independent m_i -dimensional vectors \mathbf{Y}_i , where $\mathbf{Y}_i \sim f_i(\mathbf{Y}_i|\boldsymbol{\theta}_i)$, $\boldsymbol{\theta}_i \in \Omega_i \subset \mathbb{R}^{d_i}$ and the multiple hypothesis testing problem (1.1). Following a Bayesian methodology, one needs to set prior odds simultaneously for each of n null hypotheses. However, it might not be obvious to formulate prior beliefs on individual odds especially when they are not assumed to be independent. On the other hand, one has usually some intuition on the number (proportion) of hypotheses coming from the nulls and the alternatives (respectively, true and false nulls). For example, in the analysis of microarray data, it might be hard to estimate in advance the chances for each gene to be differentially expressed but typically it is believed that the number of such genes is small. A configuration of true and false null hypotheses is uniquely determined by an n -dimensional indicator vector \mathbf{x} , where $x_i = I\{\boldsymbol{\theta}_i \in \bar{\Theta}_i\}$, $i = 1, \dots, n$. Let $k = x_1 + \dots + x_n$ be the number of hypotheses coming from the alternatives and assume some prior distribution $k \sim \pi(k) > 0$, $k = 0, \dots, n$. Several possible choices for $\pi(k)$ are discussed later in the paper. For a given k , assume all $\binom{n}{k}$ various configurations of true and false null hypotheses to be equally likely *a priori*, i.e., conditionally on k ,

$$P\left(\mathbf{x} \mid \sum_{i=1}^n x_i = k\right) = \binom{n}{k}^{-1}. \quad (2.1)$$

To complete the prior, assume

$$(\boldsymbol{\theta}_i|x_i = 0) \sim p_{0i}(\boldsymbol{\theta}_i) \quad \text{and} \quad (\boldsymbol{\theta}_i|x_i = 1) \sim p_{1i}(\boldsymbol{\theta}_i) \quad (2.2)$$

for some specified densities $p_{0i}(\boldsymbol{\theta}_i)$ and $p_{1i}(\boldsymbol{\theta}_i)$ on Θ_i and $\bar{\Theta}_i$ respectively.

2.2. MAP multiple testing procedure. The Bayesian inference in multiple hypothesis testing is based on the posterior joint distribution of null hypotheses that is uniquely defined by the posterior distribution $\pi(\mathbf{x}, k|\mathbf{Y}_1, \dots, \mathbf{Y}_n)$.

For the independent data $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ and the proposed hierarchical prior model, we have

$$\begin{aligned}
& \pi(\mathbf{x}, k | \mathbf{Y}_1, \dots, \mathbf{Y}_n) \\
& \propto \binom{n}{k}^{-1} \pi(k) I \left\{ \sum_{i=1}^n x_i = k \right\} \\
& \times \prod_{i=1}^n \left(\int_{\bar{\Theta}_i} f_i(\mathbf{Y}_i | \boldsymbol{\theta}_i) p_{1i}(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i \right)^{x_i} \left(\int_{\Theta_i} f_i(\mathbf{Y}_i | \boldsymbol{\theta}_i) p_{0i}(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i \right)^{1-x_i} \\
& \propto \binom{n}{k}^{-1} \pi(k) I \left\{ \sum_{i=1}^n x_i = k \right\} \prod_{i=1}^n \left(\frac{\int_{\bar{\Theta}_i} f_i(\mathbf{Y}_i | \boldsymbol{\theta}_i) p_{1i}(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i}{\int_{\Theta_i} f_i(\mathbf{Y}_i | \boldsymbol{\theta}_i) p_{0i}(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i} \right)^{x_i} \\
& = \binom{n}{k}^{-1} \pi(k) I \left\{ \sum_{i=1}^n x_i = k \right\} \prod_{i=1}^n (B_i^{-1})^{x_i}, \tag{2.3}
\end{aligned}$$

where B_i is the Bayes factor of H_{0i} (e.g., Berger, 1985, Section 4.3.3):

$$B_i = \frac{\int_{\Theta_i} f_i(\mathbf{Y}_i | \boldsymbol{\theta}_i) p_{0i}(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i}{\int_{\bar{\Theta}_i} f_i(\mathbf{Y}_i | \boldsymbol{\theta}_i) p_{1i}(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i}. \tag{2.4}$$

Given the posterior distribution $\pi(\mathbf{x}, k | \mathbf{Y}_1, \dots, \mathbf{Y}_n)$, a common approach is to select the most likely configuration of true and false null hypotheses, that is the posterior mode of (2.3). Such a Bayesian rule is also known as a maximum *a posteriori* (MAP) rule and corresponds to the (non-additive) “nothing or everything” multiple “0-1” loss discussed in the Introduction. Generally, to find the posterior mode of $\pi(\mathbf{x}, k | \mathbf{Y}_1, \dots, \mathbf{Y}_n)$, one should look through all 2^n various configurations of true and false null hypotheses. However, for the model at hand, the number of possible candidates is essentially reduced to $n + 1$ only. To see this, let $\hat{\mathbf{x}}(k)$ be a maximizer of (2.3) for a *fixed* k that indicates the most likely configuration with k false null hypotheses. From (2.3), the obvious solution for $\hat{\mathbf{x}}(k)$ is $\hat{x}_i(k) = 1$ for the k tests with the smallest Bayes factors B_i and zeroes for others. Thus, one needs to choose only from $n + 1$ configurations $(\hat{\mathbf{x}}(k), k) = H_{1(1)} \cap \dots \cap H_{1(k)} \cap H_{0(k+1)} \dots \cap H_{0(n)}$, $k = 0, \dots, n$, where $H_{0(i)}$ and $H_{1(i)}$ are respectively the null and alternative hypotheses corresponding to $B_{(i)}$, $B_{(1)} \leq \dots \leq B_{(n)}$. This leads to the following Bayesian MAP multiple testing procedure :

1. Calculate Bayes factors B_i , $i = 1, \dots, n$ for all individual tests and order them in increasing sequence $B_{(1)} \leq \dots \leq B_{(n)}$.

2. Find \hat{k} that maximizes

$$\hat{\pi}_k = \pi(\hat{\boldsymbol{x}}(k), k | \mathbf{Y}_1, \dots, \mathbf{Y}_n) \propto \binom{n}{k}^{-1} \pi(k) \prod_{i=0}^k B_{(i)}^{-1}, \quad (2.5)$$

where we set $B_{(0)} = 1$ to include $k = 0$.

3. Accept all null hypotheses if $\hat{k} = 0$; otherwise, reject \hat{k} null hypotheses corresponding to $B_{(1)}, \dots, B_{(\hat{k})}$ and accept others.

One can also consider stepwise versions of the above MAP procedure. In addition to the global maximum \hat{k} in (2.5), let \hat{k}_l and \hat{k}_r be respectively the leftmost and the rightmost local maxima of $\hat{\pi}_k$. Obviously, $\hat{k}_l \leq \hat{k} \leq \hat{k}_r$. Consider the successive ratios

$$\frac{\hat{\pi}_k}{\hat{\pi}_{k-1}} = \frac{k}{n - k + 1} \frac{\pi(k)}{\pi(k-1)} B_{(k)}^{-1}.$$

Evidently, \hat{k}_r is the maximal k for which $\hat{\pi}_k / \hat{\pi}_{k-1} > 1$. Hence, choosing $k = \hat{k}_r$ corresponds to the Bayesian step-up (MAP-up) multiple testing procedure starting with the least significant null hypothesis with the largest Bayes factor $B_{(n)}$ and accepting null hypotheses as long as

$$B_{(k)} > \frac{k}{n - k + 1} \frac{\pi(k)}{\pi(k-1)}. \quad (2.6)$$

Analogously, $k = \hat{k}_l$ yields the Bayesian step-down (MAP-down) procedure that starts with $B_{(1)}$ and rejects null hypotheses until for the first time (2.6) holds.

3 Some Parallels with Existing Multiple Testing Procedures

In this Section, we discuss the relations between the presented MAP multiple testing procedures with several other existing Bayesian and frequentist procedures.

3.1. Bayesian procedures. Bayesians usually treat simultaneous testing of n hypotheses within model selection framework, where configuration of true and false null hypotheses is considered as a disjoint partition of the parameter space $\Omega_1 \oplus \dots \oplus \Omega_n$ (e.g., Kass and Raftery, 1995; Berger and

Pericchi, 1996). The main obstacle of this approach is that one then generally needs to compare 2^n possible configurations. The proposed MAP multiple testing procedure can also be viewed within such a framework but, as we have already mentioned, assuming independence of $f_i(\mathbf{Y}_i|\boldsymbol{\theta}_i)$, there are only $n + 1$ candidates to look through in our case. Similar to frequentist analysis, stepwise multiple testing procedures allow further reduction of the computational cost of the problem for large n . To the best of our knowledge, the only attempt to develop a Bayesian multiple stepwise testing procedure was in a recent paper of Sarkar and Chen (2004) although their hierarchical prior model is different. Using our notations, the motivation behind their procedure can be described as follows. Let π_k and $\hat{\pi}_k$ be a prior and the corresponding posterior probability respectively of a configuration $(\hat{\mathbf{x}}(k), k)$. In particular, for our model $\pi_k = \binom{n}{k}^{-1}\pi(k)$ and $\hat{\pi}_k$ is given in (2.5). The *stepwise* Bayes factor $B^{(k)}$ of Sarkar and Chen (2004) is the Bayes factor of $(\hat{\mathbf{x}}(k), k)$ to *any* of $(\hat{\mathbf{x}}(i), i)$, $i = k + 1, \dots, n$ and is defined as

$$\begin{aligned} B^{(k)} &= \frac{P(\hat{\mathbf{x}}(k), k | \mathbf{Y}_1, \dots, \mathbf{Y}_n)}{P(\bigcup_{i=k+1}^n \{(\hat{\mathbf{x}}(i), i) | \mathbf{Y}_1, \dots, \mathbf{Y}_n\})} \cdot \frac{P(\bigcup_{i=k+1}^n \{(\hat{\mathbf{x}}(i), i)\})}{P(\hat{\mathbf{x}}(k), k)} \\ &= \frac{\hat{\pi}_k}{\sum_{i=k+1}^n \hat{\pi}_i} \cdot \frac{\sum_{i=k+1}^n \pi_i}{\pi_k}. \end{aligned}$$

One then starts from $k = 0$ and stops the first time $B^{(k)} > 1$. The main difference between such a step-down procedure and the one proposed in Section 2.2. is that the latter in effect compares $(\hat{\mathbf{x}}(k), k)$ against $(\hat{\mathbf{x}}(k + 1), k + 1)$ rather than against $\bigcup_{i=k+1}^n (\hat{\mathbf{x}}(i), i)$. The critical values for the corresponding Bayes factors also differ.

3.2. Frequentist procedures. As we have mentioned in the Introduction, Bayesian multiple testing procedures and their frequentist counterparts are conceptually different. The latter work with a series of p -values p_i while the former are typically based on a sequence of Bayes factors B_i . O'Hagan (1995) and Berger and Pericchi (1996) proposed modifications of a standard Bayes factor for the case of improper noninformative priors but these issues are beyond the scope of the paper. The dispute about whether a Bayes factor or a p -value represents the better evidence for statistical inference and attempts to "reconcile" them have a long history (see e.g., Casella and Berger, 1987; Good, 1992; Kass and Raftery, 1995; Berger *et al.* 1997; Berger, 2003; Bayarri and Berger, 2004). The differences are generally due to different philosophies behind frequentist and Bayesian approaches. Nevertheless, we believe that they are not complete strangers and may benefit from each

other. Rubin (1984) argued for the importance of frequentist analysis of Bayesian procedures for better understanding and validating their results. Similarly, Bayesian interpretation of a frequentist procedure is often helpful in providing intuition behind it.

Generally, there are no immediate connections between B_i and p_i , and direct analogies between the two approaches to hypothesis testing can be verified in some particular cases only. Consider the following, which is probably the most known case. Let Y_i have independent symmetric location distributions $f_i(|y_i - \theta_i|)$, $i = 1, \dots, n$ that have monotone likelihood ratios. Consider a set of n one-sided simultaneous tests $H_{0i} : \theta_i \leq \theta_{0i}$ against $H_{1i} : \theta_i > \theta_{0i}$. Assume some prior $\pi(k)$ on the number of false H_{0i} , (2.1) and noninformative priors $p_{0i}(\theta_i) = 1_{(-\infty, \theta_{0i})}(\theta_i)$ and $p_{1i}(\theta_i) = 1_{(\theta_{0i}, \infty)}(\theta_i)$ in (2.2). Simple straightforward calculus shows that in this case, due to the symmetry of f_i , $P(\theta_i \leq \theta_{0i} | Y_i) = p_i$, where p_i is the p -value for the likelihood ratio test (e.g., Berger, 1985, Section 4.3.3). The resulting ordered Bayes factors (2.4) are evidently $B_{(i)} = p_{(i)} / (1 - p_{(i)})$, and the stepwise MAP multiple testing procedures introduced in Section 2.2. then compare $p_{(i)}$ with a sequence of critical values $p_i^* = c_i / (1 + c_i)$, where c_i is given in the right-hand side of (2.6). Using the relation $B_i = p_i / (1 - p_i)$ and re-writing (2.6) in terms of p_i , it is possible to find corresponding priors $\pi(k)$ to mimic various stepwise frequentist one-sided testing procedures. One is, however, naturally looking for some “meaningful” priors. As an illustration, consider two particular choices for $\pi(k)$.

Let $k \sim B(n, \alpha_n)$. This binomial prior suggests that each null hypothesis independently has the same prior probability α_n (depending possibly on n) of being false. Small α_n reflects a *sparsity* assumption that there is only a small fraction of hypotheses coming from alternatives (true alternatives or, equivalently, false nulls) whereas the majority of true hypotheses arise from nulls (true nulls). The binomial prior yields independent priors on θ_i and the corresponding $p_i^* = \alpha_n$ are the same for all $p_{(i)}$ (see (2.6)). Obviously, $\hat{k}_l = \hat{k} = \hat{k}_r$ in this case, both stepwise versions coincide with the original MAP procedure and reject all nulls with p -values less than α_n . The binomial prior, therefore, corresponds to independent testing of each individual null hypothesis at the same significance level α_n that implies the expected number of erroneously rejected null hypotheses $n_0 \alpha_n$, where n_0 is the (unknown) number of true nulls. For a *super-sparse* case $E(k) = n \alpha_n < 1$, and it is equivalent to the Bonferroni procedure with the significance level $\alpha = n \alpha_n$.

For the truncated geometric prior $G_n^*(1-q)$, where $\pi(k) = (1-q)q^k/(1-q^{n+1})$, $k = 0, \dots, n$; $0 < q < 1$, the ordered $p_{(i)}$'s are successively compared with

$$p_i^* = \frac{i}{n - i(1-q) + 1} q, \quad 1 \leq i \leq h,$$

which coincide with the critical values from the adaptive step-down procedure of Benjamini *et al.* (2006). The authors gave the motivation for such a procedure and reported on simulation results showing that its FDR level does not exceed q although no rigorous proof was given.

The priors $\pi(k)$ corresponding to both Bonferroni and adaptive FDR procedures are sparse. In fact, this is not so surprising. The “nothing or everything” multiple “0-1” loss behind the MAP testing procedure places equal penalties for a wrong decision in both directions. On the other hand, traditional frequentist criteria in multiple testing, like FWE and FDR, mainly concern erroneous rejections of null hypotheses. Thus, both approaches are likely to yield similar results when a Bayesian believes *a priori* that the number of hypotheses arising from alternatives is small, and the total loss for their wrong rejections is, therefore, relatively low. Finally, note that for both cases, the resulting FWE and FDR levels respectively for the MAP procedure are defined by the parameters of the prior $\pi(k)$ instead of being fixed in advance at some standard level (e.g., .01 or .05). In practice, these parameters are rarely known *a priori* and should be estimated from the data (see Section 5.1.), which makes the corresponding error levels *data-adaptive*.

4 Normal Data

In this Section, we discuss the proposed MAP multiple testing procedure and its stepwise versions for the normal data. Let Y_{i1}, \dots, Y_{im_i} , $i = 1, \dots, n$ be independent random samples of sizes m_i from $N(\mu_i, \sigma_i^2)$, where the variance σ_i^2 is assumed to be known. Testing one-sided hypotheses $H_{0i} : \mu_i \leq \mu_{0i}$ against $H_{1i} : \mu_i > \mu_{0i}$ is essentially a particular case of a general one-sided testing problem considered in Section 3.2, and all the results developed there can be applied directly to normal data. In this section, we focus on multiple testing of *point* null hypotheses $H_{0i} : \mu_i = \mu_{0i}$, which is probably more relevant in practice, against the two-sided alternative $H_{1i} : \mu_i \neq \mu_{0i}$. Unlike one-sided testing, frequentist and Bayesian inferences are different in this case. There are several important issues in imposing priors for point null hypotheses. To test a point null hypothesis one cannot use a continuous prior density $\pi_i(\mu_i)$ on μ_i since it yields a zero prior (and, hence, posterior)

probability for $\mu_i = \mu_{0i}$. A usual approach in this case is to set p_{0i} to be a probability atom $\delta(\mu_{0i})$ at μ_{0i} and a density $p_{1i}(\mu_i)$ for $\mu_i \neq \mu_{0i}$ in (2.2), where, unlike the case of one-sided hypotheses, $p_{1i}(\mu_i)$ should be a *proper* density to avoid Lindley's paradox (e.g., Berger, 1985, Section 4.3.3). A simple and common choice is a conjugate prior $N(\mu_{0i}, \tau^2)$ (see e.g., Berger *et al.* 1997). Then straightforward calculus implies that

$$B_i = \sqrt{1 + \gamma_i} \exp\left\{-\frac{Z_i^2}{2(1 + 1/\gamma_i)}\right\}, \tag{4.1}$$

where $Z_i = \sqrt{m_i}(\bar{Y}_i - \mu_{0i})/\sigma_i$ is the standardized sample average and $\gamma_i = m_i\tau^2/\sigma_i^2$ is the variance ratio (see also (4.15)–(4.17) of Berger, 1985). On the other hand, the corresponding p -value for the i -th test is $p_i = 2\tilde{\Phi}(|Z_i|)$, where $\tilde{\Phi}(\cdot) = 1 - \Phi(\cdot)$ and $\Phi(\cdot)$ is the standard normal cumulative distribution function. From (4.1), one has $B_i \leq \sqrt{1 + \gamma_i}$ and

$$p_i = 2\tilde{\Phi}\left(\left(2(1 + 1/\gamma_i) \ln\left\{\sqrt{1 + \gamma_i}/B_i\right\}\right)^{1/2}\right), \tag{4.2}$$

which is an increasing function of B_i . Hence, imposing a prior $\pi(k)$ on the number of false null hypotheses, (2.6) in the stepwise MAP multiple testing procedures becomes equivalent to comparing $p_{(i)}$ with the critical value p_i^* , where

$$p_i^* = 2\tilde{\Phi}\left(\left(2(1 + 1/\gamma_i) \ln_+\left\{\sqrt{1 + \gamma_i} \cdot \frac{n - i + 1}{i} \cdot \frac{\pi(i - 1)}{\pi(i)}\right\}\right)^{1/2}\right), \tag{4.3}$$

and $\ln_+(x) = \max(0, \ln(x))$.

Consider several choices for $\pi(k)$ with the same mean λ_n depending possibly on n . For simplicity of exposition, assume that all samples have the same variance σ^2 and are of an equal size m so that $\gamma_i = \gamma = m\tau^2/\sigma^2$. For the binomial prior $B(n, \xi_n)$, where $\xi_n = \lambda_n/n$, (4.3) yields the same critical value for all tests

$$p^* = 2\tilde{\Phi}\left(\left(2(1 + 1/\gamma) \ln_+\left\{\sqrt{1 + \gamma} \frac{1 - \xi_n}{\xi_n}\right\}\right)^{1/2}\right). \tag{4.4}$$

From a frequentist viewpoint, it corresponds to testing each individual null hypothesis at the significance level p^* and, therefore, controlling the expected number of erroneously rejected null hypotheses at the level n_0p^* . As in Section 3.2., p^* is defined by the parameters of the prior and the noise level

instead of being fixed in advance at .01 or .05, and, in practice, it is typically estimated from the data (see Section 5.1.). Assuming $\lambda_n = o(n)$ (sparsity) for sufficiently large m , one gets a simpler approximation p_a^* for p^* :

$$\begin{aligned} p_a^* &\sim 2\tilde{\Phi} \left(\left(2 \ln \left\{ \sqrt{\gamma} \frac{1 - \xi_n}{\xi_n} \right\} \right)^{1/2} \right) \\ &\sim \frac{\xi_n}{1 - \xi_n} \left(\pi \gamma \ln \left\{ \sqrt{\gamma} \frac{1 - \xi_n}{\xi_n} \right\} \right)^{-1/2} \\ &\sim \frac{\xi_n}{\sqrt{\pi \gamma \ln(\sqrt{\gamma}/\xi_n)}}, \end{aligned} \quad (4.5)$$

where we exploited the well known asymptotic relation $\tilde{\Phi}(x) \sim \phi(x)/x$ for large x , where $\phi(\cdot)$ is the standard normal density (e.g., Barndorff-Nielsen and Cox, 1989, p.56). If, in addition, $\lambda_n < \sqrt{\pi \gamma \ln n}$ (super-sparse case), it follows immediately that the MAP multiple testing procedure with the binomial prior is a Bonferroni procedure with the FWE controlling level

$$\alpha_n = np^* \sim \frac{\lambda_n}{\sqrt{\pi \gamma \ln(\sqrt{\gamma}/\xi_n)}} (< 1) \quad (4.6)$$

(see (4.5)).

Consider now the truncated Poisson distribution $Pois^*(\lambda_n)$, where

$$\pi(k) = \frac{\lambda_n^k/k!}{\sum_{j=0}^n \lambda_n^j/j!}, \quad k = 0, \dots, n. \quad (4.7)$$

LEMMA 4.1. *Suppose k has the truncated Poisson distribution $Pois^*(\lambda_n)$ (4.7), where $\lambda_n = o(n)$. Then,*

- (i) $E(k) = \lambda_n(1 - \delta_n)$, where a positive sequence $\delta_n = o(1)$, and
- (ii) $k = o(n)$ almost surely.

The proof of Lemma 4.1 is given in the Appendix. From (4.3) and Lemma 4.1, one has

$$\begin{aligned} p_i^* &= 2\tilde{\Phi} \left(\left(2(1 + 1/\gamma) \ln_+ \left\{ \sqrt{1 + \gamma} \cdot \frac{n - i + 1}{\lambda_n} \right\} \right)^{1/2} \right) \\ &\sim \frac{\lambda_n}{n - i + 1} \cdot \frac{1}{\sqrt{\pi \gamma \ln(\sqrt{\gamma}/\xi_n)}}. \end{aligned}$$

Let again $\lambda_n < \sqrt{\pi\gamma \ln n}$. The corresponding step-up and step-down procedures in this case are, in fact, the FWE controlling procedures of Hochberg (1988) and Holm (1979) respectively with the same FWE level α_n in (4.6). Therefore, the original MAP procedure sandwiched between its two stepwise versions also asymptotically controls the FWE at the level α_n .

Consider also a “reflected” truncated Poisson distribution, where

$$\pi(k) = \frac{(n - \lambda_n)^{n-k} / (n - k)!}{\sum_{j=0}^n (n - \lambda_n)^{n-j} / (n - j)!}, \quad k = 0, \dots, n. \tag{4.8}$$

LEMMA 4.2. *Let $\pi(k)$ be given by (4.8), where $\lambda_n = o(n)$ but $\lambda_n / \sqrt{n \ln n} \rightarrow \infty$. Then,*

- (i) $E(k) = \lambda_n(1 + \delta_n)$, where a positive sequence $\delta_n = o(1)$, and
- (ii) $k = \lambda_n(1 + o(1))$ almost surely.

The proof of Lemma 4.2 is given in the Appendix. Lemma 4.2 and (4.3) imply

$$\begin{aligned} p_i^* &= 2\tilde{\Phi} \left(\left(2(1 + 1/\gamma) \ln_+ \left\{ \sqrt{1 + \gamma} \cdot \frac{n - \lambda_n}{i} \right\} \right)^{1/2} \right) \\ &\sim \frac{i}{n - \lambda_n} \cdot \frac{1}{\sqrt{\pi\gamma \ln(\sqrt{\gamma}/\xi_n)}} = \frac{i}{n} \cdot q_n, \end{aligned} \tag{4.9}$$

where $q_n = (1 - \xi_n)^{-1}(\pi\gamma \ln(\sqrt{\gamma}/\xi_n))^{-1/2}$. Within a frequentist framework, the corresponding step-up procedure mimics the well-known step-up FDR controlling procedure of Benjamini and Hochberg (1995) with the FDR parameter q_n . They proved that its FDR level does not exceed $(n_0/n)q_n \leq q_n$. Later, Benjamini and Yekutieli (2001) established that the FDR level is, in fact, exactly $(n_0/n)q_n$. Sarkar (2002) showed that it is also true for the corresponding step-down procedure. In addition, Lemma 4.2 indicates that for sparse cases typically $\hat{k}_l \sim \hat{k} \sim \hat{k}_r \sim \lambda (\leq E(k))$ and all three versions of the MAP testing procedure yield similar results (see also Section 5). The above examples indicate again that the traditional frequentist procedures focusing mainly on the control of erroneous rejections of null hypotheses, correspond to *sparse* priors, where the expected proportion of false null hypotheses tends to zero. For extremely conservative FWE controlling procedures, the expected number of false null hypotheses grows at most at logarithmic rate (*super-sparse* priors), while for less stringent FDR controlling procedures it grows as n^β , $0 < \beta < 1$. We have discussed the reasons for such sparsity of priors in Section 3.2.

5 Simulated Example

To test the performance of the proposed procedure, we considered a simulated example, where the data was generated according to the normal model considered in Section 4 :

$$Y_{i,j} = \mu_i + \epsilon_{i,j} \quad i = 1, \dots, n; \quad j = 1, \dots, m, \quad (5.1)$$

where the $\epsilon_{i,j}$'s are i.i.d. $N(0, \sigma^2)$ variables. Such an example is similar to that of Ishwaran and Rao (2003) motivated by the analysis of microarray data. In microarray experiments, expression levels of thousands of genes present in a biological sample are simultaneously measured to identify a small proportion of differentially expressed genes. The model (5.1) is supposed to mimic an experiment where two groups of biological samples (control and treatment) are marked using two dyes (usually Cy5 and Cy3 for red and green respectively) and then hybridized on several cDNA microarrays. In such experimental context, the value $Y_{i,j}$ represents the log intensity ratio between the control and the treatment groups for the i -th gene on the j -th array after a suitable normalization procedure. To detect differentially expressed genes, one simultaneously tests $H_{0i} : \mu_i = 0$ against $H_{1i} : \mu_i \neq 0$ for $i = 1, \dots, n$.

Let $\pi(k)$ be a specified prior distribution on the number of differentially expressed genes with $E(k) = \lambda$. Given k , all $\binom{n}{k}$ possible configurations were assumed to be equally likely. Unexpressed genes naturally had a zero expression level while the expression levels μ_i of differentially expressed genes were simulated as independent $N(0, \tau^2)$ variables. A similar simulation model was considered in Ishwaran and Rao (2003) but with the same *fixed* μ for all false nulls. We believe that our model is more realistic for the microarray data where gene expression levels may be different.

It should be noted that model (5.1) might be, in fact, somewhat idealized for real microarray data. In particular, it ignores co-regulations between gene expressions. Storey and Tibshirani (2001) argue that these co-regulations are commonly of "clumpy" type, where different genes can be clustered in small independent groups having within group correlations. Several simulation studies (see e.g. Storey and Tibshirani, 2001; Ishwaran and Rao, 2003; Reiner, Yekutieli and Benjamini, 2003) indicate though that clumpy dependencies have a relatively minor effect when group sizes are small compared to the total number of tested hypotheses — a very reasonable scenario with microarray data. In addition, the assumption of equal variances is not always valid. The proposed MAP testing procedure can still be adapted for

different σ_i^2 (see Section 4) but estimation of the parameters of the model described in Section 5.1 below should be modified in this case. Alternatively, sometimes the data can be transformed first through a variance-stabilization transformation (e.g. logarithmic).

5.1. Estimation of parameters. To apply the developed MAP procedures for a chosen $\pi(k)$, one still needs to specify the noise variance σ^2 and the prior variance τ^2 or, equivalently, the variance ratio $\gamma = m\tau^2/\sigma^2$. These parameters are rarely known *a priori* in practice and should be estimated from the data in the spirit of empirical Bayes. Moreover, we also assume that $\pi(k)$ has a prescribed parametric form with an unknown expectation λ that should be estimated as well.

Straightforward calculus yields the following marginal likelihood of the observed data $\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_n)'$, $i = 1, \dots, n$, where $\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,m})'$:

$$L(\mathbf{Y}; \lambda, \sigma^2, \gamma) = \sum_{k=0}^n \pi(k) \binom{n}{k}^{-1} \sum_{\mathbf{x}_j: \sum_i x_{ji}=k} \prod_{i=1}^n f_1(\mathbf{Y}_i)^{x_{ji}} f_2(\mathbf{Y}_i)^{1-x_{ji}},$$

where \mathbf{x}_j is the indicator vector, and f_1 and f_2 are m -variate normal densities with zero means and variance matrices $V = \sigma^2(I_m + (\gamma/m)\mathbf{1}\mathbf{1}')$ and $\sigma^2 I_m$ respectively. Following the ideas of empirical Bayes approach, we estimate the unknown parameters λ , σ^2 and γ by the corresponding marginal maximum likelihood estimators (MLEs). There are no closed form solutions and the EM algorithm is applied to obtain the MLEs numerically.

Regard the indicator vector \mathbf{x} as a “missing” data. The complete log-likelihood for the “augmented” data (\mathbf{Y}, \mathbf{x}) is then

$$l(\mathbf{Y}, \mathbf{x}; \lambda, \sigma^2, \gamma) = \ln \pi(k) - \ln \binom{n}{k} + \sum_{i=1}^n x_i \ln f_1(\mathbf{Y}_i) + \sum_{i=1}^n (1 - x_i) \ln f_2(\mathbf{Y}_i),$$

where $k = x_1 + \dots + x_n$.

At the h -th iteration, the E-step consists of computing the conditional expectation

$$\begin{aligned} \hat{l}_{[h]} &= E \left(l(\mathbf{Y}, \mathbf{x}) | \mathbf{Y}, \lambda_{[h]}, \sigma_{[h]}^2, \gamma_{[h]} \right) \\ &= E \left(\ln \left(\pi(k) \binom{n}{k}^{-1} \right) \middle| \mathbf{Y}, \lambda_{[h]} \right) + \sum \eta_{i[h]} \ln f_1(\mathbf{Y}_i) \\ &\quad + \sum_{i=1}^n (1 - \eta_{i[h]}) \ln f_2(\mathbf{Y}_i), \end{aligned}$$

where $\eta_{i[h]} = E\left(x_i | \mathbf{Y}_i, \lambda_{[h]}, \sigma_{[h]}^2, \gamma_{[h]}\right) = \left(1 + B_{i[h]} \frac{n - \lambda_{[h]}}{\lambda_{[h]}}\right)^{-1}$ and the Bayes factor $B_{i[h]} = (1 + \gamma_{[h]})^{1/2} \exp\left\{-m\bar{Y}_i^2 / \left(2\sigma_{[h]}^2 (1 + 1/\gamma_{[h]})\right)\right\}$ (cf. (4.1)).

Straightforward calculus shows that, regardless of the prior $\pi(k)$, maximizing $\hat{l}_{[h]}$ with respect to $\hat{\sigma}_{[h+1]}^2$ and $\hat{\gamma}_{[h+1]}$ under the positivity constraints on the M-step results in the following solutions :

$$\hat{\gamma}_{[h+1]} = \max\left(0, \frac{m \sum_{i=1}^n \eta_{i[h]} \bar{Y}_i^2}{\hat{\sigma}_{0[h+1]}^2 \sum_{i=1}^n \eta_{i[h]}} - 1\right),$$

$$\hat{\sigma}_{[h+1]}^2 = \begin{cases} \hat{\sigma}_{0[h+1]}^2 & \text{if } \hat{\gamma}_{[h+1]} > 0 \\ \sum_{i=1}^n \sum_{j=1}^m Y_{i,j}^2 / (nm) & \text{if } \hat{\gamma}_{[h+1]} = 0, \end{cases}$$

where

$$\hat{\sigma}_{0[h+1]}^2 = \frac{\sum_{i=1}^n \eta_{i[h]} \sum_{j=1}^m (Y_{i,j} - \bar{Y}_i)^2 + \sum_{i=1}^n (1 - \eta_{i[h]}) \sum_{j=1}^m Y_{i,j}^2}{nm - \sum_{i=1}^n \eta_{i[h]}}.$$

The solution for $\hat{\lambda}_{[h+1]}$ is a maximizer of

$$m(\lambda) = E\left\{\ln\left(\pi(k) \binom{n}{k}^{-1}\right) \mid \mathbf{Y}_1, \dots, \mathbf{Y}_n, \lambda\right\}$$

and it depends on the particular type of $\pi(k)$. Consider the three priors $\pi(k)$ discussed in the Section 4. For the binomial prior $B(n, \lambda/n)$, one has

$$\begin{aligned} m(\lambda) &= E(k | \mathbf{Y}_1, \dots, \mathbf{Y}_n, \lambda) \ln \frac{\lambda}{n - \lambda} + n \ln \frac{n - \lambda}{n} \\ &= \sum_{i=1}^n \eta_i \ln \frac{\lambda}{n - \lambda} + n \ln \frac{n - \lambda}{n}, \end{aligned}$$

and solving the equation $m'(\lambda) = 0$ immediately yields $\hat{\lambda}_{[h+1]} = \sum_{i=1}^n \eta_{i[h]}$.

For the truncated Poisson prior $Pois^*(\lambda)$ in (4.7),

$$m'(\lambda) = \frac{\sum_{i=1}^n \eta_i - \lambda}{\lambda} + \frac{\lambda^n / n!}{\sum_{j=0}^n \lambda^j / j!}.$$

After finding the root of $m'(\lambda)$ under the conditions of Lemma 4.1, one has $\hat{\lambda}_{[h+1]} = \sum_{i=1}^n \eta_{i[h]} (1 + o(1)) \sim \sum_{i=1}^n \eta_{i[h]}$.

Similarly, for the reflected truncated Poisson prior (4.8),

$$m'(\lambda) = \frac{\sum_{i=1}^n \eta_i - \lambda}{n - \lambda} + \frac{(n - \lambda)^n / n!}{\sum_{j=0}^n (n - \lambda)^j / (n - j)!}$$

and, under the conditions of Lemma 4.2, again $\hat{\lambda}_{[h+1]} = \sum_{i=1}^n \eta_{i[h]}(1+o(1)) \sim \sum_{i=1}^n \eta_{i[h]}$. Hence, for all the above three priors $\pi(k)$, the EM algorithm results in essentially the same MLEs for γ , σ^2 and λ .

5.2. *The results.* Data were generated according to the model (5.1), where the μ_i 's were randomly sampled from $N(0, \tau^2)$ in ξ percent of the cases, and set to zero in the remaining cases. This corresponds to a scenario where only $\xi\%$ of genes are differentially expressed. We set $\sigma = 1$ and $\tau = 2$. To mimic a typical microarray experiment, we choose a large number of tested hypotheses (genes) $n = 10000$ with a small proportion $\xi = .05$ of false nulls (differentially expressed genes) corresponding to a sparse (though not super-sparse) case. Several values of m were also tried. The number of replications was 1000.

The true values of σ^2 , τ^2 and ξ were assumed to be unknown in simulations and were estimated by the EM-algorithm described in the previous Section 5.1. We investigated two MAP multiple testing procedures corresponding to the binomial prior $B(n, \xi)$ (Binomial) and the reflected truncated Poisson prior (4.8), where $\pi(k) \propto (n - n\xi)^{n-k} / (n - k)!$ (Pois2). For both priors $\lambda = E(k) \sim n\xi$ and the EM-estimates for the parameters are the same (see Section 5.1.). The entire study was carried out using the MATLAB programming environment.

Table 1 summarizes the EM estimates of ξ , σ^2 and $\gamma = m\tau^2/\sigma^2$ averaged over 1000 replications and their standard deviations. As expected, for the given variance ratio τ^2/σ^2 , the accuracy of estimation improves as γ increases but even for moderate γ it is quite satisfactory. Similar results were obtained for $\xi = .01, .1, .2, .3$.

TABLE 1. MEANS AND STANDARD DEVIATIONS FOR THE EM-ESTIMATES OF THE UNKNOWN PARAMETERS ξ , σ^2 AND γ AVERAGED OVER 1000 REPLICATIONS. THE TRUE VALUES ARE $\xi = .05$, $\sigma = 1$ AND $\tau^2 = 4$.

m	γ	$\hat{\xi}$	$\hat{\sigma}^2$	$\hat{\gamma}$
5	20	.04859 (.00261)	1.00060 (.00657)	21.60540 (1.66995)
10	40	.04960 (.00194)	1.00030 (.00457)	41,37186 (2.91361)
20	80	.04988 (.00152)	1.00013 (.00312)	81.54536 (5.52843)
30	120	.04988 (.00138)	1.00011 (.00257)	121.49824 (8.52774)

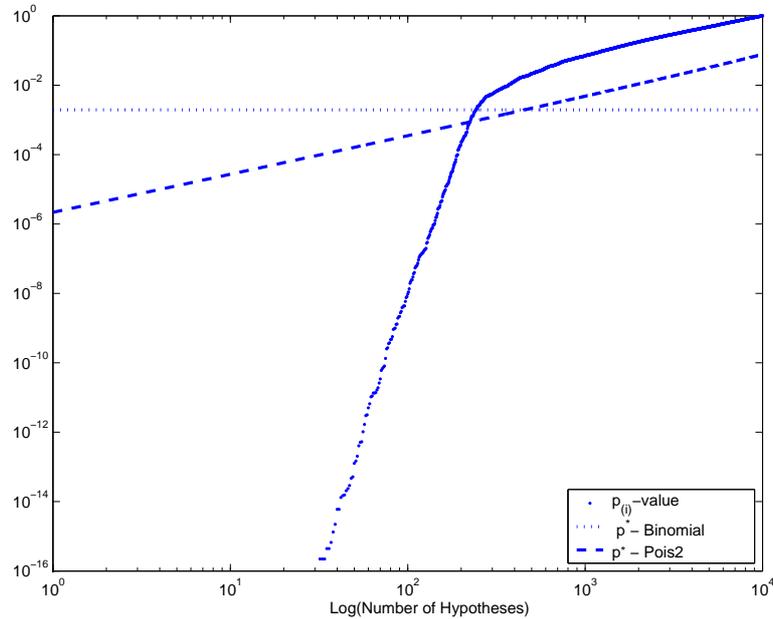


Figure 1. Ordered p -values $p_{(i)}$ and the corresponding critical values p_i^* for Binomial and Pois2 on the log-log scale for a single realization from (5.1) with $n = 10000$, $m = 5$, $\sigma = 1$ and $\xi_n = .05$.

Figure 1 shows the ordered p -values as in (4.2) and the corresponding critical values p_i^* as in (4.4) and (4.9) for Binomial and Pois2 priors respectively on the log-log scale evaluated from a single realization from (5.1). In most realizations, there was a single intersection point for Pois2 prior, and therefore all the three versions of the MAP procedure were identical.

The main simulation results on the performance of MAP multiple testing procedures are presented in Table 2 and Table 3 for $m = 5$ (small samples) and $m = 20$ (fairly large samples) respectively. The tables give the total number of detections (rejected nulls), the total number of erroneous decisions (misclassifications), the proportions of erroneous rejections of nulls and alternatives (Type I and Type II errors), and the proportions of erroneously rejected null hypotheses among those being rejected (false discovery rates, FDR) and erroneously nonrejected null hypotheses among those being accepted (false negative rates, FNR) for the original MAP procedures and for their stepwise versions. They are also compared with the two well-known multiple testing procedures: the Bonferroni procedure with familywise error levels $\alpha = .01$ and $.05$ and the step-up FDR controlling procedure of Benjamini and Hochberg (1995) with the FDR parameters $q = .01$ (BH.01) and

.05 (BH.05). The results in Table 2 and Table 3 are averaged over 1000 replications. Standard errors are given in the brackets. The unknown parameters were estimated by the EM-algorithm.

TABLE 2. PERFORMANCES OF MAP, BONFERRONI AND BENJAMINI-HOCHBERG MULTIPLE TESTING PROCEDURES ($m = 5, \xi = .05$).

Method	Detect	Miss-classified	Type I Error	Type II Error	FDR	FNR
Binomial	270 (.45124)	269 (.38654)	.00207 (.00002)	.49874 (.00078)	.07249 (.00052)	.02563 (.00004)
Pois2 (step-down)	244 (.43028)	274 (.39577)	.00095 (.00001)	.52938 (.00080)	.03657 (.00041)	.02713 (.00004)
Pois2 (step-up)	244 (.43015)	274 (.39571)	.00095 (.00001)	.52935 (.00080)	.03658 (.00041)	.02713 (.00004)
Pois2 (global)	244 (.43022)	274 (.39573)	.00095 (.00001)	.52937 (.00080)	.03657 (.00041)	.02713 (.00004)
Bonferroni (.05)	160 (.33067)	341 (.33058)	.00000 (.00000)	.68096 (.00066)	.00024 (.00004)	.03460 (.00003)
Bonferroni (.01)	143 (.31796)	357 (.31800)	.00000 (.00000)	.71397 (.00064)	.00006 (.00002)	.03622 (.00003)
BH.05	253 (.42965)	271 (.38663)	.00125 (.00001)	.51854 (.00078)	.04690 (.00044)	.02660 (.00004)
BH.01	212 (.38999)	292 (.38262)	.00021 (.00000)	.58081 (.00077)	.00941 (.00022)	.02967 (.00004)

TABLE 3. PERFORMANCES OF MAP, BONFERRONI AND BENJAMINI-HOCHBERG MULTIPLE TESTING PROCEDURES ($m = 20, \xi = .05$).

Method	Detect	Miss-classified	Type I Error	Type II Error	FDR	FNR
Binomial	372 (.35794)	152 (.33236)	.00123 (.00001)	.28011 (.00065)	.03131 (.00030)	.01454 (.00003)
Pois2 (step-down)	364 (.35089)	152 (.33121)	.00087 (.00001)	.28806 (.00065)	.02270 (.00026)	.01495 (.00003)
Pois2 (step-up)	364 (.35077)	152 (.33098)	.00087 (.00001)	.28805 (.00065)	.02271 (.00026)	.01495 (.00003)
Pois2 (global)	364 (.3584)	152 (.33111)	.00087 (.00001)	.28806 (.00065)	.02270 (.00026)	.01495 (.00003)
Bonferroni (.05)	306 (.34086)	194 (.34124)	.00001 (.00000)	.38774 (.00068)	.00017 (.00002)	.02000 (.00003)
Bonferroni (.01)	294 (.34688)	206 (.34683)	.00000 (.00000)	.41204 (.00069)	.00003 (.00001)	.02127 (.00003)
BH.05	383 (.36714)	153 (.33741)	.00191 (.00002)	.26946 (.00064)	.04711 (.00036)	.01401 (.00003)
BH.01	349 (.33686)	157 (.33830)	.00034 (.00001)	.30845 (.00066)	.00928 (.00017)	.01598 (.00003)

The Bonferroni procedures, as expected, were extremely conservative and resulted in the lowest numbers of detections and poor power properties. To better understand the results in Table 2 and Table 3, recall that for the binomial prior, the expected Type I error of the resulting MAP procedure is p^* given by (4.4). Table 4 compares Type I error from simulations with p^* and its simpler approximation (for large γ and small ξ) p_a^* from (4.5) calculated using the true parameters. The Type I error reported in the table is the average over 1000 replications, and the unknown parameters are estimates with the procedure described in Section 5.1.. One sees that simulated Type I errors and p^* match nicely. BH.01 and BH.05 controlled FDR similar to the expected levels $.95 \times .01 = .0095$ and $.95 \times .05 = .0475$ respectively. For all cases, the simulated FDR levels of Pois2 procedures were between those of BH.01 and BH.05, and that explains why the results of Pois2 procedures were sandwiched between more conservative BH.01 and less stringent BH.05. In fact, from the results of Section 4, for large γ and small ξ , the step-up and step-down Pois2 should behave similar to the FDR controlling procedures with the (adaptive) FDR level of about $(n_0/n)q_n \sim 1/\sqrt{\pi\gamma \ln(\sqrt{\gamma}/\xi)}$ (see (4.9)).

TABLE 4. SIMULATED TYPE I ERRORS FOR THE BINOMIAL PRIOR AND THE CORRESPONDING p^* AND p_a^* .

ξ	$m = 5$ ($\gamma = 20$)			$m = 20$ ($\gamma = 80$)		
	simul.	p^*	p_a^*	simul.	p^*	p_a^*
.01	.00032	.00038	.00051	.00021	.00020	.00024
.05	.00207	.00219	.00202	.01230	.00125	.00138
.1	.00491	.00519	.00438	.02850	.00285	.00297
.2	.01265	.01345	.01431	.00699	.00706	.00647
.3	.02404	.02570	.02302	.01285	.01302	.01027

Table 5 shows simulated FDR for the Pois2 prior and the corresponding approximated expected FDR from (4.9). The results show that for $m = 20$ ($\gamma = 80$), simulated FDR of Pois2 approaches the expected level. Binomial MAP procedure tends to reject more null hypotheses than Pois2. However, fairly large differences in Type I errors in favour of Pois2 and relatively negligible differences in Type II errors indicate that most of the additional detections made by Binomial were false. In fact, this is not surprising since for sparse (but not super-sparse) cases, Binomial essentially is not supposed to control any multiplicity criterion (FWE or FDR) while Pois2 is similar to the FDR controlling procedure. Finally, note that the performance of all the considered procedures improves with increasing γ .

TABLE 5. SIMULATED FDR FOR THE POIS2 PRIOR AND THE CORRESPONDING APPROXIMATED EXPECTED FDR FROM (4.9).

ξ	procedure	$m = 5$ ($\gamma = 20$)		$m = 20$ ($\gamma = 80$)	
		simul.	approx.	simul.	approx.
.01	Step-Down	.02945		.01957	
	Step-Up	.02957	.05106	.01957	.02419
	Global	.02948		.01957	
.05	Step-Down	.03657		.02270	
	Step-Up	.03658	.05951	.02271	.02769
	Global	.03657		.02270	
.1	Step-Down	.04095		.02497	
	Step-Up	.04097	.06471	.02497	.02975
	Global	.04096		.02497	
.2	Step-Down	.04702		.02719	
	Step-Up	.04703	.07156	.02719	.03235
	Global	.04702		.02719	
.3	Step-Down	.05226		.02918	
	Step-Up	.05227	.07675	.02918	.03423
	Global	.05227		.02918	

6 Concluding Remarks

The paper considered the multiple hypothesis testing within a Bayesian framework. We proposed a hierarchical prior model, where a prior distribution $\pi(k)$ is imposed on the number of hypotheses arising from alternatives, and then applied the maximum *a posteriori* (MAP) rule to find the most plausible configuration of true and false null hypotheses. In the case of independent likelihoods, the resulting MAP procedure and its closely related step-up and step-down versions are intuitively clear and computationally inexpensive. They compare ordered Bayes factors of individual hypotheses with a sequence of critical values depending on the prior and, in this sense, are similar in spirit to frequentist multiple testing procedures based on ordered p -values. By choosing different $\pi(k)$, one can mimic various existing frequentist testing procedures. In particular, for the normal data, the specific choices of $\pi(k)$ lead to several known FWE and FDR controlling procedures. We showed that the FDR controlling procedures are related to sparse $\pi(k)$, while their much more conservative FWE counterparts correspond to super-sparse priors. Furthermore, the resulting FWE and FDR levels for the Bayesian procedures are defined by the parameters of the prior rather than being fixed in advance at some traditional level (e.g., .01 or .05). In practice, the prior parameters are typically unknown and should be estimated from the data, which makes the corresponding error levels data-adaptive.

The Bayesian MAP testing procedures demonstrated good performance in a simulated example that mimics a microarray experiment.

The proposed general Bayesian approach can be used as a tool in a wide range of multiple hypothesis testing set-ups for various types of data though further analysis is needed for each specific problem at hand. Its use for model selection is another interesting topic for further research.

Appendix: Proofs of Lemmas from Section 4

6.1. Proof of Lemma 4.1.

(i) Let $a_k = \lambda_n^k/k!$. We have

$$E(k) = \frac{\sum_{k=0}^n k a_k}{\sum_{k=0}^n a_k} = \lambda_n \frac{\sum_{k=0}^{n-1} a_k}{\sum_{k=0}^n a_k} = \lambda_n \left(1 - \frac{a_n}{\sum_{k=0}^n a_k} \right) < \lambda_n. \quad (6.1)$$

On the other hand, the ratio $a_n/\sum_{k=0}^n a_k$ in (6.1) satisfies

$$\frac{a_n}{\sum_{k=0}^n a_k} < \frac{a_n}{\max_k a_k} = \frac{a_n}{a_{\lambda_n}} = \frac{\lambda_n^n \lambda_n!}{n! \lambda_n^{\lambda_n}}. \quad (6.2)$$

Exploiting the Stirling formula, for every k , one has

$$k^k e^{-k} \sqrt{2\pi k} < k! < k^k e^{-k + \frac{1}{12k}} \sqrt{2\pi k}, \quad (6.3)$$

and, hence, for $\lambda_n = o(n)$, (6.2) implies

$$\frac{a_n}{\sum_{k=0}^n a_k} < \left(\frac{\lambda_n}{n} e^{1-\lambda_n/n} \right)^n \sqrt{\frac{\lambda_n}{n}} e^{1/(12\lambda_n)} \rightarrow 0$$

as n increases.

(ii) Consider an arbitrary sequence $\epsilon_n = o(n)$ and let $S_n = P(k \geq \lambda_n + \epsilon_n)$. In what follows, we will find a particular ϵ_n such that the series $\sum_n S_n$ will converge. The first Borel-Cantelli lemma will imply then that, for $\lambda_n = o(n)$, $k = o(n)$ almost surely. One has

$$\begin{aligned} S_n &= \frac{\sum_{k=\lambda_n+\epsilon_n}^n a_k}{\sum_{k=0}^n a_k} < \frac{a_{\lambda_n+\epsilon_n} (n - \lambda_n - \epsilon_n)}{a_{\lambda_n}} \\ &= \frac{\lambda_n^{\lambda_n+\epsilon_n}}{\lambda_n^{\lambda_n}} \cdot \frac{\lambda_n! (n - \lambda_n - \epsilon_n)}{(\lambda_n + \epsilon_n)!} < n \lambda_n^{\epsilon_n} \frac{\lambda_n!}{(\lambda_n + \epsilon_n)!}. \end{aligned} \quad (6.4)$$

Applying (6.3) to get respectively the upper and lower bounds for the factorials in the numerator and denominator of (6.4), straightforward calculus yields

$$S_n < n\kappa_n e^{\frac{1}{12\lambda_n}} < Cn\kappa_n$$

for some $C > 0$, where

$$\kappa_n = \frac{e^{\epsilon_n}}{(1 + \epsilon_n/\lambda_n)^{\lambda_n + \epsilon_n + 1/2}}.$$

Hence,

$$\ln \kappa_n < \epsilon_n \left(1 - \ln \frac{\epsilon_n}{\lambda_n} \right).$$

Consider $\epsilon_n = \max(\ln n, e^\zeta \lambda_n) = o(n)$, where $\zeta > 3$. Then,

$$S_n < Cne^{-\epsilon_n(\zeta-1)} < Cn^{-(\zeta-2)},$$

and, therefore, the series $\sum_n S_n$ converges. □

6.2. Proof of Lemma 4.2.

(i) Let $b_k = (n - \lambda_n)^k/k!$. After simple algebra one has

$$E(k) = \frac{\sum_{k=0}^n (n-k)b_k}{\sum_{k=0}^n b_k} = \lambda_n + \frac{(n - \lambda_n)b_n}{\sum_{k=0}^n b_k} > \lambda_n.$$

On the other hand, under the conditions of Lemma 4.2 on λ_n ,

$$\lambda_n + \frac{(n - \lambda_n)b_n}{\sum_{k=0}^n b_k} < \lambda_n + \frac{n - \lambda_n}{\lambda_n} = \lambda_n (1 - 1/\lambda_n + n/\lambda_n^2) < \lambda_n(1 + \delta_n),$$

where $\delta_n = o(1)$ and positive.

(ii) Consider the sequence b_k defined above. For any positive sequence ϵ_n define $S_{n1} = P(k \geq \lambda_n + \epsilon_n)$ and $S_{n2} = P(k \leq \lambda_n - \epsilon_n)$. Similar to the ideas of the proof of Lemma 4.1, we will find a particular $\epsilon_n = o(\lambda_n)$ guaranteeing the convergence of both $\sum_n S_{n1}$ and $\sum_n S_{n2}$ under the conditions of Lemma 4.2, and, hence, by the first Borel-Cantelli lemma, we have $k = \lambda_n(1 + o(1))$ with probability one. Note that $\max_k b_k = b_{n-\lambda_n}$, b_k increases for $k \leq n - \lambda_n$ and decreases otherwise. One has

$$S_{n1} = \frac{\sum_{k=0}^{n-\lambda_n-\epsilon_n} b_k}{\sum_{k=0}^n b_k} < \frac{nb_{n-\lambda_n-\epsilon_n}}{b_{n-\lambda_n}} = \frac{n}{(n - \lambda_n)^{\epsilon_n}} \cdot \frac{(n - \lambda_n)!}{(n - \lambda_n - \epsilon_n)!}.$$

Applying (6.3) and straightforward calculus, one gets

$$S_{n1} < n\kappa_{n1} \sqrt{\frac{n - \lambda_n}{n - \lambda_n - \epsilon_n}} e^{1/(12(n-\lambda_n))}, \tag{6.5}$$

where

$$\kappa_{n1} = \left(\left(1 + \frac{\epsilon}{n - \lambda_n - \epsilon_n} \right)^{(n-\lambda_n-\epsilon_n)/\epsilon_n} e^{-1} \right)^{\epsilon_n}.$$

The Taylor expansion of $\ln \kappa_{n1}$ implies

$$\ln \kappa_{n1} = -\frac{1}{2} \cdot \frac{\epsilon_n^2}{n - \lambda_n - \epsilon_n} (1 + o(1)).$$

Consider $\epsilon_n = \sqrt{cn \ln n}$. Under the conditions of Lemma 4.2 on λ_n , $\epsilon_n = o(\lambda_n)$, and one has

$$\kappa_{n1} < e^{-\frac{c}{2} \cdot \frac{n \ln n}{n} (1+o(1))} = n^{-\frac{c}{2}(1+o(1))}.$$

Thus, for any $c > 4$, $\sum_n n\kappa_{n1} < \infty$ and, therefore, $\sum_n S_{n1} < \infty$ since the factor $e^{1/(12(n-\lambda_n))} \sqrt{\frac{n-\lambda_n}{n-\lambda_n-\epsilon_n}}$ in (6.5) tends to one. Similarly, for S_{n2} , one has

$$\begin{aligned} S_{n2} &= \frac{\sum_{k=n-\lambda_n+\epsilon_n}^n b_k}{\sum_{k=0}^n b_k} < \frac{nb_{n-\lambda_n+\epsilon_n}}{b_{n-\lambda_n}} = \frac{n(n - \lambda_n)^{\epsilon_n} (n - \lambda_n)!}{(n - \lambda_n + \epsilon_n)!} \\ &< n\kappa_{n2} \sqrt{\frac{n - \lambda_n}{n - \lambda_n + \epsilon_n}} e^{1/(12(n-\lambda_n))}, \end{aligned}$$

where

$$\kappa_{n2} = \left(\left(1 - \frac{\epsilon}{n - \lambda_n + \epsilon_n} \right)^{(n-\lambda_n+\epsilon_n)/\epsilon_n} e \right)^{\epsilon_n},$$

and, therefore,

$$\ln \kappa_{n2} = -\frac{1}{2} \frac{\epsilon_n^2}{n - \lambda_n + \epsilon_n} (1 + o(1)).$$

Repeating the same arguments for κ_{n2} , for $\epsilon_n = \sqrt{cn \ln n}$, one has

$$\kappa_{n2} < n^{-\frac{c}{2}(1+o(1))},$$

and $\sum_n S_{n2}$ also converges for $c > 4$. □

Acknowledgment. The authors would like to thank Yoav Benjamini and Vadim Grinshtein for fruitful discussions, and the two anonymous referees for their helpful comments.

References

- BARNDORFF-NIELSEN, O.E. and COX, D.R. (1989). *Asymptotic Techniques for Use in Statistics*. Chapman and Hall, London.
- BAYARRI, M.J. and BERGER, J.O. (2004). The interplay of Bayesian and frequentist analysis. *Statist. Science*, **19**, 58–80.
- BENJAMINI, Y. and HOCHBERG, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing, *J. Roy. Statist. Soc. Ser. B* **57**, 289–300.
- BENJAMINI, Y., KRIEGER, A.M. and YEKUTIELI, D. (2006). Adaptive linear step-up and step-down procedures that control the false discovery rate, *Biometrika* **93**, 491–507.
- BENJAMINI, Y. and YEKUTIELI, D. (2001). The control of the false discovery rate in multiple testing under dependency, *Ann. Statist.* **29**, 1165–1188.
- BERGER, J.O. (1985). *Statistical Decision Theory and Bayesian Analysis* (2nd ed.). Springer, New York.
- BERGER, J.O. (2003). Could Fisher, Jeffreys and Neymann have agreed on testing? (with discussion), *Statist. Science* **18**, 1–32.
- BERGER, J.O., BOUKAI, B. and WANG, Y. (1997). Unified frequentist and Bayesian testing of a precise hypothesis, *Statist. Science* **12**, 133–160.
- BERGER, J.O. and PERICCHI, L.R. (1996). The intrinsic Bayes factor for model selection and prediction, *J. Amer. Statist. Assoc.* **91**, 109–122.
- BERRY, D.A. and HOCHBERG, Y. (1999). Bayesian perspectives on multiple comparisons, *J. Statist. Plann. Infer.* **82**, 215–227.
- CASELLA, G. and BERGER, R.L. (1987). Reconciling Bayesian and frequentist evidence in the one-sided testing problem, *J. Amer. Statist. Assoc.* **82**, 106–111.
- GOOD, I.J. (1992). The Bayes/non-Bayes compromise: A brief review, *J. Amer. Statist. Assoc.* **87**, 597–606.
- GOPALAN, R. and BERRY, D.A. (1998). Bayesian multiple comparisons using Dirichlet process priors, *J. Amer. Statist. Assoc.* **93**, 1130–1139.
- HOCHBERG, Y. (1988). A sharper Bonferroni procedure for multiple tests of significance, *Biometrika* **75**, 800–802.
- HOCHBERG, Y. and TAMHANE, A.C. (1987). *Multiple Comparisons Procedures*. Wiley, New York.
- HOLM, S. (1979). A simple sequentially rejective multiple test procedure, *Scand. J. Statist.* **6**, 65–70.
- HSU, J.C. (1996). *Multiple Comparisons: Theory and Methods*. Chapman & Hall/CRC, Washington D.C..
- ISHWARAN, H. and RAO, J.S. (2003). Detecting differentially expressed genes in microarrays using Bayesian model selection, *J. Amer. Statist. Assoc.* **98**, 438–455.
- KASS, R.E. and RAFTERY, A.E. (1995). Bayes factors, *J. Amer. Statist. Assoc.* **430**, 773–795.
- O’HAGAN, A. (1995). Fractional Bayes factors for model comparison (with discussion), *J. Roy. Statist. Soc. Ser. B*, **56**, 99–118.
- REINER, A., YEKUTIELI, D. and BENJAMINI, Y. (2003). Identifying differentially expressed genes using false discovery rate controlling procedures, *Bioinformatics* **19**, 368–375.

- RUBIN, D.B. (1984). Bayesianly justifiable and relevant frequency calculations for the applied statistician, *Ann. Statist.* **12**, 1151–1172.
- SARKAR, S.K. (2002). Some results on false discovery rate in stepwise multiple testing procedures, *Ann. Statist.* **30**, 239–257.
- SARKAR, S.K. and CHEN, J. (2004). A Bayesian stepwise multiple testing procedure, *Technical Report, Temple University*.
- SCOTT, J.G. and BERGER, J.O. (2006). An exploration of aspects of Bayesian multiple testing, *J. Statist. Plann. Infer.* **136**, 2144–2162.
- STOREY, J.D. and TIBSHIRANI, R. (2001). Estimating false discovery rates under dependence, with applications to DNA microarrays, *Technical Report 2001-28*, Dep. of Statistics, Stanford University.
- WALLER, R.A. and DUNCAN, D.B. (1969). A Bayes rule for the symmetric multiple comparison problem, *J. Amer. Statist. Assoc.* **64**, 1184–1503.

FELIX ABRAMOVICH
DEPARTMENT OF STATISTICS
AND OPERATIONS RESEARCH
TEL AVIV UNIVERSITY
TEL AVIV 69978, ISRAEL
E-mail: felix@post.tau.ac.il

CLAUDIA ANGELINI
ISTITUTO PER LE APPLICAZIONI
DEL CALCOLO MAURO PICONE
CONSIGLIO NAZIONALE DELLE RICERCHE
80131 NAPOLI, ITALY
E-mail: c.angelini@iac.cnr.it

Paper received October 2005; revised May 2006.