

# False discovery rates: a new deal

MATTHEW STEPHENS\*

*Department of Statistics and Department of Human Genetics, University of Chicago,  
5801 S Ellis Ave, Chicago, IL 60637 USA*  
[mstephens@uchicago.edu](mailto:mstephens@uchicago.edu)

## SUMMARY

We introduce a new Empirical Bayes approach for large-scale hypothesis testing, including estimating false discovery rates (FDRs), and effect sizes. This approach has two key differences from existing approaches to FDR analysis. First, it assumes that the distribution of the actual (unobserved) effects is unimodal, with a mode at 0. This “unimodal assumption” (UA), although natural in many contexts, is not usually incorporated into standard FDR analysis, and we demonstrate how incorporating it brings many benefits. Specifically, the UA facilitates efficient and robust computation—estimating the unimodal distribution involves solving a simple convex optimization problem—and enables more accurate inferences provided that it holds. Second, the method takes as its input two numbers for each test (an effect size estimate and corresponding standard error), rather than the one number usually used (*p* value or *z* score). When available, using two numbers instead of one helps account for variation in measurement precision across tests. It also facilitates estimation of effects, and unlike standard FDR methods, our approach provides interval estimates (credible regions) for each effect in addition to measures of significance. To provide a bridge between interval estimates and significance measures, we introduce the term “local false sign rate” to refer to the probability of getting the sign of an effect wrong and argue that it is a superior measure of significance than the local FDR because it is both more generally applicable and can be more robustly estimated. Our methods are implemented in an R package `ashr` available from <http://github.com/stephens999/ashr>.

*Keywords:* Empirical Bayes; False discovery rates; Multiple testing; Shrinkage; Unimodal.

## 1. INTRODUCTION

Since its introduction in Benjamini and Hochberg (1995), the “False Discovery Rate” (FDR) has quickly established itself as a key concept in modern statistics, and the primary tool by which most practitioners handle large-scale multiple testing in which the goal is to identify the non-zero “effects” among a large number of imprecisely measured effects.

Here we consider an Empirical Bayes (EB) approach to FDR. This idea is, of course, far from new: indeed, the notion that EB approaches could be helpful in handling multiple comparisons predates introduction of the FDR (e.g. Greenland and Robins, 1991). More recently, EB approaches to the FDR have been extensively studied by several authors, especially Efron and co-workers (Efron and others, 2001;

\*To whom correspondence should be addressed.

© The Author 2016. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Efron and Tibshirani, 2002; Efron *and others*, 2003; Efron, 2008, 2010); see also Kendziorski *and others* (2003), Newton *and others* (2004), Datta and Datta (2005), and Muralidharan (2010), for example.

So what is the “New Deal” here? We introduce two simple ideas that are new (at least compared with existing widely used FDR pipelines) and can substantially affect inference. The first idea is to assume that the distribution of effects is unimodal. We provide a simple, fast, and stable computer implementation for performing EB inference under this assumption, and illustrate how it can improve inferences when the unimodal assumption (UA) is correct. The second idea is to use two numbers—effect sizes and their standard errors—rather than just one— $p$  values or  $z$  scores—to summarize each measurement. Here we use this idea to allow variation in measurement precision to be better accounted for, avoiding a problem with standard pipelines that poor-precision measurements can inflate estimated FDR. (Ploner *and others* (2006) also suggest using more than one number in FDR analyses, taking a rather different approach to the one used here.)

In addition to these two new ideas, we highlight a third idea that is old but which remains under-used in practice: the idea that it may be preferable to focus on estimation rather than on testing. In principle, Bayesian approaches can naturally unify testing and estimation into a single framework—testing is simply estimation with some positive prior probability that the effect is exactly zero. However, despite ongoing interest in this area from both frequentist (Benjamini and Yekutieli, 2005) and Bayesian (Zhao and Gene Hwang, 2012; Gelman *and others*, 2012) perspectives, in practice large-scale studies that assess many effects almost invariably focus on testing significance and controlling the FDR, and not on estimation. To help provide a bridge between FDR and estimation we introduce the term “local false sign rate” ( $lfsr$ ), which is analogous to the “local false discovery rate” ( $lfdr$ ) (Efron, 2008), but which measures confidence in the sign of each effect rather than confidence in each effect being non-zero. We show that in some settings, particularly those with many discoveries, the  $lfsr$  and  $lfdr$  can be quite different and emphasize benefits of the  $lfsr$ , particularly its increased robustness to modeling assumptions.

Although we focus here on FDR applications, the idea of performing EB inference using a flexible unimodal prior distribution is useful more generally. For example, the methods described here can be applied directly to perform shrinkage estimation for wavelet denoising (Donoho and Johnstone, 1995), an idea explored in a companion paper (Xing and Stephens, 2016). And analogous ideas can be used to perform EB inference for variances (Lu and Stephens, 2016). Importantly, and perhaps surprisingly, our work demonstrates how EB inference under a general UA is, if anything, computationally simpler than commonly used more restrictive assumptions—such as a spike and slab or Laplace prior distribution (Johnstone and Silverman, 2004)—as well as being more flexible.

We refer to our EB method as adaptive shrinkage, or ash, to emphasize its key points: using a unimodal prior naturally results in shrinkage estimation, and the shrinkage is adaptive to both the amount of signal in the data and the measurement precision of each observation. We provide implementations in an R package, `ashr`, available at <http://github.com/stephens999/ashr>. Code and instructions for reproducing analyses and figures in this paper are at <https://github.com/stephenslab/ash>.

## 2. METHODS

### 2.1. Model outline

Here we describe the simplest version of the method, before briefly discussing embellishments we have also implemented. Implementation details are given in Supplementary Information, see supplementary material available at *Biostatistics* online.

Let  $\beta = (\beta_1, \dots, \beta_J)$  denote  $J$  “effects” of interest. For example, in a genomics application  $\beta_j$  might be the difference in the mean (log) expression of gene  $j$  in two conditions. We tackle both the problem

of testing the null hypotheses  $H_j : \beta_j = 0$ , and the more general problem of estimating, and assessing uncertainty in,  $\beta_j$ .

Assume that the available data are estimates  $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_J)$  of the effects, and corresponding (estimated) standard errors  $\hat{s} = (\hat{s}_1, \dots, \hat{s}_J)$ . Our goal is to compute a posterior distribution for  $\beta$  given  $\hat{\beta}, \hat{s}$ , which by Bayes theorem is

$$p(\beta|\hat{\beta}, \hat{s}) \propto p(\beta|\hat{s})p(\hat{\beta}|\beta, \hat{s}). \quad (2.1)$$

For  $p(\beta|\hat{s})$  we assume that the  $\beta_j$  are independent from a unimodal distribution  $g$ . This UA is a key assumption that distinguishes our approach from previous EB approaches to FDR analysis. A simple way to implement the UA is to assume that  $g$  is a mixture of a point mass at 0 and a mixture of *zero-mean* normal distributions:

$$p(\beta|\hat{s}, \pi) = \prod_j g(\beta_j; \pi), \quad (2.2)$$

$$g(\cdot; \pi) = \pi_0 \delta_0(\cdot) + \sum_{k=1}^K \pi_k N(\cdot; 0, \sigma_k^2), \quad (2.3)$$

where  $\delta_0(\cdot)$  denotes a point mass on 0, and  $N(\cdot; \mu, \sigma^2)$  denotes the density of the normal distribution with mean  $\mu$  and variance  $\sigma^2$ . Here we take  $\sigma_1, \dots, \sigma_K$  to be a large and dense grid of *fixed* positive numbers spanning a range from very small to very big (so  $K$  is fixed and large). We encourage the reader to think of this grid as becoming infinitely large and dense, as a non-parametric limit, although of course in practice we use a finite grid—see Implementation Details. The mixture proportions  $\pi = (\pi_0, \dots, \pi_K)$ , which are non-negative and sum to one, are hyper-parameters to be estimated.

For the likelihood  $p(\hat{\beta}|\beta, \hat{s})$  we assume a normal approximation:

$$p(\hat{\beta}|\beta, \hat{s}) = \prod_j N(\hat{\beta}_j; \beta_j, \hat{s}_j^2). \quad (2.4)$$

This simple model features both the key ideas we want to emphasize in this paper: the UA is encapsulated in (2.3) while the different measurement precision of different observations is encapsulated in the likelihood (2.4)—specifically, observations with larger standard error will have a flatter likelihood, and therefore have less impact on inference. However, the model also has several additional assumptions that can be relaxed. Specifically,

1. The form (2.3) implies that  $g$  is symmetric about 0. More flexibility can be obtained using mixtures of uniforms ([Supplementary Information](#), equation S.1.1); indeed this allows  $g$  to approximate *any* unimodal distribution.
2. The model (2.2) assumes that the effects are identically distributed, independent of their standard errors  $\hat{s}$ . This can be relaxed [see (3.2)].
3. The normal likelihood (2.4) can be generalized to a  $t$  likelihood [see ([Supplementary Information](#), equation S.1.2)].
4. We can allow for the mode of  $g$  to be non-zero, estimated from the data by maximum likelihood.

These embellishments, detailed in [supplementary material](#) available at *Biostatistics* online, are implemented in our software. Other limitations are harder to relax, most notably the independence and conditional independence assumptions (which are also made by most existing approaches). Correlations among tests certainly arise in practice, either due to genuine correlations or due to unmeasured confounders, and their potential impact on estimated FDRs is important to consider whatever analysis methods are used ([Efron, 2007](#); [Leek and Storey, 2007](#)).

## 2.2. Fitting the model

Together, (2.2)–(2.4) imply that  $\hat{\beta}_j$  are independent with

$$p(\hat{\beta}|\hat{s}, \pi) = \prod_j \left[ \sum_{k=0}^K \pi_k N(\hat{\beta}_j; 0, \sigma_k^2 + \hat{s}_j^2) \right], \quad (2.5)$$

where we define  $\sigma_0 := 0$ .

The usual EB approach to fitting this model would involve two steps:

1. Estimate the hyper-parameters  $\pi$  by maximizing the likelihood  $L(\pi)$ , given by (2.5), yielding  $\hat{\pi} := \arg \max L(\pi)$ .
2. Compute quantities of interest from the conditional distributions  $p(\beta_j|\hat{\beta}, \hat{s}, \hat{\pi})$ . For example, the evidence against the null hypothesis  $\beta_j = 0$  can be summarized by  $p(\beta_j \neq 0|\hat{\beta}, \hat{s}, \hat{\pi})$ .

Both steps are straightforward. Step 1 is a convex optimization problem, and can be solved quickly and reliably using interior point methods (Boyd and Vandenberghe, 2004; Koenker and Mizera, 2013). (Alternatively a simple EM algorithm can also work well, particularly for modest  $J$ ; see <http://stephenslab.github.io/ash/analysis/IPvsEM.html>.) And the conditional distributions  $p(\beta_j|\hat{\beta}_j, \hat{s}_j, \hat{\pi})$  in Step 2 are analytically available, each a mixture of a point mass on zero and  $K$  normal distributions. The simplicity of Step 1 is due to our use of a fixed grid for  $\sigma_k$  in (2.3), instead of estimating  $\sigma_k$ . This simple device may be useful in other applications.

Here we slightly modify this usual procedure: instead of obtaining  $\hat{\pi}$  by maximizing the likelihood, we maximize a penalized likelihood [see (Supplementary Information, equation S.2.5)], where the penalty encourages  $\hat{\pi}_0$  to be as big as possible whilst remaining consistent with the observed data. We introduce this penalty because in FDR applications it is considered desirable to avoid underestimating  $\pi_0$  so as to avoid underestimating the FDR.

Our R package implementation typically takes 20 s on a modern laptop for  $J = 100\,000$ , and scales linearly with  $J$ .

## 2.3. The local FDR and local false sign rate

As noted above, the posterior distributions  $p(\beta_j|\hat{\beta}, \hat{s}, \hat{\pi})$  have a simple analytic form. In practice it is common, and desirable, to summarize these distributions to convey the “significance” of each observation  $j$ . One natural measure of the significance of observation  $j$  is its “local FDR” (Efron, 2008)

$$\text{lfdr}_j := \Pr(\beta_j = 0|\hat{\beta}, \hat{s}, \hat{\pi}). \quad (2.6)$$

In words,  $\text{lfdr}_j$  is the probability, given the observed data, that effect  $j$  would be a false discovery, if we were to declare it a discovery.

The  $\text{lfdr}$ , like most measures of significance, is rooted in the hypothesis testing paradigm which focuses on whether or not an effect is exactly zero. This paradigm is popular, despite the fact that many statistical practitioners have argued that it is often inappropriate because the null hypothesis  $H_j : \beta_j = 0$  is often implausible. For example, Tukey (1991) argued that “All we know about the world teaches us that the effects of  $A$  and  $B$  are always different—in some decimal place—for any  $A$  and  $B$ . Thus asking ‘Are the effects different?’ is foolish.” Instead, Tukey (Tukey, 1962, page 32) suggested that one should address

... the more meaningful question: “is the evidence strong enough to support a belief that the observed difference has the correct sign?”

Along the same lines, Gelman and co-workers (Gelman and Tuerlinckx, 2000; Gelman *and others*, 2012) suggest focusing on “type S errors,” meaning errors in sign, rather than the more traditional type I errors.

Motivated by these suggestions, we define the “local False Sign Rate” for effect  $j$ ,  $lfsr_j$ , to be the probability that we would make an error in the sign of effect  $j$  if we were forced to declare it either positive or negative. Specifically,

$$lfsr_j := \min[\Pr(\beta_j \geq 0 | \hat{\pi}, \hat{\beta}, s), \Pr(\beta_j \leq 0 | \hat{\pi}, \hat{\beta}, s)]. \quad (2.7)$$

To illustrate, suppose that

$$\Pr(\beta_j < 0 | \hat{\beta}, s, \hat{\pi}) = 0.95,$$

$$\Pr(\beta_j = 0 | \hat{\beta}, s, \hat{\pi}) = 0.03,$$

$$\Pr(\beta_j > 0 | \hat{\beta}, s, \hat{\pi}) = 0.02.$$

Then from (2.7)  $lfsr_j = \min(0.05, 0.98) = 0.05$  (and, from (2.6),  $lfdr_j = 0.03$ ). This  $lfsr$  corresponds to the fact that, given these results, we would guess that  $\beta_j$  is negative, with probability 0.05 of being wrong.

As our notation suggests,  $lfsr_j$  is analogous to  $lfdr_j$ : whereas small values of  $lfdr_j$  indicate that we can be *confident that  $\beta_j$  is non-zero*, small values of  $lfsr_j$  indicate that we can be *confident in the sign of  $\beta_j$* . Of course, being confident in the sign of an effect logically implies that we are confident it is non-zero, and this is reflected in the fact that  $lfsr_j \geq lfdr_j$  [this follows from the definition because both the events  $\beta_j \geq 0$  and  $\beta_j \leq 0$  in (2.7) include the event  $\beta_j = 0$ ]. In this sense  $lfsr$  is a more conservative measure of significance than  $lfdr$ . More importantly,  $lfsr$  is more robust to modeling assumptions (see Results).

From these “local” measures of significance, we can also compute average error rates over subsets of observations  $\Gamma \subset \{1, \dots, J\}$ . For example,

$$\widehat{FSR}(\Gamma) := (1/|\Gamma|) \sum_{j \in \Gamma} lfsr_j \quad (2.8)$$

is an estimate of the total proportion of errors made if we were to estimate the sign of all effects in  $\Gamma$  [an error measure analogous to the usual (tail) FDR]. And, we can define the “ $s$ -value”

$$s := \widehat{FSR}(\{k : lfsr_k \leq lfsr_j\}) \quad (2.9)$$

as a measure of significance analogous to Storey’s  $q$ -value (Storey, 2003). [Replacing  $lfsr_j$  with  $lfdr_j$  in (2.8) and (2.9) gives estimates for the usual FDR( $\Gamma$ ), and the  $q$ -values respectively.]

## 2.4. Related work

**2.4.1. Previous approaches focused on FDR** Among previous methods that explicitly consider FDR, our work seems naturally compared with the EB methods of Efron (2008) and Muralidharan (2010) (implemented in the R packages `locfdr` and `mixfdr`, respectively) and with the widely used methods from Storey (2003) (implemented in the R package `qvalue`), which although not formally an EB approach, shares some elements in common.

There are two key differences between our approach and these existing methods. First, whereas these existing methods take as input a single number—either a  $z$  score (`locfdr` and `mixfdr`), or a  $p$  value (`qvalue`)—for each effect, we instead work with two numbers ( $\hat{\beta}_j, \hat{s}_j$ ). Here we are building on Wakefield (2009), which develops Bayes Factors testing individual null hypotheses in a similar way; see also

Efron (1993). Using two numbers instead of one clearly has the potential to be more informative, as illustrated in Results (Figure 4).

Second, our UA that the effects are unimodal about zero, is an assumption not made by existing methods, and one that we argue to be both plausible and beneficial in many contexts. Although the UA will not hold in all settings, it will often be reasonable, especially in FDR-related contexts that focus on rejecting the null hypotheses  $\beta_j = 0$ . This is because if “ $\beta_j = 0$ ” is a plausible null hypothesis then “ $\beta_j$  very near 0” should also be plausible. Further, it seems reasonable to expect that larger effects become decreasingly plausible, and so the distribution of the effects will be unimodal about 0. To paraphrase Tukey, “All we know about the world teaches us that large effects are rare, whereas small effects abound.” We emphasize that the UA relates to the distribution of *all* effects, and not only the *detectable* effects (i.e. those that are significantly different from zero). It is very likely that the distribution of *detectable* non-zero effects will be multimodal, with one mode for detectable positive effects and another for detectable negative effects, and the UA does not contradict this.

In further support of the UA, note that large-scale regression methods almost always make the UA for the regression coefficients, which are analogous to the “effects” we estimate here. Common choices of unimodal distribution for regression coefficients include the spike and slab, Laplace, *t*, normal-gamma, normal-inverse-gamma, or horseshoe priors (Carvalho and others, 2010). These are all less flexible than our approach, which provides for general unimodal distributions, and it may be fruitful to apply our methods to the regression context; indeed see Erbe and others (2012) for work in this vein. Additionally, the UA can be motivated by its effect on point estimates, which is to “shrink” the estimates towards the mode—such shrinkage is desirable from several standpoints for improving estimation accuracy. Indeed most model-based approaches to shrinkage make parametric assumptions that obey the UA (e.g. Johnstone and Silverman, 2004).

Besides its plausibility, the UA has two important practical benefits: it facilitates more accurate estimates of FDR-related quantities, and it yields simple algorithms that are both computationally and statistically stable (see Results).

**2.4.2. Other work** There is also a very considerable literature that does not directly focus on the FDR problem, but which involves similar ideas and methods. Among these, a paper about deconvolution (Cordy and Thomas, 1997) is most similar, methodologically, to our work here: indeed, this paper includes all the elements of our approach outlined above, except for the point mass on 0 and corresponding penalty term. However, the focus is very different: Cordy and Thomas (1997) focuses entirely on estimating  $g$ , whereas our primary focus is on estimating  $\beta_j$ . Also, they provide no software implementation.

More generally, the related literature is too large to review comprehensively, but relevant key-words include “empirical Bayes,” “shrinkage,” “deconvolution,” “semi-parametric,” “shape-constrained,” and “heteroskedastic.” Some pointers to recent papers in which other relevant citations can be found include. Xie and others (2012), Sarkar and others (2014), Koenker and Mizera (2014), and Efron (2016). Much of the literature focusses on the homoskedastic case (i.e.  $\hat{s}_j$  all equal) whereas we allow for heteroskedasticity. And much of the recent shrinkage-oriented literature focuses only on point estimation of  $\beta_j$ , whereas for FDR-related applications measures of uncertainty are essential. Several recent papers consider more flexible non-parametric assumptions on  $g$  than the UA assumption we make here. In particular, Jiang and Zhang (2009) and Koenker and Mizera (2014) consider the unconstrained non-parametric maximum likelihood estimate (NPMLE) for  $g$ . These methods may be useful in settings where the UA assumption is considered too restrictive. However, the NPMLE for  $g$  is a discrete distribution, which will induce a discrete posterior distribution on  $\beta_j$ , and so although the NPMLE may perform well for point estimation, it may not adequately reflect uncertainty in  $\beta_j$ . To address this some regularization on  $g$  (e.g. as in Efron, 2016) may be necessary. Indeed, one way of thinking about the UA is as a way to regularize  $g$ .

### 3. RESULTS

We compare results of `ash` with existing FDR-based methods implemented in the R packages `qvalue` (v2.1.1), `locfdr` (v1.1-8), and `mixfdr` (v1.0, from <https://cran.r-project.org/src/contrib/Archive/mixfdr/>). In all our simulations we assume that the test statistics follow the expected theoretical distribution under the null, and we indicate this to `locfdr` using `nulltype=0` and to `mixfdr` using `theonull=TRUE`. Otherwise all packages were used with default options.

#### 3.1. Effects of the UA

Here we consider the effects of making the UA. To isolate these effects we consider the simplest case, where every observation has the same standard error,  $s_j = 1$ . That is,  $\hat{\beta}_j | \beta_j \sim N(\beta_j, 1)$  and  $\hat{s}_j = s_j = 1$ . In this case the  $z$  scores  $z_j := \hat{\beta}_j / \hat{s}_j = \hat{\beta}_j$ , so modelling the  $z$  scores is the same as modelling the  $\hat{\beta}_j$ . Thus the primary difference among methods in this setting is that `ash` makes the UA and other methods do not.

To briefly summarize the results in this section:

1. The UA can produce quite different results from existing methods.
2. The UA can yield conservative estimates of the proportion of true nulls,  $\pi_0$ , and hence conservative estimates of  $lfdr$  and FDR.
3. The UA yields a procedure that is numerically and statistically stable, and is somewhat robust to deviations from unimodality.

**3.1.1. The UA can produce quite different results from existing methods.** We illustrate the effects of the UA with a simple simulation, with effects  $\beta_j \sim N(0, 1)$  [so with  $s_j = 1$ ,  $\hat{\beta}_j \sim N(0, 2)$ ]. Though no effects are null, there are many  $p$  values near 1 and  $z$  scores near 0 (Figure 1). We used `qvalue`, `locfdr`, `mixfdr`, and `ash` to decompose the  $z$  scores ( $z_j = \hat{\beta}_j$ ), or their corresponding  $p$  values, into null and alternative components. Here we are using the fact that these methods all provide an estimated  $lfdr_j$  for each observation  $j$ , which implies such a decomposition; specifically the average  $lfdr_j$  within each histogram bin estimates the fraction of observations in that bin that come from the null vs the alternative component. The results (Figure 1) illustrate a clear difference between `ash` and existing methods: the existing methods have a “hole” in the alternative  $z$  score distribution near 0, whereas `ash`, due to the UA, has a mode near 0. (Of course the null distribution also has a peak at 0, and the  $lfdr$  under the UA is still smallest for  $z$  scores that are far from zero—i.e. large  $z$  scores remain the “most significant.”)

This qualitative difference among methods is quite general, and also occurs in simulations where most effects are null (e.g. [http://stephenslab.github.io/ash/analysis/referee\\_uaza.html](http://stephenslab.github.io/ash/analysis/referee_uaza.html)). To understand why the alternative distribution of  $z$  scores from `locfdr` and `qvalue` has a hole at zero, note that neither of these methods explicitly models the alternative distribution: instead they simply subtract a null distribution (of  $z$  scores or  $p$  values) from the observed empirical distribution, letting the alternative distribution be defined implicitly, by what remains. In deciding how much null distribution to subtract—that is, in estimating the null proportion,  $\pi_0$ —both methods assume that all  $z$  scores near zero (or, equivalently, all  $p$  values near 1) are null. The consequence of this is that their (implicitly defined) distribution for the alternative  $z$  scores has a “hole” at 0—quite different from our assumption of a mode at zero. (Why `mixfdr` exhibits similar behavior is less clear, since it does explicitly model the alternative distribution; however we believe it may be due to the default choice of penalty term  $\beta$  described in Muralidharan (2010).)

Figure 1 is also helpful in understanding the interacting role of the UA and the penalty term (Supplementary Information, equation S.2.5) that attempts to make  $\pi_0$  as “large as possible” while remaining consistent with the UA. Specifically, consider the panel that shows `ash`’s decomposition of  $z$  scores, and imagine increasing  $\pi_0$  further. This would increase the null component (dark blue) at the expense

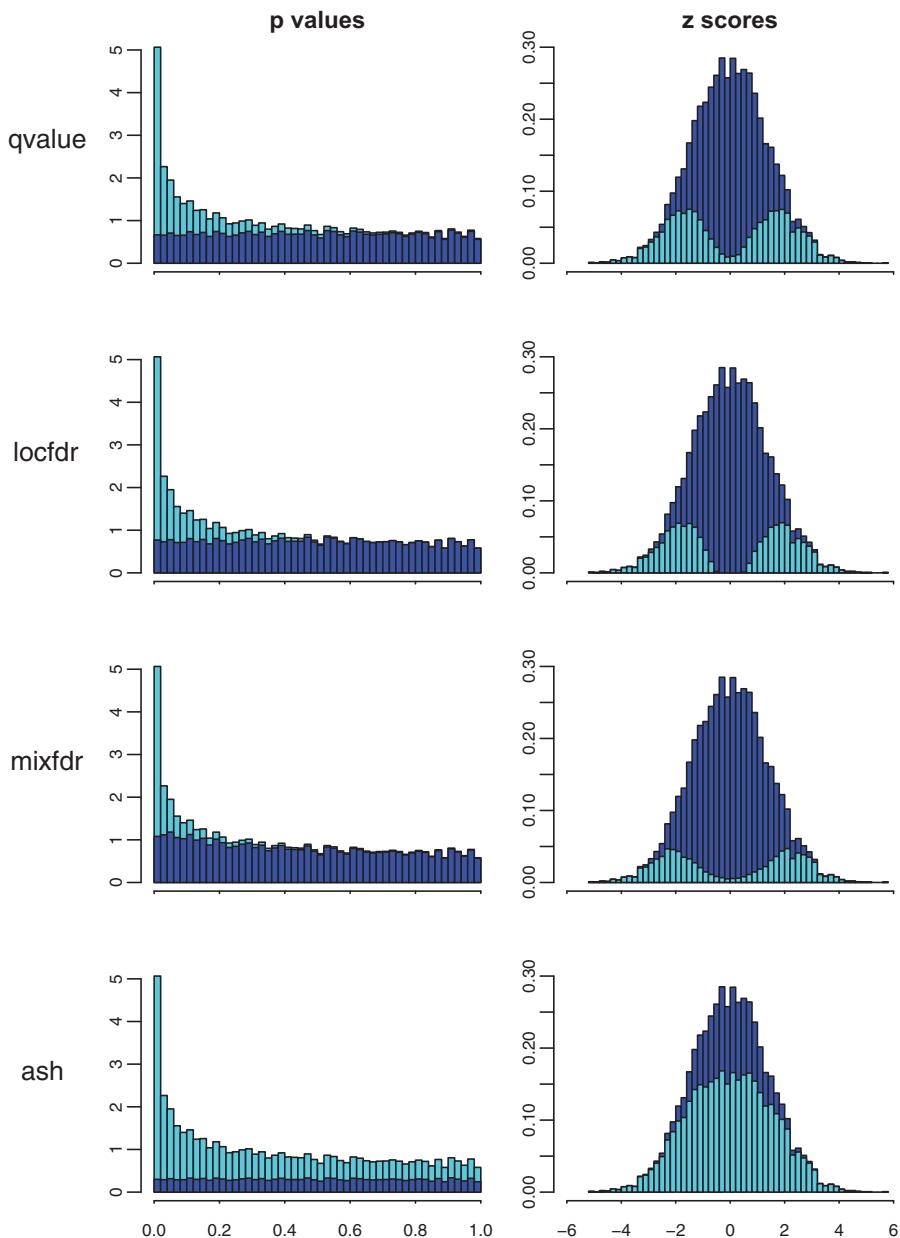


Fig. 1. Illustration that the UA in ash can produce very different results from existing methods. The figure shows, for a single simulated dataset, the way different methods decompose  $p$  values (left) and  $z$  scores (right) into a null component (dark blue) and an alternative component (cyan). In the  $z$  score space the alternative distribution is placed on the bottom to highlight the differences in its shape among methods. The three existing methods (qvalue, locfdr, mixfdr) all produce a “hole” in the alternative  $z$  score distribution around 0. In contrast ash makes the UA—that the effect sizes, and thus the  $z$  scores, have a unimodal distribution about 0—which yields a very different decomposition. (In this case the ash decomposition is closer to the truth: the data were simulated under a model where all of the effects are non-zero, so the “true” decomposition would make everything cyan.)

of the alternative component (light blue). Because the null component is  $N(0, 1)$ , and so is biggest at 0, this would eventually create a “dip” in the light-blue histogram at 0. The role of the penalty term is to push the dark blue component as far as possible, right up to (or, to be conservative, just past) the point where this dip appears. In contrast the existing methods effectively push the dark blue component until the light-blue component *disappears* at 0. See <https://stephens999.shinyapps.io/unimodal/unimodal.Rmd> for an interactive demonstration.

**3.1.2. The UA can produce conservative estimates of  $\pi_0$**  Figure 1 suggests that the UA will produce smaller estimates of  $\pi_0$  than existing methods. Consequently `ash` will estimate smaller  $lfdrs$  and FDRs, and so identify more significant discoveries at a given threshold. This is desirable, provided that these estimates remain conservative: that is, provided that  $\hat{\pi}_0$  does not underestimate the true  $\pi_0$  and  $\hat{lfdr}$  does not underestimate the true  $lfdr$ . The penalty term (Supplementary Information, equation S.2.5) aims to ensure this conservative behavior. To check its effectiveness we performed simulations under various alternative scenarios (i.e. various distributions for the non-zero effects, which we denote  $g_1$ ), and values for  $\pi_0$ . The alternative distributions are shown in Figure 2(a), with details in supplementary Table 1 available at *Biostatistics* online. They range from a “spiky” distribution—where many non-zero  $\beta$  are too close to zero to be reliably detected, making reliable estimation of  $\pi_0$  essentially impossible—to a much flatter distribution, which is a normal distribution with large variance (“big-normal”)—where most non-zero  $\beta$  are easily detected making reliable estimation of  $\pi_0$  easier. We also include one asymmetric distribution (“skew”), and one clearly bimodal distribution (“bimodal”), which, although we view as generally unrealistic, we include to assess robustness of `ash` to deviations from the UA.

For each simulation scenario we simulated 100 independent data sets, each with  $J = 1000$  observations. For each data set we simulated data as follows:

1. Simulate  $\pi_0 \sim U[0, 1]$ .
2. For  $j = 1, \dots, J$ , simulate  $\beta_j \sim \pi_0\delta_0 + (1 - \pi_0)g_1(\cdot)$ .
3. For  $j = 1, \dots, J$ , simulate  $\hat{\beta}_j | \beta_j \sim N(\beta_j, 1)$ .

Figure 2(b) compares estimates of  $\pi_0$  from `qvalue`, `locfdr`, `mixfdr`, and `ash` ( $y$ -axis) with the true values ( $x$ -axis). For `ash` we show results for  $g_1$  modelled as a mixture of normal components (“ash.n”) and as a mixture of symmetric uniform components (“ash.u”). (Results using the asymmetric uniforms, which we refer to as “half-uniforms,” and denote “ash.hu” in subsequent sections, are here generally similar to `ash.u` and omitted to avoid over-cluttering figures.) The results show that `ash` provides the smallest, most accurate, estimates for  $\pi_0$ , while remaining conservative in all scenarios where the UA holds. When the UA does not hold (“bimodal” scenario) the `ash` estimates can be slightly anti-conservative. We view this as a minor concern in practice, since we view such a strong bimodal scenario as unlikely in most applications where FDR methods are used. (In addition, the effects on  $lfsr$  estimates turn out to be relatively modest; see below.)

**3.1.3. The  $lfsr$  is more robust than  $lfdr$**  The results above show that `ash` can improve on existing methods in producing smaller, more accurate, estimates of  $\pi_0$ , which will lead to more accurate estimates of FDR. Nonetheless, in many scenarios `ash` continues to substantially over-estimate  $\pi_0$  (see the “spiky” scenario, for example). This is because these scenarios include an appreciable fraction of “small non-null effects” that are essentially indistinguishable from 0, making accurate estimation of  $\pi_0$  impossible. Put another way, and as is well known,  $\pi_0$  is not identifiable: the data can effectively provide an upper bound on plausible values of  $\pi_0$ , but not a lower bound (because the data cannot rule out that everything is non-null, but with minuscule effects). To obtain conservative behavior we must estimate  $\pi_0$  by this upper bound, which can be substantially larger than the true value.

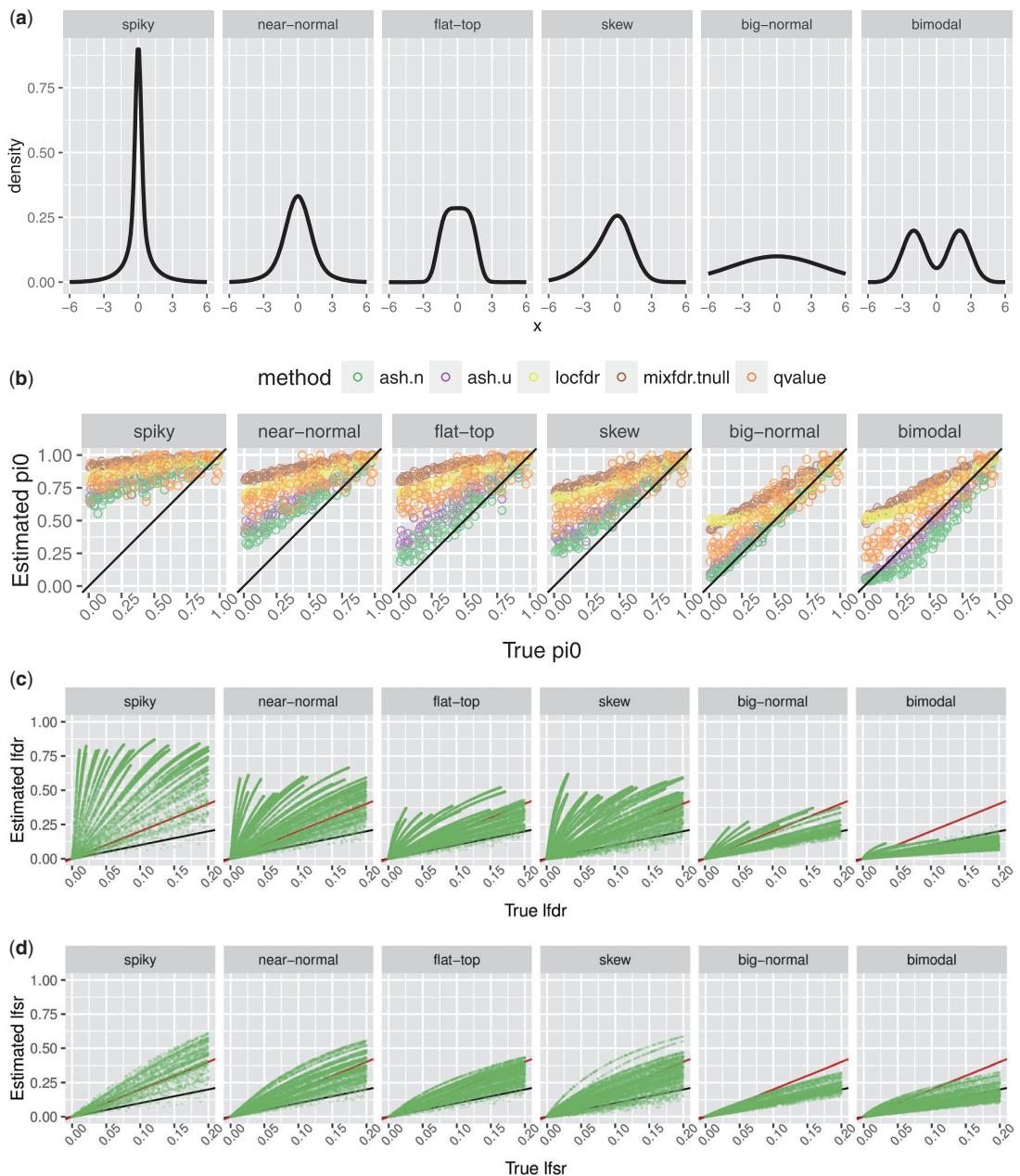


Fig. 2. Results of simulation studies (constant precision  $s_j = 1$ ). (a) Densities of non-zero effects,  $g_1$ , used in simulations. (b) Comparison of true and estimated values of  $\pi_0$ . When the UA holds all methods typically yield conservative (over)-estimates for  $\pi_0$ , with ash being least conservative, and hence most accurate. qvalue is sometimes anti-conservative when  $\pi_0 \approx 1$ . When the UA does not hold ("bimodal" scenario) the ash estimates are slightly anti-conservative. (c) Comparison of true and estimated  $Ifdr$  from ash (ash.n). Black line is  $y = x$  and red line is  $y = 2x$ . Estimates of  $Ifdr$  are conservative when UA holds, due to conservative estimates of  $\pi_0$ . (d) As in (c), but for  $Ifsr$  instead of  $Ifdr$ . Estimates of  $Ifsr$  are consistently less conservative than  $Ifdr$  when UA holds, and also less anti-conservative in bimodal scenario.

Table 1. Empirical coverage for nominal 95% lower credible bounds (all observations)

	Spiky	Near-normal	Flat-top	Skew	Big-normal	Bimodal
ash.n	0.90	0.94	0.95	0.94	0.96	0.96
ash.u	0.87	0.93	0.94	0.93	0.96	0.96
ash.hu	0.88	0.93	0.94	0.94	0.96	0.96

Coverage rates are generally satisfactory, except for the extreme “spiky” scenario. This is due to the penalty term ([Supplementary Information](#), equation S.2.5) which tends to cause over-shrinking towards zero. Removing this penalty term produces coverage rates closer to the nominal levels for uniform and normal methods (Table 2). Removing the penalty in the half-uniform case is not recommended (see online [Supplementary Information](#) for discussion)

Table 2. Empirical coverage for nominal 95% lower credible bounds (significant negative discoveries)

	Spiky	Near-normal	Flat-top	Skew	Big-normal	Bimodal
ash.n	0.94	0.94	0.94	0.86	0.95	0.96
ash.u	0.93	0.93	0.93	0.84	0.95	0.95
ash.hu	0.92	0.92	0.93	0.92	0.95	0.95

Coverage rates are generally satisfactory, except for the uniform-based methods in the spiky and near-normal scenarios, and the normal-based method in the flat-top scenario. These results likely reflect inaccurate estimates of the tails of  $g$  due to a disconnect between the tail of  $g$  and the component distributions in these cases. For example, the uniform methods sometimes substantially underestimate the length of the tail of  $g$  in these long-tailed scenarios, causing over-shrinkage of the tail toward 0

Table 3. Empirical coverage for nominal 95% lower credible bounds (significant positive discoveries)

	Spiky	Near-normal	Flat-top	Skew	Big-normal	Bimodal
ash.n	0.94	0.94	0.94	0.86	0.95	0.96
ash.u	0.93	0.93	0.93	0.84	0.95	0.95
ash.hu	0.92	0.92	0.93	0.92	0.95	0.95

Coverage rates are generally satisfactory, except for the symmetric methods under the asymmetric (“skew”) scenario

Since FDR-related quantities depend quite sensitively on  $\pi_0$ , the consequence of this overestimation of  $\pi_0$  is corresponding overestimation of FDR (and  $lfdr$ , and  $q$  values). To illustrate, Figure 2(c) compares the estimated  $lfdr$  from ash.n with the true value (computed using Bayes rule from the true  $g_1$  and  $\pi_0$ ). As predicted,  $lfdr$  is overestimated, especially in scenarios which involve many non-zero effects that are very near 0 (e.g. the spiky scenario with  $\pi_0$  small) where  $\pi_0$  can be grossly overestimated. Of course other methods will be similarly affected by this: those that more grossly overestimate  $\pi_0$ , will more grossly overestimate  $lfdr$  and FDR/ $q$ -values.

The key point we want to make here is that estimation of  $\pi_0$ , and the accompanying identifiability issues, become substantially less troublesome if we use the local false sign rate  $lfsr$  (2.7), rather than  $lfdr$ , to measure significance. This is because  $lfsr$  is less sensitive to the estimate of  $\pi_0$ . To illustrate, Figure 2(d) compares the estimated  $lfsr$  from ash.n with the true value: although the estimated  $lfsr$  continue to be conservative, overestimating the truth, the overestimation is substantially less pronounced than for the  $lfdr$ , especially for the “spiky” scenario. Further, in the bi-modal scenario, the anti-conservative behavior is less pronounced in  $lfsr$  than  $lfdr$ .

Compared with previous debates, this section advances an additional reason for focusing on the sign of the effect, rather than just testing whether it is 0. In previous debates authors have argued against testing whether an effect is 0 because it is *implausible that effects are exactly 0*. Here we add that *even if one*

*believes that some effects may be exactly zero*, it is still better to focus on the sign, because generally *the data are more informative about that question* and so inferences are more robust to, say, the inevitable mis-estimation of  $\pi_0$ . To provide some intuition, consider an observation with a  $z$  score of 0. The  $lfd$  of this observation can range from 0 (if  $\pi_0 = 0$ ) to 1 (if  $\pi_0 = 1$ ). But, assuming a symmetric  $g$ , the  $lfsr > 0.5$  whatever the value of  $\pi_0$ , because the observation  $z = 0$  says nothing about the sign of the effect. Thus, there are two reasons to use the  $lfsr$  instead of the  $lfd$ : it answers a question that is more generally meaningful (e.g. it applies whether or not zero effects truly exist), and estimation of  $lfsr$  is more robust.

Given that we argue for using  $lfsr$  rather than  $lfd$ , one might ask whether we even need a point mass on zero in our analysis. Indeed, one advantage of the  $lfsr$  is that it makes sense even if no effect is exactly zero. And, if we are prepared to assume that no effects are exactly zero, then removing the point mass yields smaller and more accurate estimates of  $lfsr$  when that assumption is true ([supplementary Figure 2a](#) available at *Biostatistics* online). However, there is “no free lunch”: if in fact some effects are exactly zero then the analysis with no point mass will tend to be anti-conservative, underestimating  $lfsr$  ([supplementary Figure 2b](#) available at *Biostatistics* online). We conclude that if ensuring a “conservative” analysis is important then one should allow for a point mass at 0.

**3.1.4. The UA helps provide reliable estimates of  $g$**  An important advantage of our EB approach based on modelling the effects  $\beta_j$ , rather than  $p$  values or  $z$  scores, is that it can estimate the effects  $\beta_j$ . Specifically, it provides a posterior distribution for each  $\beta_j$ , which can be used to construct interval estimates, etc. Further, because the posterior distribution is, by definition, conditional on the observed data, interval estimates based on posterior distributions are also valid Bayesian inferences for any subset of the effects that have been selected based on the observed data. This kind of “post-selection” validity is much harder to achieve in the frequentist paradigm. In particular the posterior distribution solves the (Bayesian analogue of the) “False Coverage Rate” problem posed by [Benjamini and Yekutieli \(2005\)](#) which [Efron \(2008\)](#) summarizes as follows: “having applied FDR methods to select a set of non-null cases, how can confidence intervals be assigned to the true effect size for each selected case”? [Efron \(2008\)](#) notes the potential for EB approaches to tackle this problem, and [Zhao and Gene Hwang \(2012\)](#) considers in detail the case where the non-null effects are normally distributed.

The ability of the EB approach to provide valid “post-selection” interval estimates is extremely attractive in principle. But its usefulness in practice depends on reliably estimating the distribution  $g$ . Estimating  $g$  is a “deconvolution problem,” which are notoriously difficult in general. Indeed, Efron emphasizes the difficulties of implementing a stable general algorithm, noting in his rejoinder “the effort foundered on practical difficulties involving the perils of deconvolution ... Maybe I am trying to be overly non-parametric ... but it is hard to imagine a generally satisfactory parametric formulation ...” ([Efron, 2008](#) rejoinder, page 46). We argue here that the UA can greatly simplify the deconvolution problem, producing both computationally and statistically stable estimates of  $g$ .

To illustrate, we compare the estimated  $g$  from `ash` (under the UA) with the non-parametric maximum likelihood estimate (NPMLE) for  $g$  (i.e. estimated entirely non-parametrically without the unimodal constraint). The NPMLE is straightforward to compute in R using the `REBayes::GLmix` function ([Koenker and Mizera, 2014](#)). Figure 3 shows results under six different scenarios. The estimated cdf from `ash` is generally closer to the truth, even in the bi-modal scenario. [`ash` tends to systematically overestimate the mass of  $g$  near zero; this can be avoided by removing the penalty term ([Supplementary Information, equation S.2.5](#)); [supplementary Figure 1](#) available at *Biostatistics* online.] Furthermore, the estimate from `ash` is also substantially more “regular” than the NPMLE, which has several almost-vertical segments indicative of a concentration of density in the estimated  $g$  at those locations. Indeed the NPMLE is a discrete distribution ([Koenker and Mizera, 2014](#)), so this kind of concentration

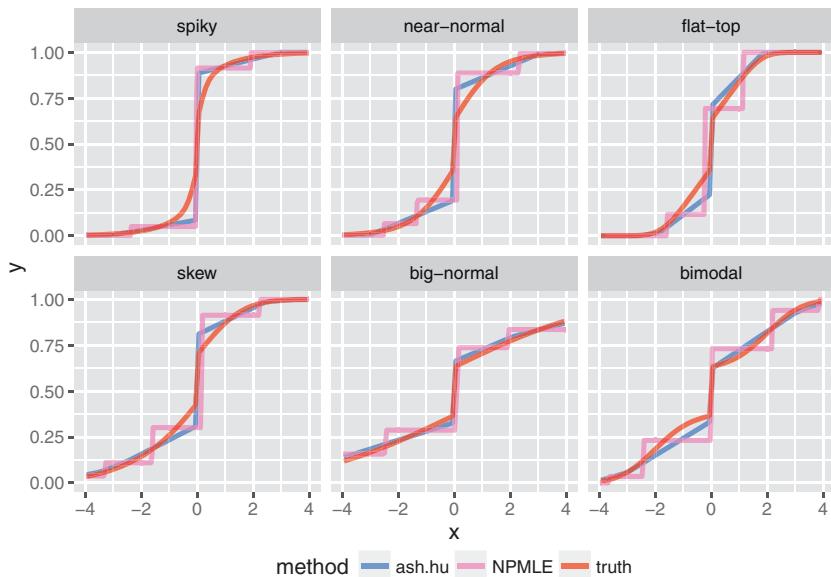


Fig. 3. Comparison of estimated cdfs from `ash` and the NPMLE. Different `ash` methods perform similarly, so only `ash.hu` is shown for clarity. Each panel shows results for a single example data set, one for each scenario in Figure 2(a). The results illustrate how the UA made by `ash` regularizes the estimated cdfs compared with the NPMLE.

will always occur. The UA prevents this concentration, effectively regularizing the estimated  $g$ . While the UA is not the only way to achieve this (e.g. Efron, 2016), we view it as attractive and widely applicable.

**3.1.5. Calibration of posterior intervals** To quantify the effects of errors in estimates of  $g$  we examine the calibration of the resulting posterior distributions (averaged over 100 simulations in each Scenario). Specifically we examine the empirical coverage of nominal lower 95% credible bounds for (i) all observations; (ii) significant negative discoveries; (iii) significant positive discoveries. We examine only lower bounds because the results for upper bounds follow by symmetry (except for the one asymmetric scenario). We separately examine positive and negative discoveries because the lower bound plays a different role in each case: for negative discoveries the lower bound is typically large and negative and limits how big (in absolute value) the effect could be; for positive discoveries the lower bound is positive, and limits how small (in absolute value) the effect could be. Intuitively, the lower bound for negative discoveries depends on the accuracy of  $g$  in its tail, whereas for positive discoveries it is more dependent on the accuracy of  $g$  in the center.

The results are shown in Tables 1–3. Most of the empirical coverage rates are in the range 0.92–0.96 for nominal coverage of 0.95, which we view as adequate for practical applications. The strongest deviations from nominal rates are noted and discussed in the table captions. One general issue is that the methods based on mixtures of uniform distributions often slightly curtail the tail of  $g$ , causing the probability of very large outlying effects to be understated; see also <http://stephenslab.github.io/ash/analysis/efron.fcr.html>.

### 3.2. Differing measurement precision across units

We turn now to the second important component of our work: allowing for varying measurement precision across units. The key to this is the use of a likelihood, (2.4) (or, more generally, in

Supplementary Information, equation (S.1.2)), that explicitly incorporates the measurement precision (standard error) of each  $\hat{\beta}_j$ .

To illustrate, we conduct a simulation where half the measurements are quite precise (standard error  $s_j = 1$ ), and the other half are very poor ( $s_j = 10$ ). In both cases, we assume that half the effects are null and the other half are normally distributed with standard deviation 1:

$$p(\beta) = 0.5\delta_0(\beta) + 0.5N(\beta; 0, 1). \quad (3.1)$$

In this setting, the poor-precision measurements ( $s_j = 10$ ) tell us very little, and any sane analysis should effectively ignore them. However, this is not the case in standard FDR-type analyses (Figure 4). This is because the poor measurements produce  $p$  values that are approximately uniform [Figure 4(a)], which, when combined with the good-precision measurements, dilute the overall signal (e.g. they reduce the density of  $p$  values near 0). This is reflected in the results of FDR methods like `qvalue` and `locfdr`: the estimated error rates ( $q$ -values, or  $lfdr$  values) for the good-precision observations increase when the low-precision observations are included in the analysis [Figure 4(b)]. In contrast, the results from `ash` for the good-precision observations are unaffected by including the low-precision observations in the analysis [Figure 4(b)].

Another consequence of incorporating measurement precision into the likelihood is that `ash` can re-order the significance of observations compared with the original  $p$  values or  $z$  scores. Effectively `ash` downweights the poor-precision observations by assigning them a higher  $lfsr$  than good precision measurements that have the same  $p$  value [Figure 4(c)]. The intuition is that, due to their poor precision, these measurements contain very little information about the sign of the effects (or indeed any other aspect of the effects), and so the  $lfsr$  for these poor-precision measurements is always high; see Guan and Stephens (2008) for related discussion. Here this re-ordering results in improved performance: `ash` identifies more true positive effects at a given level of false positives [Figure 4(d)].

**3.2.1. Dependence of  $\beta_j$  on  $s_j$**  Although downweighting low-precision observations may seem intuitive, we must now confess that the issues are more subtle than our treatment above suggests. Specifically, it turns out that the downweighting behavior depends on an assumption that we have made up to now, that  $\beta_j$  is independent of  $s_j$  (2.2). In practice this assumption may not hold. For example, in gene expression studies, genes with higher biological variance may tend to exhibit larger effects  $\beta_j$  (because they are less constrained). These genes will also tend to have larger  $s_j$ , inducing a dependence between  $\beta_j$  and  $s_j$ .

Motivated by this, we generalize the prior (2.2) to

$$\frac{\beta_j}{\hat{s}_j^\alpha} \mid \hat{s}_j \sim g(\cdot; \pi), \quad (3.2)$$

where  $\alpha$  is to be estimated or specified. Setting  $\alpha = 0$  yields (2.2). Setting  $\alpha > 0$  implies that the effects with larger standard error tend to be larger (in absolute value). Fitting this model for any  $\alpha$  is straightforward using the same methods as for  $\alpha = 0$  (see [supplementary material](#) available at *Biostatistics* online).

The case  $\alpha = 1$  in (3.2) is of special interest because it corresponds most closely to existing methods. Specifically, it can be shown that with  $\alpha = 1$  the  $lfsr_j$  values from `ash.n` are monotonic in the  $p$  values: effects with smaller  $p$  values have smaller  $lfsr_j$ . This result generalizes a result in Wakefield (2009), who referred to a prior that produces the same ranking as the  $p$  values as a “ $p$  value prior.” Of course, if the  $lfsr_j$  are monotonic in the  $p$  values then the downweighting and reordering of significance illustrated in Figure 4 will not occur. The intuition is that under  $\alpha = 1$  the poor precision observations have larger effects sizes, and consequently the same power as the high-precision observations—under these conditions the poor precision observations are not “contaminating” the high precision observations, and so downweighting

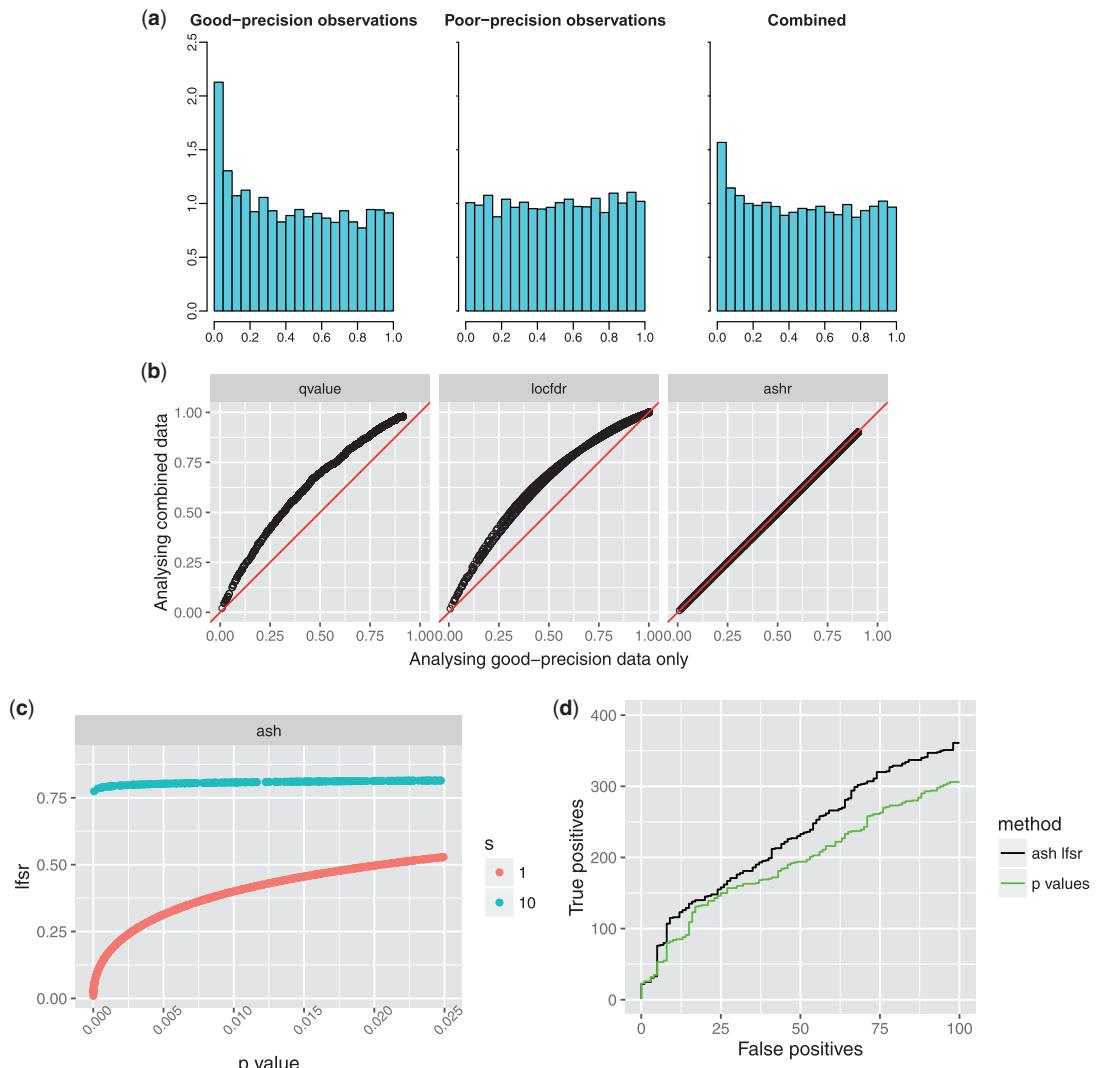


Fig. 4. Simulations showing how, with existing methods, but not `ash`, poor-precision observations can contaminate signal from good-precision observations. (a) Density histograms of  $p$  values for good-precision, poor-precision, and combined observations. The combined data show less signal than the good-precision data, due to the contamination effect of the poor-precision measurements. (b) Results of different methods applied to good-precision observations only (x-axis) and combined data (y-axis). Each point shows the “significance” ( $q$  values from `qvalue`;  $lfsr$  for `locfdr`;  $lfsr$  for `ash`) of a good-precision observation under the two different analyses. For existing methods including the poor-precision observations reduces significance of good-precision observations, whereas for `ash` the poor-precision observations have little effect (because they have a very flat likelihood). (c) The relationship between  $lfsr$  and  $p$ -value is different for good-precision ( $s = 1$ ) and low-precision ( $s = 10$ ) measurements: `ash` assigns the low-precision measurements a higher  $lfsr$ , effectively downweighting them. (d) Trade-off between true positives (y) vs false positives (x) as the significance threshold ( $lfsr$  or  $p$  value) is varied. By downweighting the low-precision observations `ash` re-orders the significance of observations, producing more true positives at a given number of false positives. It is important to note that this behaviour of `ash` depends on choice of  $\alpha$ . See Section 3.2.1 for discussion.

them is unnecessary. Thus running `aash` with  $\alpha = 1$  will produce the same significance ranking as existing methods. Nonetheless, it is not equivalent to them, and indeed still has the benefits outlined previously: due to the UA `aash` can produce less conservative estimates of  $\pi_0$  and  $lfdr_j$ ; and because `aash` models the  $\beta_j$  it can produce interval estimates.

As an aside, we note that with  $\alpha = 1$  the `aash` estimates of both  $g$  and  $lfsr_j$  depend on the pairs  $(\hat{\beta}_j, \hat{s}_j)$  only through the  $z$  scores  $z_j = \hat{\beta}_j/\hat{s}_j$  [though interval estimates for  $\beta_j$  still depend on  $(\hat{\beta}_j, \hat{s}_j)$ ]. This means that `aash` can be run (with  $\alpha = 1$ ) in settings where the  $z_j$  are available but the  $(\hat{\beta}_j, \hat{s}_j)$  are not. It also opens the intriguing possibility of running `aash` on  $p$  values obtained from any test (e.g. a permutation test), by first converting each  $p$  value to a corresponding  $z$  score. However, the meaning of and motivation for the UA may be unclear in such settings, and caution seems warranted before proceeding along these lines.

In practice, appropriate choice of  $\alpha$  will depend on the actual relationship between  $\beta_j$  and  $s_j$ , which will be dataset-specific. Further, although we have focused on the special cases  $\alpha = 0$  and  $\alpha = 1$ , there seems no strong reason to expect that either will necessarily be optimal in practice. Following the logic of the EB approach we suggest selecting  $\alpha$  by maximum likelihood, which is implemented in our software using a simple 1-d grid search (implementation due to C. Dai).

#### 4. DISCUSSION

We have presented an Empirical Bayes approach to large-scale multiple testing that emphasizes two ideas. First, we emphasize the potential benefits of using two numbers ( $\hat{\beta}_j$  and its standard error) rather than just one number (a  $p$  value or  $z$  score) to summarize the information on each test. While requiring two numbers is slightly more onerous than requiring one, in many settings these numbers are easily available and if so we argue it makes sense to use them. Second, we note the potential benefits—both statistical and computational—of assuming that the effects come from a unimodal distribution, and provide flexible implementations for performing inference under this assumption. We also introduce the “false sign rate” as an alternative measure of error to the FDR, and illustrate its improved robustness to errors in model fit, particularly mis-estimation of the proportion of null tests,  $\pi_0$ .

Multiple testing is often referred to as a “problem” or a “burden.” In our opinion, EB approaches turn this idea on its head, treating multiple testing as an *opportunity*: specifically, an opportunity to learn about the prior distributions, and other modelling assumptions, to improve inference and make informed decisions about significance (see also Greenland and Robins, 1991). This view also emphasizes that, what matters in multiple testing settings is *not* the number of tests, but the *results* of the tests. Indeed, the FDR at a given fixed threshold does not depend on the number of tests: as the number of tests increases, both the true positives and false positives increase linearly, and the FDR remains the same. [If this intuitive argument does not convince, see Storey (2003), and note that the FDR at a given  $p$  value threshold does not depend on the number of tests  $m$ .] Conversely, the FDR *does* depend on the overall distribution of effects, and particularly on  $\pi_0$ , for example. The EB approach captures this dependence in an intuitive way: if there are lots of strong signals then we infer  $\pi_0$  to be small, and the estimated FDR (or  $lfdr$ , or  $lfsr$ ) at a given threshold may be low, even if a large number of tests were performed; and conversely if there are no strong signals then we infer  $\pi_0$  to be large and the FDR at the same threshold may be high, even if relatively few tests were performed. More generally, overall signal strength is reflected in the estimated  $g$ , which in turn influences the estimated FDR.

Two important practical issues that we have not addressed here are correlations among tests, and the potential for deviations from the theoretical null distributions of test statistics. These two issues are connected: specifically, unmeasured confounding factors can cause both correlations among tests and deviations from the theoretical null (Efron, 2007; Leek and Storey, 2007). And although there are certainly other factors that could cause dependence among tests, unmeasured confounders are perhaps the most

worrisome in practice because they can induce strong correlations among large numbers of tests and profoundly impact results, ultimately resulting in too many hypotheses being rejected and a failure to control FDR. We are acutely aware that, because our method is less conservative than existing methods, it may unwittingly exacerbate these issues if they are not adequately dealt with. Approaches to deal with unmeasured confounders can be largely divided into two types: those that simply attempt to correct for the resulting inflation of test statistics (Devlin and Roeder, 1999; Efron, 2004), and those that attempt to infer confounders using clustering, principal components analysis, or factor models (Pritchard *and others*, 2000; Price *and others*, 2006; Leek and Storey, 2007; Gagnon-Bartsch and Speed, 2012), and then correct for them in computation of the test statistics (in our case,  $\hat{\beta}, \hat{s}$ ). When these latter approaches are viable they provide perhaps the most satisfactory solution, and are certainly a good fit for our framework. Alternatively, our methods could be modified to allow for test statistic inflation, an idea that may be worth pursuing in future work.

Another important practical issue is the challenge of small sample sizes. For example, in genomics applications researchers sometimes attempt to identify differences between two conditions based on only a handful of samples in each. In such settings the normal likelihood approximation (2.4) will be inadequate. And, although the  $t$  likelihood (Supplementary Information, equation (S.1.2)) partially addresses this issue, it is also, it turns out, not entirely satisfactory. The root of the problem is that, with small sample sizes, raw estimated standard errors  $\hat{s}_j$  can be horribly variable. In genomics it is routine to address this issue by applying EB methods (Smyth, 2004) to “moderate” (i.e. shrink) variance estimates, before computing  $p$  values from “moderated” test statistics. We are currently investigating how our methods should incorporate such “moderated” variance estimates to make it applicable to small sample settings.

Our approach involves compromises between flexibility, generality, and simplicity on the one hand, and statistical efficiency and principle on the other. For example, in using an EB approach that uses a point estimate for  $g$ , rather than a fully Bayes approach that accounts for uncertainty in  $g$ , we have opted for simplicity over statistical principle. And in summarizing every test by two numbers and making a normal or  $t$  approximation to the likelihood, we have aimed to produce generic methods that can be applied whenever such summary data are available—just as qvalue can be applied to any set of  $p$  values, for example—although possibly at the expense of statistical efficiency compared with developing multiple tailored approaches based on context-specific likelihoods. Any attempt to produce generic methods will involve compromise between generality and efficiency. In genomics, many analyses—not only FDR-based analyses—involve first computing a series of  $p$  values before subjecting them to some further downstream analysis. An important message here is that working with two numbers ( $\hat{\beta}_j, \hat{s}_j$ ), rather than one ( $p_j$  or  $z_j$ ), can yield substantial gains in functionality (e.g. estimating effect sizes, as well as testing; accounting for variations in measurement precision across units) while losing only a little in generality. We hope that our work will encourage development of methods that exploit this idea in other contexts.

#### SUPPLEMENTARY MATERIAL

Supplementary material is available at <http://biostatistics.oxfordjournals.org>.

#### ACKNOWLEDGMENTS

I thank J. Lafferty for pointing out the convexity of the likelihood function, which lead to an improved implementation with interior point methods. Statistical analyses were conducted in the R programming language (R Core Team, 2012), Figures produced using the ggplot2 package (Wickham, 2009), and text prepared using L<sup>A</sup>T<sub>E</sub>X. Development of the methods in this paper was greatly enhanced by the use of the knitr package (Xie, 2013) within the RStudio GUI, and git and github. The ashR R package is

available from <http://github.com/stephens999/ashr> and includes contributions from Chaoxing (Rick) Dai, Mengyin Lu, Nan Xiao, and Tian Sen. *Conflict of Interest:* None declared.

#### FUNDING

National Institutes of Health (HG02585) and a grant from the Gordon and Betty Moore Foundation (GBMF #4559).

#### 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)* **57**, 289–300.
- Benjamini, Y. and Yekutieli, D. (2005). False discovery rate-adjusted multiple confidence intervals for selected parameters. *Journal of the American Statistical Association* **100**, 71–81.
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge, UK: Cambridge University Press.
- Carvalho, C. M., Polson, N. G. and Scott, J. G. (2010). The horseshoe estimator for sparse signals. *Biometrika* **97**, asq017.
- Cordy, C. B. and Thomas, D. R. (1997). Deconvolution of a distribution function. *Journal of the American Statistical Association* **92**, 1459–1465.
- Datta, S. and Datta, S. (2005). Empirical Bayes screening of many p-values with applications to microarray studies. *Bioinformatics* **21**, 1987–1994.
- Devlin, B. and Roeder, K. (1999). Genomic control for association studies. *Biometrics* **55**, 997–1004.
- Donoho, D. and Johnstone, I. (1995). Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association* **90**, 1200–1224.
- Efron, B. (1993). Bayes and likelihood calculations from confidence intervals. *Biometrika* **80**, 3–26.
- Efron, B. (2004). Large-scale simultaneous hypothesis testing: the choice of a null hypothesis. *Journal of the American Statistical Association* **99**, 96–104.
- Efron, B. (2007). Correlation and large-scale simultaneous significance testing. *Journal of the American Statistical Association* **102**, 93–107.
- Efron, B. (2008). Microarrays, empirical Bayes and the two-groups model. *Statistical Science* **23**, 1–22.
- Efron, B. (2010). *Large-scale Inference: Empirical Bayes Methods for Estimation, Testing, and Prediction*, Volume 1. Cambridge, UK: Cambridge University Press ([http://statweb.stanford.edu/~ckirby/brad/LSI/monograph\\_CUP.pdf](http://statweb.stanford.edu/~ckirby/brad/LSI/monograph_CUP.pdf)).
- Efron, B. (2016). Empirical Bayes deconvolution estimates. *Biometrika* **103**, 1–20.
- Efron, B. (2003). Robbins, empirical Bayes and microarrays. *The Annals of Statistics* **31**, 366–378.
- Efron, B. and Tibshirani, R. (2002). Empirical Bayes methods and false discovery rates for microarrays. *Genetic Epidemiology* **23**, 70–86.
- Efron, B., Tibshirani, R., Storey, J. D. and Tusher, V. (2001). Empirical Bayes analysis of a microarray experiment. *Journal of the American Statistical Association* **96**, 1151–1160.
- Erbe, M., Hayes, B. J., Matukumalli, L. K., Goswami, S., Bowman, P. J., Reich, C. M., Mason, B. A. and Goddard, M. E. (2012). Improving accuracy of genomic predictions within and between dairy cattle breeds with imputed high-density single nucleotide polymorphism panels. *Journal of Dairy Science* **95**, 4114–29.
- Gagnon-Bartsch, J. A. and Speed, T. P. (2012). Using control genes to correct for unwanted variation in microarray data. *Biostatistics* **13**, 539–552.

- Gelman, A. and Tuerlinckx, F. (2000). Type s error rates for classical and Bayesian single and multiple comparison procedures. *Computational Statistics* **15**, 373–390.
- Gelman, A., Hill, J. and Yajima, M. (2012). Why we (usually) don't have to worry about multiple comparisons. *Journal of Research on Educational Effectiveness* **5**, 189–211.
- Greenland, S. and Robins, J. M. (1991). Empirical-Bayes adjustments for multiple comparisons are sometimes useful. *Epidemiology* **2**, 244–251.
- Guan, Y. and Stephens, M. (2008). Practical issues in imputation-based association mapping. *PLoS Genetics* **4**.
- Jiang, W. and Zhang, C.-H. (2009). General maximum likelihood empirical Bayes estimation of normal means. *The Annals of Statistics* **37**, 1647–1684.
- Johnstone, I. M. and Silverman, B. W. (2004). Needles and straw in haystacks: empirical Bayes estimates of possibly sparse sequences. *Annals of Statistics* **32**, 1594–1649.
- Kendziora, C. M., Newton, M. A., Lan, H. and Gould, M. N. (2003). On parametric empirical Bayes methods for comparing multiple groups using replicated gene expression profiles. *Statistics in Medicine* **22**, 3899–3914.
- Koenker, R. and Mizera, I. (2013). Convex optimization in R. *Journal of Statistical Software* **60**, 1–23.
- Koenker, R. and Mizera, I. (2014). Convex optimization, shape constraints, compound decisions, and empirical Bayes rules. *Journal of the American Statistical Association* **109**, 674–685.
- Leek, J. T. and Storey, J. D. (2007). Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genetics* **3**, 1724–35.
- Lu, M. and Stephens, M. (2016). Variance adaptive shrinkage (vash): flexible empirical Bayes estimation of variances. *bioRxiv*, 048660, In press.
- Muralidharan, O. (2010). An empirical Bayes mixture method for effect size and false discovery rate estimation. *The Annals of Applied Statistics* **4**, 422–438.
- Newton, M. A., Noueiry, A., Sarkar, D. and Ahlquist, P. (2004). Detecting differential gene expression with a semiparametric hierarchical mixture method. *Biostatistics* **5**, 155–176.
- Ploner, A., Calza, S., Gusnanto, A. and Pawitan, Y. (2006). Multidimensional local false discovery rate for microarray studies. *Bioinformatics* **22**, 556–565.
- Price, A. L., Patterson, N. J., Plenge, R. M., Weinblatt, M. E., Shadick, N. A. and Reich, D. (2006). Principal components analysis corrects for stratification in genome-wide association studies. *Nature Genetics* **38**, 904–9.
- Pritchard, J. K., Stephens, M., Rosenberg, N. A. and Donnelly, P. (2000). Association mapping in structured populations. *American Journal of Human Genetics* **67**, 170–181.
- R Core Team. (2012). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org/>, Accessed June 3, 2013.
- Sarkar, A., Mallick, B. K., Staudenmayer, J., Pati, D. and Carroll, R. J. (2014). Bayesian semiparametric density deconvolution in the presence of conditionally heteroscedastic measurement errors. *Journal of Computational and Graphical Statistics* **23**, 1101–1125.
- Smyth, G. K. (2004). Linear models and empirical Bayes methods for assessing differential expression in microarray experiments. *Statistical Applications in Genetics and Molecular Biology* **3**, Article3.
- Storey, J. D. (2003). The positive false discovery rate: a Bayesian interpretation and the q-value. *The Annals of Statistics* **31**, 2013–2035.
- Tukey, J. W. (1962). The future of data analysis. *The Annals of Mathematical Statistics* **33**, 1–67.

- Tukey, J. W. (1991). The philosophy of multiple comparisons. *Statistical Science* **6**, 100–116.
- Wakefield, J. (2009). Bayes factors for genome-wide association studies: comparison with p-values. *Genetic Epidemiology* **33**, 79–86.
- Wickham, H. (2009). *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer.
- Xie, Y. (2013). *Dynamic Documents with R and Knitr*, Volume 29. Boca Raton, FL: CRC Press.
- Xie, X., Kou, S. C. and Brown, L. D. (2012). Sure estimates for a heteroscedastic hierarchical model. *Journal of the American Statistical Association* **107**(500), 1465–1479.
- Xing, Z. and Stephens, M. (2016). Smoothing via adaptive shrinkage (smash): denoising poisson and heteroskedastic Gaussian signals. arXiv preprint arXiv:1605.07787.
- Zhao, Z. and Gene Hwang, J. T. (2012). Empirical Bayes false coverage rate controlling confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **74**, 871–891.

[Received June 14, 2016; revised June 14, 2016; accepted for publication June 21, 2016]