

# Genome-wide Significance Levels and Weighted Hypothesis Testing

Kathryn Roeder and Larry Wasserman \*

Department of Statistics  
Carnegie Mellon University  
Pittsburgh, PA

Address for correspondence and reprints:

Kathryn Roeder, Department of Statistics, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213. E-mail: roeder@stat.cmu.edu

---

\*Research supported by National Institute of Mental Health grant MH057881 and NSF Grant AST 0434343. The authors thank Jamie Robins for helping us to clarify several issues.

## Abstract

Genetic investigations often involve the testing of vast numbers of related hypotheses simultaneously. To control the overall error rate, a substantial penalty is required, making it difficult to detect signals of moderate strength. To improve the power in this setting a number of authors have considered using weighted p-values, with the motivation often based upon the scientific plausibility of the hypotheses. We review this literature, derive optimal weights and show that the power is remarkably robust to misspecification of these weights. We consider two methods for choosing weights in practice. The first, external weighting, is based on prior information. The second, estimated weighting, uses the data to choose weights.

Key Words: Bonferroni correction, Multiple testing, Weighted p-values.

## 1 Introduction

Testing for association between genetic variation and a complex disease typically requires scanning hundreds of thousands of genetic polymorphisms. In a multiple testing situation, such as a genome-wide association study (GWAS), the null hypothesis is rejected for any test that achieves a p-value less than a predetermined threshold (usually on the order of  $10^{-8}$ ). Data from these investigations has renewed interest in the multiple testing problem. The introduction of the false discovery rate and a procedure to control it by Benjamini and Hochberg (1995) inspired hope that this would be an effective way to control error while increasing power (Storey and Tibshirani, 2003; Sabatti et al., 2003). To further bolster power, recent statistical methods have been proposed that up-weight and down-weight hypotheses, based on prior likelihood of association with the phenotype (Genovese et al., 2006; Roeder et al., 2006, 2007; Wang et al., 2007). Such prior information is often available in practice.

Weighted procedures multiply the threshold by the weight  $w$ , for each test, raising the threshold when  $w > 1$  and lowering it if  $w < 1$ . To control the overall rate of false positives, a budget must be imposed on the weighting scheme, so that the average weight is one. If the weights are informative, the procedure improves power substantially, but, if the weights are uninformative, the loss in power is usually small. Surprisingly, aside from this budget requirement, any set of non-negative weights

is valid (Genovese et al., 2006). While desirable in some respects, this flexibility makes it difficult to select weights for a particular analysis.

The first such weighting scheme appears to be Holm (1979). Related ideas can be found in Benjamini and Hochberg (1997), Chen et al. (2000), Genovese et al. (2006), Kropf et al. (2004), Rosenthal and Rubin (1983), Schuster et al. (2004), Westfall and Krishen (2001), Westfall et al. (2004), Blanchard and Roquain (2008b), and Roquain and van de Wiel (2008) among others. Several of these approaches use data dependent weights and yet maintain familywise error control. There are, of course, other ways to improve power aside from weighting. Some notable recent approaches include Rubin et al. (2006), Storey (2007), Donoho and Jin (2004), Signoravitch (2006), Westfall et al. (1998), Westfall and Soper (2001), Efron (2007) and Sun and Cai (2007). Of these, our approach is closest to Rubin et al. (2006).

In some cases, the optimal weights can be estimated from the data. An approach developed by Westfall et al. (2004) utilizes quadratic forms to construct such weights; however, this approach assumes the individual measurements are normally distributed. This approach is suited to applications such as microarray data for which the observations are approximately normally distributed. We are interested in applications such as tests for genetic association. In this setting the individual observations are discrete, but the test statistics are approximately normally distributed.

In general,  $p$ -value weighting raises several important questions. How much power can we gain if we guess well in the weight assignment? How much power can we lose if we guess poorly? In this paper we show that the optimal weights have a simple parametric form and we investigate various approaches for estimating these weights. We also show the power is very robust to misspecification of the weights. In particular, in Section 3 we show that (i) sparse weights (few large weights and minimum weight close to 1) lead to huge power gains for well specified weights, but minute power loss for poorly specified weights; and (ii) in the non-sparse case, under weak conditions, the worst case power for poorly specified weights is typically better than the power obtained using equal weights.

We consider two methods for choosing the weights: (i) external weights, where prior information (based on scientific knowledge or prior data) singles out specific hypotheses (Section 4) and

(ii) estimated weights where the data are used to construct weights (Section 5). External weights are prone to bias while estimated weights are prone to variability. The two robustness properties reduce concerns about bias and variance.

To motivate this work consider an example (Figure 1) of external weighting that arises in genetic epidemiology. To identify variants of genes that induce greater susceptibility to disease, two types of studies (linkage and association) are often performed. Whole genome linkage analysis has been conducted for most major diseases. These data can be summarized by a linkage trace, a smooth stochastic process  $\{Z(s) : s \in [0, L]\}$  where each  $s$  corresponds to a location on the genome. At points that correspond to a variant of a gene of interest, the mean of the process  $\mu(s) = E(Z(s))$  is a large positive value; however, due to extensive spatial correlation in the process,  $\mu(s)$  is also non-zero in the vicinity of the variant. Tests for association between genetic polymorphisms and disease status for each of many genetic markers across the genome are also of interest. Like linkage analysis, the association statistics  $\{T_j : j = 1, \dots, m\}$  map to spatial locations  $\{s_j : j = 1, \dots, m\}$  on the genome. The number of tests  $m$  can be large, on the order of 1,000,000. Until recently, whole genome association analysis was prohibitively expensive, but technological advances have now made such studies feasible. Due to the multiple testing correction, it is difficult to achieve sufficient power to obtain definitive results in these studies. The linkage trace provides one obvious source of information from which the weights can be constructed; see Section 6 for further elaboration. Unlike linkage analysis, however, the spatial correlation in association tests is weak. For this reason, other choices such as genetic pathways could offer a more promising source for weights in the future.

## 2 Background

### 2.1 Multiple Testing

Consider a multiple testing situation in which  $m$  tests are being performed. Suppose  $m_0$  of the null hypotheses are true and  $m_1 = m - m_0$  null hypotheses are false. We can categorize the  $m$  tests as in Table 1. In this notation  $F$  is the number of false positives. To control the familywise error

	$H_0$ Rejected	$H_0$ Not Rejected	Total
$H_0$ True	$F$	$m_0 - F$	$m_0$
$H_0$ False	$T$	$m_1 - T$	$m_1$
Total	$S$	$m - S$	$m$

**Table 1.**  $2 \times 2$  classification of  $m$  hypothesis tests.

rate it is traditional to bound  $P(F > 0)$  at  $\alpha$ . When the tests are independent, the simplest way to control this probability is to reject only those tests for which the p-value is less than  $\alpha/m$ ; this is called the Bonferroni procedure.

In 1995 Benjamini and Hochberg (BH) introduced a new approach to multiple hypothesis testing that controls the false discovery rate (FDR), defined as the expected fraction of false rejections among those hypotheses rejected. Let  $P_{(1)} < \dots < P_{(m)}$  be the ordered  $p$ -values from  $m$  hypothesis tests, with  $P_{(0)} \equiv 0$ . Then, the BH procedure rejects any null hypothesis for which  $P \leq T$  with

$$T = \max \left\{ P_{(i)} : P_{(i)} \leq \frac{\alpha i}{m} \right\}.$$

This quantity is of more scientific relevance than the overall type I error rate in GWAS. Also, the procedure is more powerful than the Bonferroni method. Adaptive variants of the procedure can increase power further at little additional computational expense; see Benjamini et al. (2006) and Storey (2002).

BH controls the false discovery rate at level  $\alpha m_0/m$ , where  $m_0$  is the number of true null hypotheses. With certain dependence assumptions on the p-values, this is true regardless of how many nulls are true and regardless of the distribution of the p-values under the alternatives (Benjamini and Yekutieli, 2001; Blanchard and Roquain, 2008a; Sarkar, 2002). Under some distributional assumptions, Genovese and Wasserman (2002) show that, asymptotically, the BH method corresponds to rejecting all p-values less than a particular p-value threshold  $u^*$ . Specifically  $u^*$  is the solution to the equation  $H(u) = \beta u$  and  $\beta = (\frac{1}{\alpha} - A_0)/(1 - A_0)$ , where  $A_0 = m_0/m$  and  $H$  is the (common) distribution of the p-value under the alternative. The key result is that  $\alpha/m \leq u^* \leq \alpha$  which shows that the BH method is intermediate between Bonferroni (corresponding to  $\alpha/m$ ) and uncorrected testing (corresponding to  $\alpha$ ). If  $A_0$  is close to 0, however, as it usually is in GWA,

then  $\beta$  is a very large quantity the power of the FDR is not much improved over the Bonferroni procedure.

The power of the BH method can be improved with adaptations. Blanchard and Roquain (2008b) have given numerical comparisons of different adaptive procedures under dependence. Romano et al. (2008) have considered improving the adaptive procedure of Benjamini et al. (2006) using the bootstrap. Sarkar and Heller (2008) have noted that the adaptive procedure of Benjamini et al. may not perform well compared to Storey's (2002) procedure for certain parameter choices.

## 2.2 Weighted Multiple Testing

We are given hypotheses  $H = (H_1, \dots, H_m)$  and standardized test statistics  $T = (T_1, \dots, T_m)$  where  $T_j \sim N(\xi_j, 1)$ . Likewise  $T_j^2 \sim \chi_1^2(\xi_j^2)$ . For a two-sided hypothesis,  $H_j = 1$  if  $\xi_j \neq 0$  and  $H_j = 0$  otherwise. For the sake of parsimony, unless otherwise noted, results will be stated for a one-sided test where  $H_j = 1$  if  $\xi_j > 0$  although the results extend easily to the two-sided case. Let  $\theta = (\xi_1, \dots, \xi_m)$  denote the vector of means.

The p-values associated with the tests are  $P = (P_1, \dots, P_m)$  where  $P_j = \overline{\Phi}(T_j)$ ,  $\overline{\Phi} = 1 - \Phi$  and  $\Phi$  denotes the standard Normal CDF. Let  $P_{(1)} \leq \dots \leq P_{(m)}$  denote the sorted p-values and let  $T_{(1)} \geq \dots \geq T_{(m)}$  denote the sorted test statistics.

A *rejection set*  $\mathcal{R}$  is a subset of  $\{1, \dots, m\}$ . Say that  $\mathcal{R}$  *controls familywise error at level*  $\alpha$  if  $\mathbb{P}(\mathcal{R} \cap \mathcal{H}_0) \leq \alpha$  where  $\mathcal{H}_0 = \{j : H_j = 0\}$ . The *Bonferroni rejection set* is  $\mathcal{R} = \{j : P_j < \alpha/m\} = \{j : T_j > z_{\alpha/m}\}$  where we use the notation  $z_\beta = \overline{\Phi}^{-1}(\beta)$ .

The weighted Bonferroni procedure (Rosenthal and Rubin, 1983; Genovese et al., 2006) is as follows. Specify nonnegative weights  $w = (w_1, \dots, w_m)$  and reject hypothesis  $H_j$  if

$$j \in \mathcal{R} = \left\{ j : \frac{P_j}{w_j} \leq \frac{\alpha}{m} \right\}. \quad (1)$$

In the following lemma we show that as long as  $m^{-1} \sum_j w_j \equiv \overline{w} = 1$  the rejection set  $\mathcal{R}$  controls familywise error at level  $\alpha$ . The second lemma includes a simple modification that will be needed later.

**Lemma 2.1** *If  $\overline{w} = 1$ , then  $\mathcal{R}$  controls familywise error at level  $\alpha$ .*

**Lemma 2.2** *Suppose that  $W_j = g(V_j, c)$ ,  $j = 1, \dots, m$  for some random variables  $V_1, \dots, V_m$ , some constant  $c$  and some function  $g$ . Further, suppose that  $V_j$  has a known distribution  $H$  whenever  $j \in \mathcal{H}_0$  and that  $P_j$  is independent of  $V_j$  for all  $j \in \mathcal{H}_0$ . The rule that rejects when  $P_j \leq \alpha W_j/m$  controls familywise error at level  $\alpha$  if  $c$  is chosen to satisfy  $\mathbb{E}_H(g(V_j, c)) \leq 1$ .*

Genovese et al. (2006) also showed that false discovery methods benefit by weighting. Recall that the false discovery proportion (FDP) is

$$\text{FDP} = \frac{\text{number of false rejections}}{\text{number of rejections}} = \frac{|\mathcal{R} \cap \mathcal{H}_0|}{|\mathcal{R}|}$$

where the ratio is defined to be 0 if the denominator is 0. The false discovery rate (FDR) is  $\text{FDR} = \mathbb{E}(\text{FDP})$ . Benjamini and Hochberg (1995) proved  $\text{FDR} \leq \alpha$  if  $\mathcal{R} = \{j : P_{(j)} \leq T\}$  where  $T = \max\{j : P_{(j)} \leq j\alpha/m\}$ . Genovese et al. (2006) showed that  $\text{FDR} \leq \alpha$  if the  $P'_j$ 's are replaced by  $Q_j = P_j/w_j$  provided  $\bar{w} = 1$ . This paper focuses on familywise error using the weighted procedure (1). Similar results hold for FDR, and other familywise controlling procedures such as Holm's test.

### 3 Power, Robustness and Optimality

The optimal weights, derived below, can be re-expressed as optimal cutoffs for testing. Specifically, rejecting when  $P_j/w_j \leq \alpha/m$  is the same as rejection when  $T_j > \xi_j/2 + c/\xi_j$ . This result can be obtained from Spjøtvoll (1972) and is identical to the result in Rubin et al. (2006) obtained independently. The remainder of the paper, which shows some good properties of the weighted method, can thus also be considered as providing support for their method for selecting test specific cutoffs. In particular, Rubin et al. (2007)'s simulations indicate that even poorly specified estimates of the cutoffs  $\xi_j/2 + c/\xi_j$  can still perform well. In this section we provide insight into why this is true.

The power of a single, one-sided alternative in the unweighted case ( $w_j = 1$ ) is

$$\pi(\xi_j, 1) = \mathbb{P}(T_j > z_{\alpha/m}) = \bar{\Phi}(z_{\alpha/m} - \xi_j).$$

The power<sup>1</sup> in the weighted case is

$$\pi(\xi_j, w_j) = \mathbb{P}\left(P_j < \frac{\alpha w_j}{m}\right) = \mathbb{P}\left(T_j > \bar{\Phi}^{-1}\left(\frac{\alpha w_j}{m}\right)\right) = \bar{\Phi}\left(z_{\alpha w_j/m} - \xi_j\right). \quad (2)$$

Weighting increases the power when  $w_j > 1$  and decreases the power when  $w_j < 1$  for the  $j$ 'th alternative.

Given  $\theta = (\xi_1, \dots, \xi_m)$  and  $w = (w_1, \dots, w_m)$  we define the *average power*

$$\frac{1}{m_1} \sum_{j=1}^m \pi(\xi_j, w_j) I(\xi_j > 0),$$

where  $m_1 = \sum_{j=1}^m I(\xi_j > 0)$ . More generally, if  $\xi$  is drawn from a distribution  $Q$  and  $w = w(\xi)$  is a weight function we define the average power  $\int \pi(\xi, w(\xi)) I(\xi > 0) dQ(\xi) / \int I(\xi > 0) dQ(\xi)$ . If we take  $Q$  to be the empirical distribution of  $(\xi_1, \dots, \xi_m)$  then this reduces to the previous expression. In this formulation we require  $w(\xi) \geq 0$  and  $\int w(\xi) dQ(\xi) = 1$ .

In the following theorem we see that the set of optimal weight functions form a one parameter family indexed by a constant  $c$ .

**Theorem 3.1** *Given  $\theta = (\xi_1, \dots, \xi_m)$ , the optimal weight vector  $w = (w_1, \dots, w_m)$  that maximizes the average power subject to  $w_j \geq 0$  and  $\bar{w} = 1$  is  $w = (\rho_c(\xi_1), \dots, \rho_c(\xi_m))$  where*

$$\rho_c(\xi) = \left(\frac{m}{\alpha}\right) \bar{\Phi}\left(\frac{\xi}{2} + \frac{c}{\xi}\right) I(\xi > 0), \quad (3)$$

and  $c \equiv c(\theta)$  is defined by the condition

$$\frac{1}{m} \sum_{j=1}^m \rho_c(\xi_j) = 1. \quad (4)$$

The proof, essentially a special case of Spjøtvoll (1972), is in the appendix. Figure 2 displays the function  $\rho_c(\xi)$  for various values of  $c$  (the function is normalized to have maximum 1 for easier visualization). The result generalizes to the case where the alternative means are random variables with distribution  $Q$  in which case  $c$  is defined by  $\int \rho_c(\xi) dQ(\xi) = 1$ .

From (2) and (3) we have immediately:

---

<sup>1</sup>For a two-sided alternative the power is

$$\pi(\xi_j, w_j) = \bar{\Phi}\left(z_{\alpha w_j/2m} - \xi_j\right) + \bar{\Phi}\left(z_{\alpha w_j/2m} + \xi_j\right).$$

**Lemma 3.2** *The power at an alternative with mean  $\xi$  under optimal weights is  $\bar{\Phi}(c/\xi - \xi/2)$ . The average power under optimal weights, which we call the **oracle power**, is*

$$\frac{1}{m_1} \sum_{j=1}^m \bar{\Phi} \left( \frac{c}{\xi_j} - \frac{\xi_j}{2} \right) I(\xi_j > 0)$$

where  $m_1 = \sum_j I(\xi_j > 0)$ . The oracle power is not attainable since the optimal weights depend on  $\theta = (\xi_1, \dots, \xi_m)$ . In practice, the weights will either be chosen by prior information or by estimating the  $\xi$ 's. This raises the following question: how sensitive is the power to correct specification of the weights? Now we show that the power is very robust to weight misspecification.

**Property I:** Sparse weights (minimum weight close to 1) are highly robust. If most weights are less than 1 and the minimum weight is close to 1 then correct specification (large weights on alternatives) leads to large power gains but incorrect specification (large weights on nulls) leads to little power loss.

**Property II:** Worst case analysis. Weighted hypothesis testing, even with poorly chosen weights, typically does as well or better than Bonferroni except when the the alternative means are large, in which both have high power.

Let us now make these statements precise. Also, see Genovese et al. (2006) and Roeder et al. (2006) for other results on the effect of weight misspecification.

**Property I.** Consider first the case where the weights take two distinct values and the alternatives have a common mean  $\xi$ . Let  $\epsilon$  denote the fraction of hypotheses given the larger of the two values of the weights  $B$ . Then, the weight vector  $w$  is proportional to

$$\left( \underbrace{B, \dots, B}_{k \text{ terms}}, \underbrace{1, \dots, 1}_{m-k \text{ terms}} \right)$$

where  $k = \epsilon m$  and  $B > 1$  and hence the normalized weights are

$$w = \left( \underbrace{w_1, \dots, w_1}_{k \text{ terms}}, \underbrace{w_0, \dots, w_0}_{m-k \text{ terms}} \right)$$

where

$$w_1 = \frac{B}{\epsilon B + (1 - \epsilon)}, \quad w_0 = \frac{1}{\epsilon B + (1 - \epsilon)}.$$

We say that the weights are sparse if  $\epsilon$  is small. Provided  $B$  is considerably less than  $1/\epsilon$ , most weights are near 1 in the sparse case.

Rather than investigate the average power, we focus on a single alternative with mean  $\xi$ . The power gain by up-weighting this hypothesis is the power under weight  $w_1$  minus the unweighted power  $\pi(\xi, w_1) - \pi(\xi, 1)$ . Similarly, the power loss for down-weighting is  $\pi(\xi, 1) - \pi(\xi, w_0)$ . The gain minus the loss, which we call the robustness function, is

$$\begin{aligned} R(B, \epsilon) &\equiv \left( \pi(\xi, w_1) - \pi(\xi, 1) \right) + \left( \pi(\xi, 1) - \pi(\xi, w_0) \right) \\ &= \bar{\Phi} \left( z_{\alpha w_1/m} - \xi \right) + \bar{\Phi} \left( z_{\alpha w_0/m} - \xi \right) - 2\bar{\Phi} \left( z_{\alpha/m} - \xi \right). \end{aligned}$$

The gain outweighs the loss if and only if  $R(B, \epsilon) > 0$ .

In the sparse weighting scenario  $k$  is small and  $w_0 \approx 1$  by assumption, consequently, an analysis of  $R(B, \epsilon)$  sheds light on the effect of weighting on power, without the added complications involved in a full analysis of average power.

**Theorem 3.3** *Fix  $B > 1$ . Then,  $\lim_{\epsilon \rightarrow 0} R(B, \epsilon) > 0$ . Moreover, there exists  $\epsilon^*(B) > 0$  such that  $R(B, \epsilon) > 0$  for all  $\epsilon < \epsilon^*(B)$ .*

We can generalize this beyond the two-valued case as follows. Let  $w$  be any weight vector such that  $\bar{w} = 1$ . Now define the (worst case) robustness function

$$R(\xi) \equiv \min_{\{j: w_j > 1, H_j = 1\}} \{ \pi(\xi, w_j) - \pi(\xi, 1) \} - \max_{\{j: w_j < 1, H_j = 1\}} \{ \pi(\xi, 1) - \pi(\xi, w_j) \}.$$

We will see that  $R(\xi) > 0$  under weak conditions and that the maximal robustness is obtained for  $\xi$  near the Bonferroni cutoff  $z_{\alpha/m}$ .

**Theorem 3.4** *A necessary and sufficient condition for  $R(\xi) > 0$  is*

$$R_{b,B}(\xi) \equiv \Phi \left( z_{\alpha B/m} - \xi \right) + \Phi \left( z_{\alpha b/m} - \xi \right) - 2\Phi \left( z_{\alpha/m} - \xi \right) \leq 0 \quad (5)$$

where  $B = \min\{w_j : w_j > 1\}$ ,  $b = \min\{w_j\}$ . Moreover,

$$R_{b,B}(\xi) = -\Delta(\xi) + O(1 - b)$$

where

$$\Delta(\xi) = \left( \Phi\left(z_{\alpha/m} - \xi\right) - \Phi\left(z_{\alpha B/m} - \xi\right) \right) > 0$$

and, as  $b \rightarrow 1$ ,  $\mu(\{\xi : R(\xi) < 0\}) \rightarrow 0$  and  $\inf_{\xi} R(\xi) \rightarrow 0$ .

Based on the theorem we see that there is overwhelming robustness as long as the minimum weight is near 1. Even in the extreme case  $b = 0$ , there is still a safe zone, an interval of values of  $\xi$  over which  $R(\xi) > 0$ .

**Lemma 3.5** *Suppose that  $B \geq 2$ . Then there exists  $\xi_* > 0$  such that  $R_{B,b}(\xi) > 0$  for all  $0 \leq \xi \leq \xi_*$  and all  $b$ . An upper bound on  $\xi_*$  is  $z_{\alpha/m} - 1/(z_{\alpha/m} - z_{B\alpha/m})$ .*

**Property II.** Even if the weights are not sparse, the power of the weighted test tends to be acceptable.

The result holds even though the weights themselves can be very sensitive to changes in  $\theta$ . Consider the following example. Suppose that  $\theta = (\xi_1, \dots, \xi_m)$  where each  $\xi$  is equal to either 0 or some fixed number  $\xi$ . The empirical distribution of the  $\xi_j$ 's is thus  $Q = (1 - a)\delta_0 + a\delta_{\xi}$  where  $\delta$  denotes a point mass and  $a$  is the fraction of nonzero means. The optimal weights are 0 for  $\xi_j = 0$  and  $1/a$  for  $\xi_j = \xi$ . Let  $\tilde{Q} = (1 - a - \gamma)\delta_0 + \gamma\delta_u + a\delta_{\xi}$  where  $u$  is a small positive number. Since we have only moved the mass at 0 to  $u$ , and  $u$  is small, we would hope that  $w(\xi)$  will not change much. But this is not the case. Set  $\xi = A + \sqrt{A^2 - 2c}$ ,  $u = B - \sqrt{B^2 - 2c}$  where

$$A = \bar{\Phi}^{-1}\left(\frac{\alpha}{(m(\gamma K + a))}\right), \quad B = \bar{\Phi}^{-1}\left(\frac{K\alpha}{(m(\gamma K + a))}\right).$$

This arrangement yields weights  $w_0$  and  $w_1$  on  $u$  and  $\xi$  such that  $w_0/w_1 = K$ . For example, if  $m = 1000$ ,  $\alpha = 0.05$ ,  $a = .1$ ,  $\gamma = .1$ ,  $K = 1000$ , and  $c = .1$ , then  $u = .03$  and  $\xi = 9.8$ . The optimal weight on  $\xi$  under  $Q$  is 10 but under  $\tilde{Q}$  it is .00999 and so is reduced by a factor of 1001.

More generally we have the following result which shows that the weights are, in a certain sense, a discontinuous function of  $\theta$ .

**Lemma 3.6** *Fix  $\alpha$  and  $m$ . For any  $\delta > 0$  and  $\epsilon > 0$  there exists  $Q = (1 - a)\delta_0 + a\delta_\xi$  and  $\tilde{Q} = (1 - a - \gamma)\delta_0 + \gamma\delta_u + a\delta_\xi$  such that  $d(Q, \tilde{Q}) < \delta$ , and  $\tilde{\rho}(\xi)/\rho(\xi) < \epsilon$ , where  $a = \alpha/4$ ,  $d(Q, \tilde{Q}) = \sup_\xi |Q(-\infty, \xi] - \tilde{Q}(-\infty, \xi]|$  is the Kolmogorov-Smirnov distance,  $\rho$  is the optimal weight function for  $Q$  and  $\tilde{\rho}$  is the optimal weight function for  $\tilde{Q}$ .*

Fortunately, this feature of the weight function does not pose a serious hurdle in practice because it is possible to have high power even with poor weights. In Figure 4 the plots on the left show the power as a function of the alternative mean  $\xi$ . The dark solid line shows the lowest possible power assuming the weights were estimated as poorly as possible (under conditions specified below). The lighter solid line is the power of the unweighted (Bonferroni) method. The dotted line shows the power under theoretically optimal weights. The worst case weighted power is typically close to or larger than the Bonferroni power except for large  $\xi$  when they are both large.

To begin formal analysis assume that each mean is either equal to 0 or  $\xi$  for some fixed  $\xi > 0$ . Thus, the empirical distribution is  $Q = (1 - a)\delta_0 + a\delta_\xi$  where  $\delta$  denotes a point mass and  $a$  is the fraction of nonzero  $\xi_j$ 's. The optimal weights are  $1/a$  for hypotheses whose mean is  $\xi$ . To study the effect of misspecification error, consider the case where  $\gamma m$  nulls are mistaken for alternatives with mean  $u > 0$ . This corresponds to misspecifying  $Q$  to be  $\tilde{Q} = (1 - a - \gamma)\delta_0 + \gamma\delta_u + a\delta_\xi$ . We will study the effect of varying  $u$  so let  $\pi(u)$  denote the power at the true alternative  $\xi$  as a function of  $u$ . Also, let  $\pi_{\text{Bonf}}$  denote the power using equal weights (Bonferroni). Note that changing  $Q = (1 - a)\delta_0 + a\delta_\xi$  to  $Q = (1 - a)\delta_0 + a\delta_{\xi'}$  for  $\xi' \neq \xi$  does not change the weights.

As the weights are a function of  $c$ , we first need to find  $c$  as a function of  $u$ . The normalization condition (4) reduces to

$$\gamma \bar{\Phi} \left( \frac{u}{2} + \frac{c}{u} \right) + a \bar{\Phi} \left( \frac{\xi}{2} + \frac{c}{\xi} \right) = \frac{\alpha}{m} \quad (6)$$

which implicitly defines the function  $c(u)$ . First we consider what happens when  $u$  is restricted to be less than  $\xi$ .

**Theorem 3.7** Assume that  $\alpha/m \leq \gamma + a \leq 1$ . Let  $Q = (1 - a)\delta_0 + a\delta_\xi$  and  $\tilde{Q} = (1 - a - \gamma)\delta_0 + \gamma\delta_u + a\delta_\xi$  with  $0 \leq u \leq \xi$ . Let  $C(\xi) = \sup_{0 \leq u \leq \xi} c(u)$  and define  $\xi_0 = z_{\alpha/(m(\gamma+a))}$ ,

1. For  $\xi \leq \xi_0$ ,  $C(\xi) = \xi\xi_0 - \xi^2/2$ . For  $\xi > \xi_0$ ,  $C(\xi)$  is the solution to

$$\gamma\bar{\Phi}(\sqrt{2c}) + a\bar{\Phi}\left(\frac{c}{\xi} + \frac{\xi}{2}\right) = \frac{\alpha}{m}.$$

In this case,  $C(\xi) = z_{\alpha/(m\gamma)}^2/2 + O(a)$ .

2. Let  $\xi_* = z_{\alpha/m} + \sqrt{z_{\alpha/m}^2 - z_q^2}$ , where  $q = \alpha(1 - a)/(m\gamma)$ . For  $\xi < \xi_*$ ,

$$\inf_{0 < u < \xi} \pi(u) \geq \pi_{\text{Bonf}}. \quad (7)$$

For  $\xi \geq \xi_*$  we have

$$\inf_{0 < u < \xi} \pi(u) \geq \bar{\Phi}\left(\frac{z_{\alpha/(m\gamma)}^2 - \xi_*^2}{2\xi_*}\right) - O(a) \quad (8)$$

$$\approx 1 - \bar{\Phi}\left(\sqrt{2 \log \frac{1-a}{\gamma}}\right) - O(a) \quad (9)$$

$$\geq 1 - \frac{\gamma}{1-a} - O(a). \quad (10)$$

The factor  $\bar{\Phi}\left(\sqrt{2 \log(1-a)/\gamma}\right) \approx \gamma/(1-a)$  is the worst case power deficit due to misspecification. Now we drop the assumption that  $u \leq \xi$ .

**Theorem 3.8** Let  $Q = (1 - a)\delta_0 + a\delta_\xi$  and let  $Q_u \equiv (1 - a - \gamma)\delta_0 + \gamma\delta_u + a\delta_\xi$ . Let  $\pi_u$  denote the power at  $\xi$  using the weights computed under  $Q_u$ .

1. The least favorable  $u$  is  $u_* \equiv \operatorname{argmin}_{u \geq 0} \pi_u = \sqrt{2c_*} = z_{\alpha/(m\gamma)} + O(a)$  where  $c_*$  solves

$$\gamma\bar{\Phi}(\sqrt{2c_*}) + a\bar{\Phi}\left(\frac{\xi}{2} + \frac{c_*}{\xi}\right) = \frac{\alpha}{m}$$

and  $c_* = z_{\alpha/(m\gamma)}^2/2 + O(a)$ .

2. The minimal power is

$$\inf_u \pi_u = \bar{\Phi}\left(\frac{c_*}{\xi} - \frac{\xi}{2}\right) = \bar{\Phi}\left(\frac{z_{\alpha/(m\gamma)}^2 - \xi^2}{2\xi}\right) + O(a).$$

3. A sufficient condition for  $\inf_u \pi_u$  to be larger than the power of the Bonferroni method is

$$\xi \geq z_{\alpha/m} + \sqrt{z_{\alpha/m}^2 - z_{\alpha/(m\gamma)}^2} + O(a).$$

## 4 Choosing External Weights

One approach to choosing external weights (or test statistic cutoffs) is to use empirical Bayes methods to model prior information while being careful to preserve error control as in Westfall and Soper (2001) for example. Here we consider a simple method that takes advantage of the robustness properties we have discussed. We will focus here on the two-valued case. Thus,

$$w = \underbrace{(w_1, \dots, w_1)}_{k \text{ terms}}, \underbrace{(w_0, \dots, w_0)}_{m-k \text{ terms}}$$

where  $k = \epsilon m$ ,  $w_1 = B/(\epsilon B + (1 - \epsilon))$  and  $w_0 = 1/(\epsilon B + (1 - \epsilon))$ . In practice, we would typically have a fixed fraction of hypotheses  $\epsilon$  that we want to give more weight to. The question is how to choose  $B$ . We will focus on choosing  $B$  to produce weights with good properties at interesting values of  $\xi$ . Now large values of  $\xi$  already have high power. Very small values of  $\xi$  have extremely low power and benefit little by weighting. This leads us to focus on constructing weights that are useful for a *marginal effect*, defined as the alternative  $\xi_0$  that has power 1/2 when given weight 1. Thus, the marginal effect is  $\xi_0 = z_{\alpha/m}$ . In the rest of this section then we assume that all nonzero  $\xi_j$ 's are equal to  $\xi_0$ . Of course, the validity of the procedure does not depend on this assumption being true.

Fix  $0 < \epsilon < 1$  and vary  $B$ . As we increase  $B$ , we will eventually reach a point  $B_0(\epsilon)$  where  $R(B, \epsilon) < 0$  which we call turnaround point. Formally,  $B_0(\epsilon) = \sup\{B : R(B, \epsilon) > 0\}$ . The top panel in Figure 5 shows  $B_0(\epsilon)$  versus  $\epsilon$  which shows that for small  $\epsilon$  we can choose  $B$  large without loss of power. The bottom panel shows  $R(B, \epsilon)$  for  $\epsilon = 0.1$ . Ideally for a given  $\epsilon$  one chooses  $B$  near  $B_*(\epsilon)$ , the value of  $B$  that maximizes  $R(B, \epsilon)$ .

**Theorem 4.1** Fix  $0 < \epsilon < 1$ . As a function of  $B$ ,  $R(B, \epsilon)$  is unimodal and satisfies  $R(1, \epsilon) = 1$ ,  $R'(1, \epsilon) > 0$  and  $R(\infty, \epsilon) < 0$ . Hence,  $B_0(\epsilon)$  exists and is unique. Also,  $R(B, \epsilon)$  has a unique maximum at some point  $B^*(\epsilon)$  and  $R(B^*(\epsilon), \epsilon) > 0$ .

When  $\epsilon$  is very small,  $B$  can be large, provided  $w_0 \approx 1$ . For example, suppose we want to increase the chance of rejecting one particular hypothesis so that  $\epsilon = 1/m$ . Then,

$$w_1 = \frac{mB}{B + m - 1} \approx B, \quad w_0 = \frac{1}{B + m - 1} \approx 1$$

and

$$\lim_{m \rightarrow \infty} \lim_{B \rightarrow \infty} \pi(\xi_j, w_1) = 1, \quad \text{while} \quad \lim_{m \rightarrow \infty} \lim_{B \rightarrow \infty} \pi(\xi_j, w_0) = \frac{1}{2}.$$

The next results show that binary weighting schemes are optimal in a certain sense. Suppose we want to have at least a fraction  $\epsilon$  with high power  $1 - \beta$  and otherwise we want to maximize the minimum power.

**Theorem 4.2** *Consider the following optimization problem: Given  $0 < \epsilon < 1$  and  $0 < \beta < 1/2$ , find a vector  $w = (w_1, \dots, w_m)$  that maximizes  $\min_j \pi(\xi_m, w_j)$  subject  $\bar{w} = 1$ , and  $\#\{j : \pi(w_j, \xi_m) \geq 1 - \beta\}/m \geq \epsilon$ . The solution is given by  $c = \bar{\Phi}(z_{\alpha/m} + z_{1-\beta})$ ,  $B = \epsilon m(1 - \epsilon)/(\alpha - \epsilon cm)$ ,  $w_1 = B/(\epsilon B + (1 - \epsilon))$ ,  $w_0 = 1/(\epsilon B + (1 - \epsilon))$ , and  $k = \epsilon m$ .*

If our goal is to maximize the number of alternatives with high power while maintaining a minimum power loss, the solution is given as follows.

**Theorem 4.3** *Consider the following optimization problem: Given  $0 < \beta < 1/2$ , find a vector  $w = (w_1, \dots, w_m)$  that maximizes  $\#\{j : \pi(w_j, \xi_m) \geq 1 - \beta\}$  subject to  $\bar{w} = 1$ , and  $\min_j \pi(w_j, \xi_m) \geq \delta$ . The solution is*

$$w_1 = \frac{m \bar{\Phi}(z_{\alpha/m} + z_{1-\beta})}{\alpha}, \quad w_0 = \frac{m \bar{\Phi}(z_{\alpha/m} + z_{\delta})}{\alpha}, \quad \epsilon = \frac{1 - w_0}{w_1 - w_0}$$

and  $k = m\epsilon$ .

A special case that falls under this Theorem permits the minimum power to be 0. In this case  $w_0 = 0$  and  $\epsilon = 1/w_1$ .

## 5 Estimated Weights

In practice  $\xi_j$  is not known, so it must be estimated to utilize the weight function. A natural choice is to build on the two stage experimental design (Satagopan and Elston, 2003; Wang et al., 2006) and split the data into subsets, using one subset to estimate  $\xi_i$ , and hence  $w(\xi_i)$ , and the second to conduct a weighted test of the hypothesis (Rubin et al., 2006). This approach would arise naturally

in an association test conducted in stages. It does lead to a gain in power relative to unweighted testing of stage 2 data; however, it is not better than simply using the full data set without weights for the analysis (Rubin et al., 2006). These results are corroborated by Skol et al. (2006) in a related context. They showed that it is better to use stages 1 and 2 jointly, rather than using stage 2 as an independent replication of stage 1.

To gain a strong advantage with data-based weights, prior information is needed. One option is to order the tests (Rubin et al., 2006), but with a large number of tests this can be challenging. The type of prior information readily available to investigators is often non-specific. For instance, SNPs might naturally be grouped, based on features that make various candidates more promising for this disease under investigation. For a brain-disorder phenotype we might cross-classify SNPs by categorical variables such as functionality, brain expression and so forth. The SNPs in one group may seem most promising, a priori, while those in another seem least promising. Intermediate groups may be somewhat ambiguous. It is easy to imagine additional variables that further partition the SNPs into various classes that help to separate the more promising SNPs from the others. While this type of information lends itself to grouping SNPs, it does not lead directly to weights for the groups. Indeed it might not even be possible to choose a natural ordering of the groups. What is needed is a way to use the data to determine the weights, once the groups are formed.

Until recently, methods for weighted multiple-testing required that prior weights be developed independently of the data under investigation (Genovese et al., 2006; Roeder et al., 2007). Here we provide a data based estimate of weights based on results of grouped analysis. One way to implement this approach is to follow these steps:

1. Partition the tests into subsets  $\mathcal{G}_1, \dots, \mathcal{G}_K$ , with the  $k$ 'th group containing  $r_k$  elements, ensuring that  $r_k$  is at least 20-30.
2. Calculate the sample mean  $Y_k$  and variance  $S_k^2$  for the test statistics in each group.
3. Label the  $i$ 'th test in group  $k$ ,  $T_{ik}$ . At best only a fraction of the elements in each group will have a signal, hence we assume that for  $i = 1, \dots, r_k$  the distribution of the test statistics is

approximated by a mixture model

$$T_{ik} \sim (1 - \pi_k)N(0, 1) + \pi_k N(\xi_k, 1)$$

or

$$T_{ik} \sim (1 - \pi_k)\chi_1^2(0) + \pi_k \chi_1^2(\xi_k^2)$$

where  $\xi_k$  is the signal size for those tests with a signal in the  $k$ 'th group. This is an approximation because the signal is likely to vary across tests. The mixture of normals is only appropriate when the tests are one-sided. For two-sided alternatives, the  $\chi^2$  is the natural approach. This test squares the noncentrality parameter, effectively removing any ambiguity about the direction of the associations.

4. Estimate  $(\pi_k, \xi_k)$  using the method of moments estimator (for details see Appendix). Because  $\xi_k$  has no meaning when  $\pi_k = 0$ , the  $\hat{\xi}_k$  is set to 0 when  $\hat{\pi}_k$  is close to zero. For the normal model the estimators are

$$\hat{\pi}_k = Y_k^2 / (Y_k^2 + S_k^2 - 1), \quad \hat{\xi}_k = Y_k / \pi_k, \quad (11)$$

provided  $\hat{\pi}_k > 1/r_k$ ; otherwise  $\hat{\xi}_k = 0$ .

For the  $\chi^2$  model they are

$$\hat{\xi}_k^2 = \frac{(S_k^2 + Y_k^2 + 3)}{Y_k - 1}, \quad \hat{\pi}_k = \frac{Y_k - 1}{\hat{\xi}_k^2} \quad (12)$$

provided  $Y_k > 1$  and  $1/r_k < \hat{\pi}_k < (r_k - 1)/r_k$ ; otherwise  $\hat{\xi}_k = 0$ .

5. For each of the  $k$  groups, construct weights  $w(\hat{\xi}_k)$ . It is apparent in Figure 1 that if  $|\hat{\xi}_k| < \delta$ , for  $\delta$  near 0, then  $w(\hat{\xi}_k) \approx 0$  and it is unlikely that any tests in the  $k$ 'th group will be significant, regardless of the p-value. The stochastic quantity  $\delta$  depends upon the relative values of  $(\hat{\xi}_1, \dots, \hat{\xi}_K)$ , and the number of elements in each group. For this reason we have found that smoothing the weights generally improves power of the procedure. We suggest using a linear combination such as

$$\hat{w}_k = (1 - \gamma) w(\hat{\xi}_k) + \gamma K^{-1} \sum_k w(\hat{\xi}_k),$$

with  $\gamma = 0.01$  or  $0.05$ . The larger the choice of  $\gamma$ , the more evenly distributed the weights across groups. Alternatively, one could smooth the weights by using a Stein shrinkage estimator or bagging procedure to obtain a more robust estimator of  $(\xi_1, \dots, \xi_K)$  (Hastie et al., 2001). Regardless of how the weights are smoothed, one should renorm them to ensure the weights sum to  $m$ . Each test in group  $k$  receives the weight  $\hat{w}_k$ . Another effect of the smoothing is to ensure that each group gets a weight greater than 0.

This weighting scheme relies on data-based estimators of the optimal weights, but with a partition of the data sufficiently crude to preserve the control of family-wise error rate. The approach is an example of the “sieve principle” (Bickel et al., 1993). The sieve principle works because the number of parameters estimated is far less than the number of observations. Thus, many observations are used to estimate each parameter. Consequently parameters are estimated with substantially less variability than if they were estimated using only the test statistics from the particular gene under investigation. Because the weights are determined by the size of the tests in the entire cluster the probability of upweighting simply because a single test is large, due to chance, is small.

## 6 Examples

### 6.1 Binary weights

In a study of nicotine dependence, Saccone et al. (2007) used binary weights in a candidate gene study. Their study involved 3713 genetic variants (single nucleotide polymorphisms or SNPs) encompassing 348 genes. The genes were divided into two types: 52 nicotinic and dopaminergic receptor genes; and 296 other candidate genes. Each SNP associated with a gene in the first group was allocated ten times the weight of a gene in the other category. Using a generous false discovery rate ( $\alpha = .4$ ), they identified 39 SNPs; 78% of these were nicotine receptors, in contrast to the fraction of nicotine receptors overall (15%).

## 6.2 Independent data weights

For family-based study designs, tests of association are based on transmission data. In these studies, data are available from which one can compute the potential power to detect a signal at each SNP tested; see Ionita-Laza et al. (2007) for a detailed explanation of this unique feature of family-based data. Because the data used to calculate the power are independent of the test statistics for association, these data are available for construction of the weights. Motivated by this possibility, Ionita-Laza et al. (2007) developed a weighting scheme. Using independent data, they ranked the SNPs from most to least promising, in terms of power. They then constructed an exponential weighting scheme, based on simulations of genetic models. The scheme results in a small number of SNPs receiving a top weight, successively more SNPs receiving correspondingly lower weights, and finally a large number receiving the lowest weight. In their simulations they found that the power of the test can often be doubled using this procedure. Using the FHS data they apply the technique to a genome-wide association study with 116,204 SNPs and 923 participants. The phenotype of interest is height. Using their weighting scheme they obtained one significant result with weights and none without weights.

## 6.3 Linkage weights

Finding variation in the genetic code that increases the risk for complex diseases, such as Type II diabetes and schizophrenia, is critically important to the advancement of genetic epidemiology. In the Introduction we describe a means by which weights could be extracted from linkage data. Here we illustrate the idea with both data and simulations.

In the analysis of 955 cases and 1498 controls enrolled in a genome-wide association study, McQueen and colleagues (2008) used weights derived from published linkage results. They combined results from 11 linkage studies on bipolar disorder to obtain  $Z$  scores corresponding to the locations of each association test. From the linkage results they computed weighted p-values using the cumulative normal weight function (Roeder et al., 2006). Although none of their results were genome-wide significant, they obtained promising results in four regions. Three of these are obtained due to a strong p-values in combination with a linkage peak. One signal did not correspond

to a linkage peak, but continued to be in the top tier of p-values, after weights were applied.

To illustrate how binary weights could be derived from such linkage data we present a realistic synthetic example. Using the methods described in Roeder et al. (2006), we create a linkage trace that captures many of the features found in actual linkage traces. In this simulation we generate a full genome (23 chromosomes) and place 20 disease variants at random, one per chromosome. The signals from these variants were designed to yield a weak signals with a broad peaks. Next, we simulated 100,000 normally distributed association test statistics mapped to the same genome. Again, 20 of these tests were generated under the alternative hypothesis of association. These signals were also weak.

To illustrate the synthetic data, six typical chromosomes are displayed in Figure 1. Each displayed chromosome has one true signal, with the association test statistic at that location indicated by an upspike; none of the association tests generated under the null hypothesis are plotted. Without weights, only 2 of the 20 signals could be detected using a Bonferonni correction. Using binary weights, as described above, with  $\epsilon = 0.05$  and  $B = 10$  we discover 5 of the 20 signals. In left column of the figure all three signals were discovered, while in the right column none were discovered (indicated by presence of a down-spike). Comparing the top row, we see that both signals were up-weighted in the correct location, but the association signal was not strong enough in the top right chromosome to achieve significance. Alternatively, in the bottom left panel the association statistic was substantial enough to reject the null hypothesis without the benefit of up-weighting.

To examine the robustness of the procedure to choice of weights, we tried 4 choices of  $\epsilon$  (.01,.05,.1,.2) with  $1 \leq B \leq 50$ . We made no false discoveries with any of these choices. The power is displayed in Figure 6. To assist in the choice of parameters we have found it helpful to examine the number of discoveries for each choice. In this example, the number of discoveries varied between 2 (unweighted, i.e.  $B = 1$ ) to 6 ( $\epsilon = .2, B \geq 10$ ). 5 discoveries were made for a broad range of choices. In principle choosing  $(\epsilon, B)$  to maximize the number of discoveries can inflate the error rate. In our simulations we have found that seaching within the family of weights defined by 1 or 2 parameters, such as this binary weight system based upon a linkage trace, tends to provide very close to nominal protection against false discoveries.

## 7 Discussion

Several authors have explored the effect of weights on power of multiple testing procedures (e.g. Westfall et al. 2004). These investigations show that the power of multiple testing procedures can be increased by using weighted p-values. Here we derive the optimal weights for a commonly used family of tests and show that the power is remarkably robust to misspecification of these weights.

The same ideas used here can be applied to other testing methods to improve power. In particular, weights can be added to the FDR method, Holm's stepdown test, and the Donoho and Jin (2004) method. Weighting ideas can also be used for confidence intervals. Another open question is the connection with Bayesian methods which have already been developed to some extent in Efron et al. (2001).

GWAS for some phenotypes such as Type 1 diabetes have yielded exciting results Todd et al. (2007), while results for other complex diseases have been much less successful. Presumably many studies do not have sufficient power to detect the genetic variants associated with the phenotypes, even though thousands of cases and controls have been genotyped. To bolster power, we recommend up-weighting and down-weighting hypotheses, based on prior likelihood of association with the phenotype. For instance, Wang et al. (2007) describe pathway-based approaches for the analysis of GWAS.

Multiple testing arises in GWAS analyses in other contexts as well. Frequently multiple tests, assuming different genetic models, are applied to each genetic marker. Multiple markers in a neighborhood can be analyzed simultaneously to increase the signal, using haplotypes, multivariate models, and fine-mapping techniques. Data are often collected in multiple stages of the experiment, and at each stage promising markers are tested for association. In summary, many questions concerning multiple testing remain open in the context of GWAS.

## 8 Appendix

Proof of Lemma 2.1. The familywise error is

$$\begin{aligned} \mathbb{P}((\mathcal{R} \cap \mathcal{H}_0) > 0) &= \mathbb{P}\left(P_j \leq \frac{\alpha w_j}{m} \text{ for some } j \in \mathcal{H}_0\right) \\ &\leq \sum_{j \in \mathcal{H}_0} \mathbb{P}\left(P_j \leq \frac{\alpha w_j}{m}\right) = \frac{\alpha}{m} \sum_{j \in \mathcal{H}_0} w_j \leq \alpha \bar{w} = \alpha. \quad \blacksquare \end{aligned}$$

Proof of Lemma 2.2. The familywise error is

$$\begin{aligned} \mathbb{P}((\mathcal{R} \cap \mathcal{H}_0) > 0) &= \mathbb{P}\left(P_j \leq \frac{\alpha W_j}{m} \text{ for some } j \in \mathcal{H}_0\right) \\ &\leq \sum_{j \in \mathcal{H}_0} \mathbb{P}\left(P_j \leq \frac{\alpha W_j}{m}\right) = \sum_{j \in \mathcal{H}_0} \mathbb{E}_H \left( \mathbb{P}\left(P_j \leq \frac{\alpha w_j}{m} \mid W_j = w_j\right) \right) \\ &= \sum_{j \in \mathcal{H}_0} \mathbb{E}_H(\alpha W_j/m) = \frac{\alpha}{m} \sum_{j \in \mathcal{H}_0} \mathbb{E}_H(W_j) \\ &\leq \frac{m_0 \alpha}{m} \leq \alpha. \end{aligned}$$

Proof of Theorem 3.1. Let  $C$  denote the set of hypotheses with  $\xi_j > 0$ . Power is optimized if  $w_j = 0$  for  $j \notin C$ . The average power is

$$\frac{1}{m_1} \sum_{j \in C} \bar{\Phi} \left( \bar{\Phi}^{-1} \left( \frac{\alpha w_j}{m} \right) - \xi_j \right).$$

with constraint

$$\sum w_j = m.$$

Choose  $\underline{w}$  to maximize

$$\pi = \frac{1}{m_1} \sum_{j \in C} \bar{\Phi} \left( \bar{\Phi}^{-1} \left( \frac{\alpha w_j}{m} \right) - \xi_j \right) - \lambda \left( m - \sum w_i \right)$$

by setting the derivative to zero

$$\frac{\partial}{\partial w_i} \pi = -\lambda + \frac{\phi \left( \bar{\Phi}^{-1} \left( \frac{\alpha w_j}{m} \right) - \xi_j \right) \alpha}{\phi \left( \bar{\Phi}^{-1} \left( \frac{\alpha w_j}{m} \right) \right) m} = 0$$

$$\frac{m\lambda}{\alpha} = \frac{\phi\left(\bar{\Phi}^{-1}\left(\frac{\alpha w_j}{m}\right) - \xi_j\right)}{\phi\left(\bar{\Phi}^{-1}\left(\frac{\alpha w_j}{m}\right)\right)}$$

The  $\underline{w}$  that solves these equations is given in (3). Finally, solve for  $c$  such that  $\sum_i w_i = m$ . ■

**Proof of Theorem 3.4.** The first statement follows easily by noting that the worst case corresponds to choosing weight  $B$  in the first term in  $R(\xi)$  and choosing weight  $b$  in the second term in  $R(\xi)$ . The rest follows by Taylor expanding  $R_{b,B}(\xi)$  around  $b = 1$ . ■

**Proof of Lemma 3.5.** With  $b = 0$ ,  $R_{b,B}(\xi) \geq 0$  when

$$\bar{\Phi}(z_{B\alpha/m} - \xi) - 2\bar{\Phi}(z_{\alpha/m} - \xi) \geq 0. \quad (13)$$

With  $B \geq 2$ , (13) holds at  $\xi = 0$ . The left hand side is increasing in  $\xi$  for  $\xi$  near 0 but (13) does not hold at  $\xi = z_{\alpha/m}$ . So (13) must hold in the interval  $[0, \xi_*]$ . Rewrite (13) as  $\bar{\Phi}(z_{B\alpha/m} - \xi) - \bar{\Phi}(z_{\alpha/m} - \xi) \geq \bar{\Phi}(z_{\alpha/m} - \xi)$ . We lower bound the left hand side and upper bound the right hand side. The left hand side is  $\bar{\Phi}(z_{B\alpha/m} - \xi) - \bar{\Phi}(z_{\alpha/m} - \xi) = \int_{z_{B\alpha/m} - \xi}^{z_{\alpha/m} - \xi} \phi(u) du \geq (z_{\alpha/m} - z_{B\alpha/m})\phi(z_{\alpha/m} - \xi)$ . The right hand side can be bounded using Mill's ratio:  $\bar{\Phi}(z_{\alpha/m} - \xi) \leq \phi(z_{\alpha/m} - \xi)/(z_{\alpha/m} - \xi)$ . Set the lower bound greater than the upper bound to obtain the stated result. ■

**Proof of Lemma 3.6.** Choose  $K > 1$  such that  $1/(K + 1) < 1/a - \epsilon$ . Choose  $1 > \gamma > (2\alpha - a)/K$ . Choose a small  $c > 0$ . Let  $\xi = A + \sqrt{A^2 - 2c}$  and  $u = B - \sqrt{B^2 - 2c}$  where

$$A = \bar{\Phi}^{-1}\left(\frac{\alpha}{(m(\gamma K + a))}\right), \quad B = \bar{\Phi}^{-1}\left(\frac{K\alpha}{(m(\gamma K + a))}\right).$$

Then  $\rho(\xi) = 1/a$  and  $\tilde{\rho}(\xi) = 1/(K + 1)$ . Now  $d(Q, \tilde{Q}) = \gamma$ . Taking  $K$  sufficiently large and  $\gamma$  sufficiently close to  $(2\alpha - a)/K$  makes  $\gamma < \delta$ . ■

It is convenient to prove Theorem 3.8 before proving Theorem 3.7.

Proof of Theorem 3.8. Let  $c_*$  solve

$$\gamma\bar{\Phi}(\sqrt{2c_*}) + a\bar{\Phi}\left(\frac{\xi}{2} + \frac{c_*}{\xi}\right) = \frac{\alpha}{m}. \quad (14)$$

We claim first that for any  $c > c_*$ , there is no  $u$  such that the weights average to 1. Fix  $c > c_*$ . The weights average to 1 if and only if

$$\gamma\bar{\Phi}\left(\frac{c}{u} + \frac{u}{2}\right) + a\bar{\Phi}\left(\frac{\xi}{2} + \frac{c}{\xi}\right) = \frac{\alpha}{m}. \quad (15)$$

Since  $c > c_*$  and since the second term is decreasing in  $c$ , we must have

$$\bar{\Phi}\left(\frac{c}{u} + \frac{u}{2}\right) > \bar{\Phi}(\sqrt{2c_*}).$$

The function  $r(u) = \bar{\Phi}(c/u + u/2)$  is maximized at  $u = \sqrt{2c}$ . So  $r(\sqrt{2c}) \geq r(u)$ . But  $r(\sqrt{2c}) = \bar{\Phi}(\sqrt{2c})$ . Hence  $\bar{\Phi}(\sqrt{2c}) \geq r(u) \geq \bar{\Phi}(\sqrt{2c_*})$ . This implies  $c < c_*$  which is a contradiction. This establishes that  $\sup_u c(u) \leq c_*$ . On the other hand, taking  $c = c_*$  and  $u = \sqrt{2c_*}$  solves equation (15). Thus  $c_*$  is indeed the largest  $c$  that solves the equation which establishes the first claim. The second claim follows by noting that

$$\gamma\bar{\Phi}(\sqrt{2c_*}) + a\bar{\Phi}\left(\frac{\xi}{2} + \frac{c_*}{\xi}\right) = \gamma\bar{\Phi}(\sqrt{2c_*}) + O(a).$$

Now set this expression equal to  $\alpha/m$  and solve. ■

**Proof of Theorem 3.7.** Define  $c_*$  as in (14). If  $u_* = \sqrt{2c_*} \leq \xi$  then the the proof proceeds as in the previous proof. So we first need to establish for which values of  $\xi$  is this true. Let  $r(c) = \gamma\bar{\Phi}(\sqrt{2c}) + a\bar{\Phi}(\xi/2 + c/\xi)$ . We want to find out when the solution of  $r(c) = \alpha/m$  is such that  $\sqrt{2c} \leq \xi$ , or equivalently,  $c \leq \xi^2/2$ . Now  $r$  is decreasing in  $c$ . Since  $\gamma + a \geq \alpha/m$ ,  $r(-\infty) \geq \alpha/m$ . Hence there is a solution with  $c \leq \xi^2/2$  if and only if  $r(\xi^2/2) \leq \alpha/m$ . But  $r(\xi^2/2) = (\gamma + a)\bar{\Phi}(\xi)$  so we conclude that there is such a solution if and only if  $(\gamma + a)\bar{\Phi}(\xi) \leq \alpha/m$ , that is,  $\xi \geq z_{\alpha/(m(\gamma+a))} = \xi_0$ .

Now suppose that  $\xi < \xi_0$ . We need to find  $u \leq \xi$  to make  $c$  as large as possible in the equation  $v(u, c) \equiv \gamma\bar{\Phi}(u/2 + c/u) + a\bar{\Phi}(\xi/2 + c/\xi) = \alpha/m$ . Let  $u_* = \xi$  and  $c_* = \xi z_{\alpha/(m(\gamma+a))} - \xi^2/2$ .

By direct substitution,  $v(u_*, c_*) = \alpha/m$  for this choice of  $u$  and  $c$  and clearly  $u_* \leq \xi$  as required. We claim that this is the largest possible  $c_*$ . To see this, note that  $v(u, c) < v(u, c_*)$ . For  $\xi \leq \xi_0$ ,  $v(u, c_*)$  is a decreasing function of  $u$ . Hence,  $v(u, c) < v(u, c_*) \leq v(u_*, c_*) = \alpha/m$ . This contradicts the fact that  $v(u, c) = \alpha/m$ .

For the second claim, note that the power of the weighted test beats the power of Bonferroni if and only if the weight  $w = (m/\alpha)\overline{\Phi}(\xi/2 + C(\xi)/2) \geq 1$  which is equivalent to

$$C(\xi) \leq \xi z_{\alpha/m} - \xi^2/2. \quad (16)$$

When  $\xi \leq \xi_0$ ,  $C(\xi) = \xi\xi_0 - \xi^2/2$ . By assumption,  $\gamma + a \leq 1$  so that  $z_{\alpha/(m(\gamma+a))} \leq z_{\alpha/m}$  and Now suppose that  $\xi_0 < \xi \leq \xi_*$ . Then  $C(\xi)$  is the solution to  $r(c) = \gamma\overline{\Phi}(\sqrt{2c}) + a\overline{\Phi}(\xi/2 + c/\xi) = \alpha/m$ . We claim that (16) still holds. Suppose not. Then, since  $r(c)$  is decreasing in  $c$ ,  $r(\xi z_{\alpha/m} - \xi^2/2) > r(C(\xi)) = \alpha/m$ . But, by direct calculation,  $r(\xi z_{\alpha/m} - \xi^2/2) > \alpha/m$  implies that  $\xi > \xi_*$  which is a contradiction. Thus (7) holds.

Finally, we turn to (8). In this case,  $C(\xi) = z_{\alpha/(m\gamma)}^2/2 + O(a)$ . The worst case power is  $\overline{\Phi}(C(\xi)/\xi - \xi/2) = \overline{\Phi}(z_{\alpha/(m\gamma)}^2/(2\xi) - \xi/2) + O(a)$ . The latter is increasing in  $\xi$  and so is at least  $\overline{\Phi}(z_{\alpha/(m\gamma)}^2/(2\xi_*) - \xi_*/2) + O(a) = \overline{\Phi}((z_{\alpha/(m\gamma)}^2/(2\xi_*) - \xi_*^2)/(2\xi_*)) + O(a)$  as claimed. The next two equations follow from standard tail approximations for Gaussians. Specifically, a Gaussian quantile  $z_{\beta/m}$  can be written as  $z_{\beta/m} = \sqrt{2 \log(mL_m/\beta)}$  where  $L_m = c \log^a(m)$  for constants  $a$  and  $c$  Donoho and Jin (2004). Inserting this into the previous expression yields the final expression.

■

**Proof of Theorem 4.2.** Setting  $\pi(w, \xi_m) = \overline{\Phi}(\overline{\Phi}^{-1}(w\alpha/m) - \xi_m)$  equal to  $1 - \beta$  implies  $w = (m/\alpha)\overline{\Phi}(z_{1-\beta} + z_{\alpha/m})$  which is equal to  $w_1$  as stated in the theorem. The stated form of  $w_0$  implies that the weights average to 1. The stated solution thus satisfies the restriction that a fraction  $\epsilon$  have power at least  $1 - \beta$ . Increasing the weight of any hypothesis whose weight is  $w_0$  necessitates reducing the weight of another hypothesis. This either reduces the minimum power of forces a hypothesis with power  $1 - \beta$  to fall below  $1 - \beta$ . Hence, the stated solution does in fact maximize the minimum power. ■

## References

- 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, 1, 289–300.
- (1997). “Multiple hypotheses testing with weights.” *Scand. J. Statist.*, 24, 3, 407–418.
- Benjamini, Y., Krieger, A. M., and Yekutieli, D. (2006). “Adaptive linear step-up procedures that control the false discovery rate.” *Biometrika*, 93, 3, 491–507.
- Benjamini, Y. and Yekutieli, D. (2001). “The control of the false discovery rate in multiple testing under dependency.” *Ann. Statist.*, 29, 4, 1165–1188.
- Bickel, P., Klaassen, C., Ritov, Y., and Wellner, J. (1993). “Efficient and Adaptive Estimation for Semiparametric Models.” Tech. rep., Johns Hopkins Series in the Mathematical Statistics, Baltimore, Maryland.
- Blanchard, G. and Roquain, E. (2008a). “Adaptive FDR control under independence and dependence.” *J. Mach. Learn. Res.*. To appear.
- (2008b). “Two simple sufficient conditions for FDR control.” *Electron. J. Stat.*, 2, 963–992.
- Chen, J. J., Lin, K. K., Huque, M., and Arani, R. B. (2000). “Weighted  $p$ -value adjustments for animal carcinogenicity trend test.” *Biometrics*, 56, 586–592.
- Donoho, D. and Jin, J. (2004). “Higher criticism for detecting sparse heterogeneous mixtures.” *Ann. Statist.*, 32, 3, 962–994.
- Efron, B. (2007). “Simultaneous inference: When should hypothesis testing problems be combined?” *Ann. Appl. Stat.*, 2, 197–223.
- Efron, B., Tibshirani, R., Storey, J. D., and Tusher, V. (2001). “Empirical Bayes analysis of a microarray experiment.” *J. Amer. Statist. Assoc.*, 96, 456, 1151–1160.

- Genovese, C. and Wasserman, L. (2002). “Operating characteristics and extensions of the false discovery rate procedure.” *J.R. Statist. Sec. B.*, 64, 499–517.
- Genovese, C. R., Roeder, K., and Wasserman, L. (2006). “False discovery control with  $p$ -value weighting.” *Biometrika*, 93, 3, 509–524.
- Hastie, T., Tibshirani, R., and Friedman, J. (2001). *The elements of statistical learning: data mining, inference, and prediction*. Springer.
- Holm, S. (1979). “A simple sequentially rejective multiple test procedure.” *Scand. J. Statist.*, 6, 2, 65–70.
- Ionita-Laza, I., McQueen, M., Laird, N., and Lange, C. (2007). “Genomewide weighted hypothesis testing in family-based association studies, with an application to a 100K scan.” *Am. J. Hum. Genet.*, 81, 607–614.
- Kropf, S., Läuter, J., Eszlinger, M., Krohn, K., and Paschke, R. (2004). “Nonparametric multiple test procedures with data-driven order of hypotheses and with weighted hypotheses.” *J. Statist. Plann. Inference*, 125, 1-2, 31–47.
- McQueen, M. and colleagues (2008). “Personal Communication.”
- Roeder, K., Bacanu, S.-A., Wasserman, L., and Devlin, B. (2006). “Using Linkage Genome Scans to Improve Power of Association in Genome Scans.” *Am. J. Hum. Genet.*, 78, 243–252.
- Roeder, K., Wasserman, L., and Devlin, B. (2007). “Improving power in genome-wide association studies: weights tip the scale.” *Genet. Epidemiol.*, 31, 741–747.
- Romano, J. P., Shaikh, A. M., and Wolf, M. (2008). “Control of the false discovery rate under dependence using the bootstrap and subsampling.” *TEST*, 17, 3, 417–442.
- Roquain, E. and van de Wiel, M. (2008). “Multi-weighting for FDR control.” *arXiv:0807.4081*, 78.

- Rosenthal, R. and Rubin, D. (1983). “Ensemble-adjusted p-values.” *Psychol. Bull.*, 94, 540–541.
- Rubin, D., Dudoit, S., and van der Laan, M. (2006). “A method to increase the power of multiple testing procedures through sample splitting.” *Stat. Appl. Genet. Mol. Biol.*, 5, Art. 19, 20 pp. (electronic).
- Sabatti, C., Service, S., and Freimer, N. (2003). “False discovery rate in linkage and association genome screens for complex disorders.” *Genetics*, 164, 829–33.
- Saccone, S., A.L., H., Saccone, N., Chase, G., Konvicka, K., Madden, P., Breslau, N., Johnson, E., Hatsukami, D., Pomerleau, O., Swan, G., Goate, A., Rutter, J., Bertelsen, S., Fox, L., Fugman, D., Martin, N., Montgomery, G., Wang, J., Ballinger, D., Rice, J., and Bierut, L. (2007). “Cholinergic nicotinic receptor genes implicated in a nicotine dependence association study targeting 348 candidate genes with 3713 SNPs.” *Hum Mol Genet.*, 16, 36–49.
- Sarkar, S. and Heller, R. (2008). “Comments on: Control of the false discovery rate under dependence using the bootstrap and subsampling.” *TEST*, 17, 450–455.
- Sarkar, S. K. (2002). “Some results on false discovery rate in stepwise multiple testing procedures.” *Ann. Statist.*, 30, 1, 239–257.
- Satagopan, J. and Elston, R. (2003). “Optimal two-stage genotyping in population-based association studies.” *Genet Epidemiol*, 25, 149–57.
- Schuster, E., Kropf, S., and Roeder, I. (2004). “Micro array based gene expression analysis using parametric multivariate tests per gene—a generalized application of multiple procedures with data-driven order of hypotheses.” *Biom. J.*, 46, 6, 687–698.
- Signoravitch, J. (2006). “Optimal multiple testing under the general linear model.” Tech. rep., Harvard Biostatistics.
- Skol, A., Scott, L., Abecasis, G., and Boehnke, M. (2006). “Joint analysis is more efficient than replication-based analysis for two-stage genome-wide association studies.” *Nat Genet.*, 38, 390–394.

- Spjøtvoll, E. (1972). “On the optimality of some multiple comparison procedures.” *Ann. Math. Statist.*, 43, 398–411.
- Storey, J. and Tibshirani, R. (2003). “Statistical significance for genome-wide studies.” In *Proceedings of the National Academy of Sciences*, vol. 100, 9440–9445.
- Storey, J. D. (2002). “A direct approach to false discovery rates.” *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 64, 3, 479–498.
- (2007). “The optimal discovery procedure: a new approach to simultaneous significance testing.” *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 69, 3, 347–368.
- Sun, W. and Cai, T. T. (2007). “Oracle and adaptive compound decision rules for false discovery rate control.” *J. Amer. Statist. Assoc.*, 102, 479, 901–912.
- Todd, J., Walker, N., Cooper, J., Smyth, D., K., D., Plagnol, V., Bailey, R., Nejentsev, S., Field, S., Payne, F., Lowe, C., Szeszko, J., Hafler, J., Zeitels, L., Yang, J., Vella, A., Nutland, S., Stevens, H., Schuilenburg, H., Coleman, G., Maisuria, M., Meadows, W., L.J., S., Healy, B., Burren, O., Lam, A., Ovington, N., Allen, J., Adlem, E., Leung, H., Wallace, C., Howson, J., Guja, C., Ionescu-Tirgovi, C., Genetics of Type 1 Diabetes in Finland, Simmonds, M., Heward, J., Gough, S., Wellcome Trust Case Control Consortium, Dunger, D., Wicker, L., and Clayton, D. (2007). “Robust associations of four new chromosome regions from genome-wide analyses of type 1 diabetes.” *Nat. Genet.*, 39, 857–864.
- Wang, H., Thomas, D., Pe’er, I., and Stram, D. (2006). “Optimal two-stage genotyping designs for genome-wide association scans.” *Genet Epidemiol*, 30, 4, 356–68.
- Wang, K., Li, M., and Bucan, M. (2007). “Pathway-Based Approaches for Analysis of Genomewide Association Studies.” *Am J Hum Genet*, 81, 1278–1283.
- Westfall, P., Krishen, A., and Young, S. (1998). “Using prior information to allocate significance levels for multiple endpoints.” *Statistics in Medicine*, 17, 2107–2119.

Westfall, P. H. and Krishen, A. (2001). “Optimally weighted, fixed sequence and gatekeeper multiple testing procedures.” *J. Statist. Plann. Inference*, 99, 1, 25–40.

Westfall, P. H., Kropf, S., and Finos, L. (2004). “Weighted FWE-controlling methods in high-dimensional situations.” In *Recent developments in multiple comparison procedures*, vol. 47 of *IMS Lecture Notes Monogr. Ser.*, 143–154. Beachwood, OH: Inst. Math. Statist.

Westfall, P. H. and Soper, K. A. (2001). “Using priors to improve multiple animal carcinogenicity tests.” *J. Amer. Statist. Assoc.*, 96, 455, 827–834.

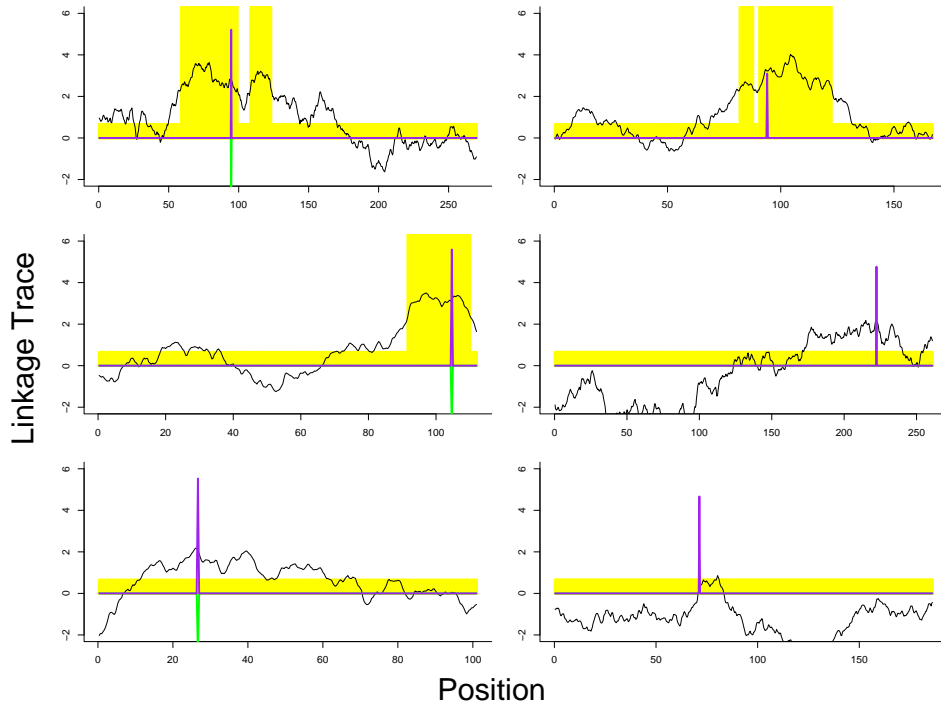


Figure 1: Linkage trace and weights for 6 chromosomes. The trace is the linkage statistic plotted as a function of position on the chromosome. The shading indicates which p-values were up/down weighted. The upspike is the association test statistic. The 3 downspikes indicate tests that were rejected using the binary weights.

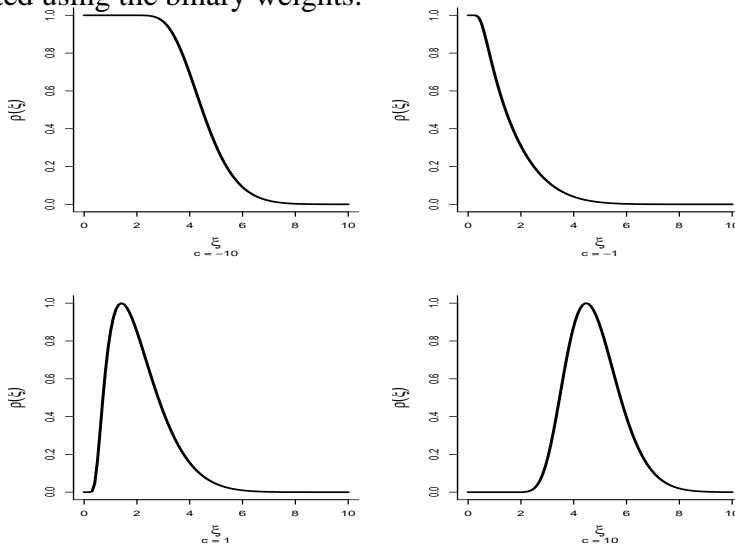


Figure 2: Optimal weight function  $\rho_c(\xi)$  for various  $c$ . In each case  $m = 1000$  and  $\alpha = 0.05$ . The functions are normalized to have maximum 1.

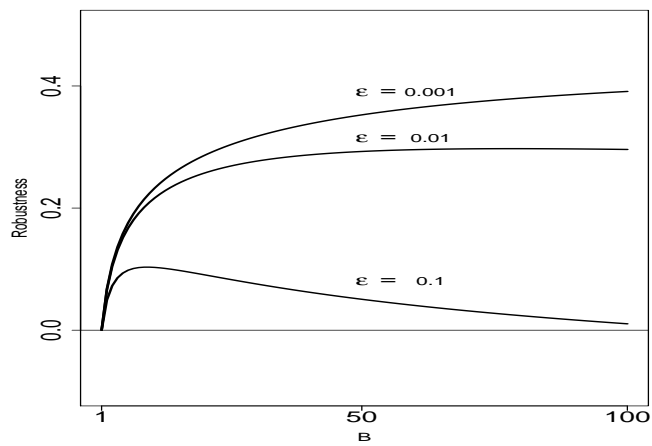


Figure 3: Robustness function for  $m = 1000$ . In this example,  $\xi = z_{\alpha/m}$  which has power  $1/2$  without weighting. The gain of correct weighting far outweighs the loss for incorrect weighting as long as the fraction of large weights  $\epsilon$  is small.

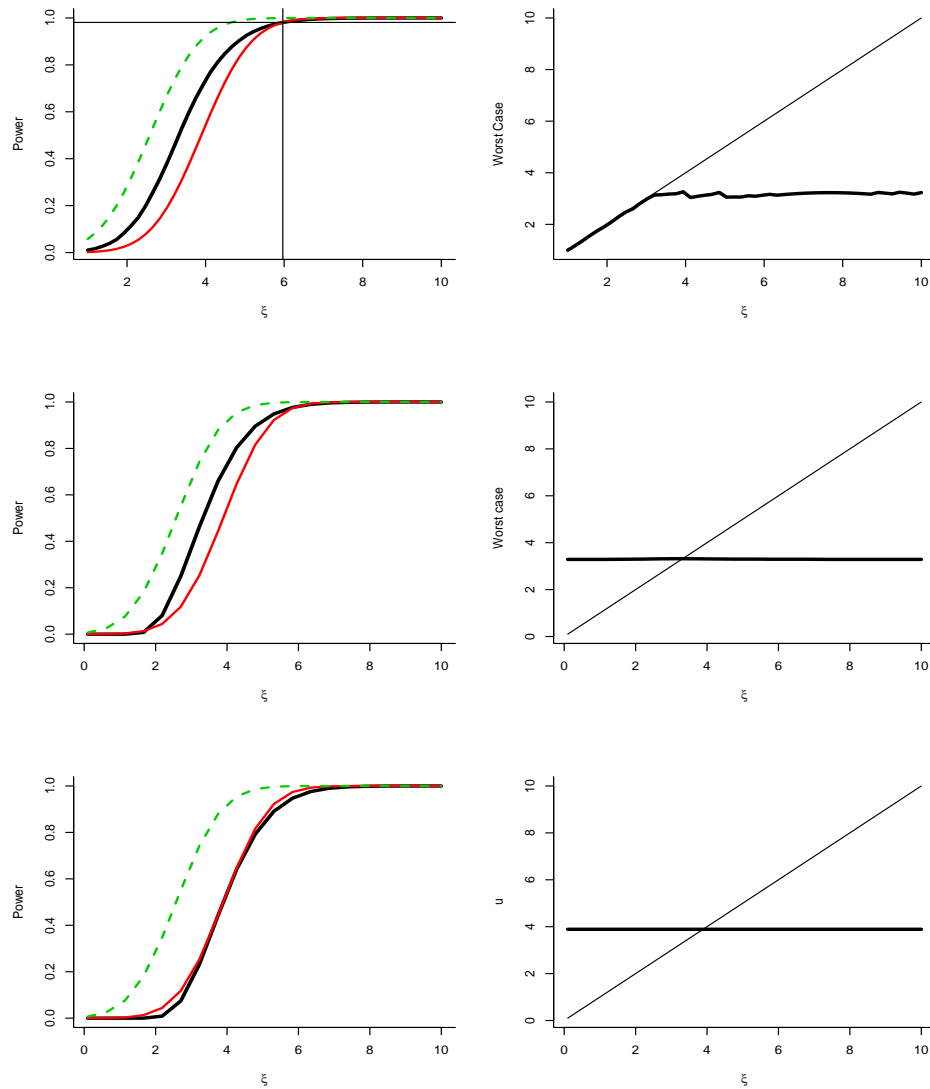


Figure 4: Power as a function of the alternative mean  $\xi$ . In these plots,  $a = .01$ ,  $m = 1000$  and  $\alpha = 0.05$ . There are  $(1 - a)m$  nulls and  $ma$  alternatives with mean  $\xi$ . The left plots shows what happens when the weights are incorrectly computed assuming that a fraction  $\gamma$  of nulls are actually alternatives with mean  $u$ . In the top plot, we restrict  $0 < u < \xi$ . In the second and third plot, no restriction is placed on  $u$ . The top and middle plot have  $\gamma = .1$  while the third plot has  $\gamma = 1 - a$  (all nulls misspecified as alternatives). The dark solid line shows the lowest possible power assuming the weights were estimated as poorly as possible. The lighter solid line is the power of the unweighted (Bonferroni) method. The dotted line is the power under the optimal weights. The vertical line in the top plot is at  $\xi_*$ . The weighted method beats unweighted for all  $\xi < \xi_*$ . The right plot shows the least favorable  $u$  as a function of  $\xi$ . That is, mistaking  $\gamma m$  nulls for alternatives with mean  $u$  leads to the worst power. Also shown is the line  $u = \xi$ .

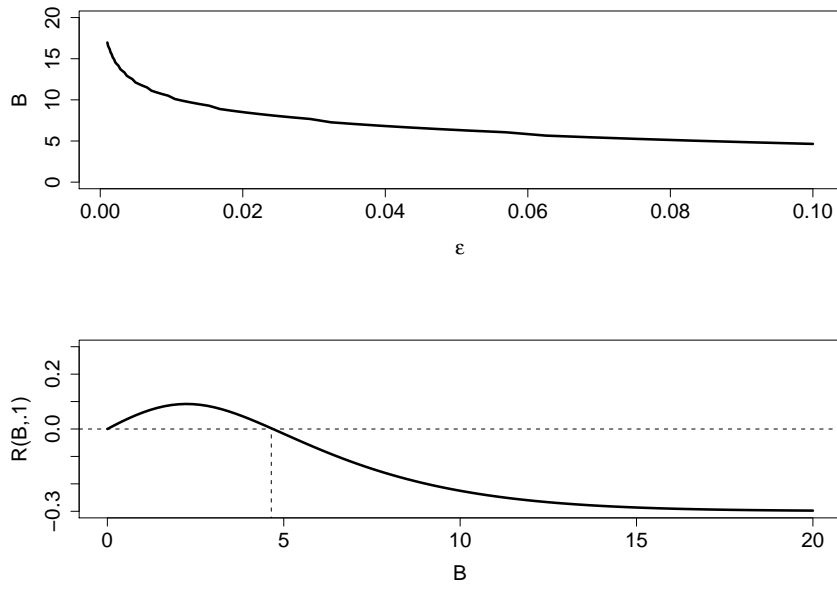


Figure 5: Top plot: turnaround point  $B_0(\epsilon)$  versus  $\epsilon$ . Bottom plot shows the robustness function  $R(B, .1)$  versus  $B$ . The turnaround point  $B_0(\epsilon)$  is shown with a vertical dotted line.

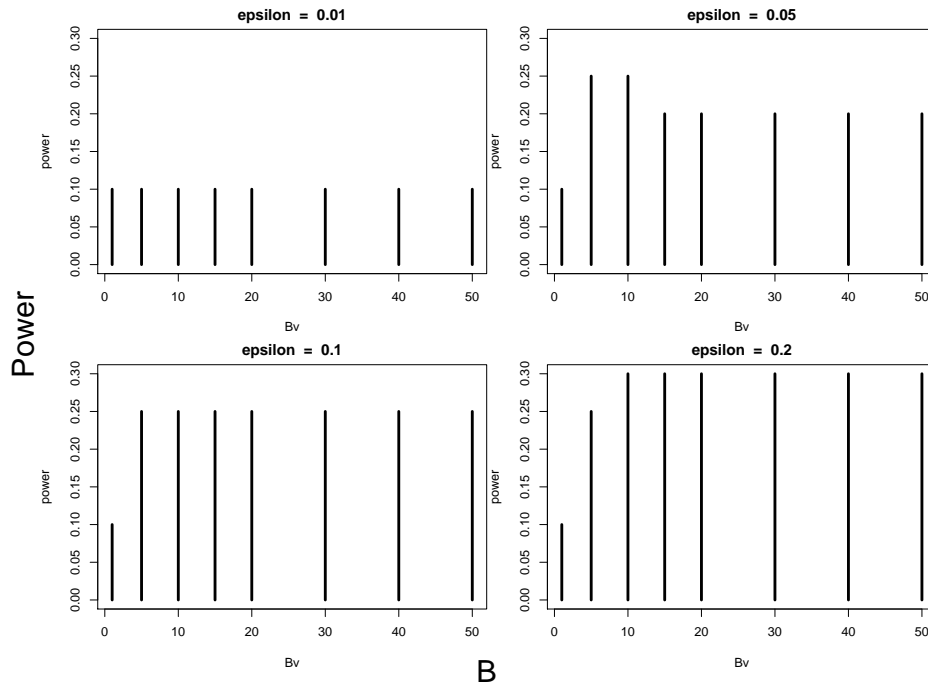


Figure 6: Power as a function of  $B$  and  $\epsilon$ .