

Large-Scale Simultaneous Hypothesis Testing: The Choice of a Null Hypothesis

Bradley Efron

Abstract

Current scientific techniques in genomics and image processing routinely produce hypothesis testing problems with hundreds or thousands of cases to consider simultaneously. This poses new difficulties for the statistician, but also opens new opportunities. In particular it allows empirical estimation of an appropriate null hypothesis. The empirical null may be considerably more dispersed than the usual theoretical null distribution that would be used for any one case considered separately. An empirical Bayes analysis plan for this situation is developed, using a local version of the false discovery rate to examine the inference issues. Two genomics problems are used as examples to show the importance of correctly choosing the null hypothesis.

Key Words: local false discovery rate, empirical Bayes, microarray analysis, empirical null hypothesis, unobserved covariates.

1. Introduction

Until recently “simultaneous inference” meant considering two or five or perhaps ten hypothesis tests at the same time, as in Miller’s classic 1981 text. Rapid progress in technology, particularly in genomics and imaging, has vastly upped the ante for simultaneous inference problems: now 500 or 5000 or even 50,000 tests may need to be evaluated at once, raising new problems for the statistician, but also opening new analytic opportunities. This paper concerns the choice of an appropriate null hypothesis in large-scale testing situations, and how this choice affects well-known inference methods such as the false discovery rate.

Simultaneous hypothesis testing begins with a collection of null hypotheses

$$H_1, H_2, \dots, H_N, \tag{1.1}$$

corresponding test statistics, possibly not independent,

$$Y_1, Y_2, \dots, Y_N, \tag{1.2}$$

and their p -values, P_1, P_2, \dots, P_N , with P_i measuring how strongly y_i , the observed value of Y_i , contradicts H_i , for instance $P_i = \text{prob}_{H_i}\{|Y_i| > |y_i|\}$. “Large-scale” means that N is a big number, say at least $N > 100$.

It is convenient though not necessary to work with z -values instead of the Y_i ’s or P_i ’s,

$$z_i = \Phi^{-1}(P_i), \quad i = 1, 2, \dots, N, \tag{1.3}$$

Φ indicating the standard normal cumulative distribution function (cdf), $\Phi^{-1}(.95) = 1.645$ etc. If H_i is exactly true then z_i will have a standard normal distribution

$$z_i|H_i \sim N(0, 1). \tag{1.4}$$

We will call (1.4) *the theoretical null hypothesis*.

Our motivating example concerns an HIV study of 1391 patients, investigating which of 6 Protease Inhibitor (“PI”) drugs cause mutations at which of 74 sites on the viral genome. Each patient provided a vector of predictors

$$\mathbf{x} = (x_1, x_2, \dots, x_6), \tag{1.5}$$

$x_j = 1$ or 0 indicating whether or not the patient used PI_j , $1 \leq \sum_1^6 x_j \leq 6$; and a vector of responses

$$\mathbf{v} = (v_1, v_2, \dots, v_{74}), \tag{1.6}$$

$v_k = 1$ or 0 indicating whether or not a mutation occurred at site k . Remark A of Section 7 describes the study in a little more detail.

For each of the 74 genomic sites, a separate logistic regression analysis was run using all 1391 cases, with that site's mutation indicators as responses and the PI indicators as predictors. Together these yielded $444 = 6 \times 74$ z -values, one for testing each null hypothesis, that drug j does not cause mutations at site k , $j = 1, 2, \dots, 6$ and $k = 1, 2, \dots, 74$. The z -values were based on the usual approximation

$$z_i = y_i/se_i, \quad i = 1, 2, \dots, 444, \quad (1.7)$$

(using a single subscript i in place of (j, k)) where y_i is the maximum likelihood estimate (MLE) of the logistic regression coefficient and se_i its approximate large-sample standard error.

Figure 1 shows a histogram of the 444 z -values, with negative z_i 's indicating greater mutational effects. The smooth curve $f(z)$ is a natural spline with seven degrees of freedom, fit to the histogram counts by Poisson regression. It emphasizes the *central peak* near $z = 0$, presumably the large majority of uninteresting drug-site combinations that have negligible mutation effects. Near its center the peak is well-described by a normal density having mean -0.35 and standard deviation 1.20 , which we will call the *empirical null hypothesis*,

$$z_i|H_i \sim N(-0.35, 1.20^2). \quad (1.8)$$

Section 3 describes the estimation methodology for (1.8), with a brief discussion of the normality assumption in Remark D of Section 7.

The difference between the theoretical null $N(0, 1)$ and empirical null $N(-0.35, 1.20^2)$ may not seem worrisome here but we will see that it substantially affects any simultaneous inference procedure. A more dramatic example is given in Section 6, for a microarray analysis where going from the theoretical to empirical null totally negates any findings of significance. Situations going in the reverse direction also occur.

In classic situations involving only a single hypothesis test we must, of necessity, employ the theoretical null hypothesis $z \sim N(0, 1)$. The main point of this paper is that large-scale testing situations permit empirical estimation of the null distribution. Sections 3 through 5 concern reasons why the empirical and theoretical null might differ, and which might be preferable in different situations.

There are scientific as well as statistical differences between small-scale and large-scale

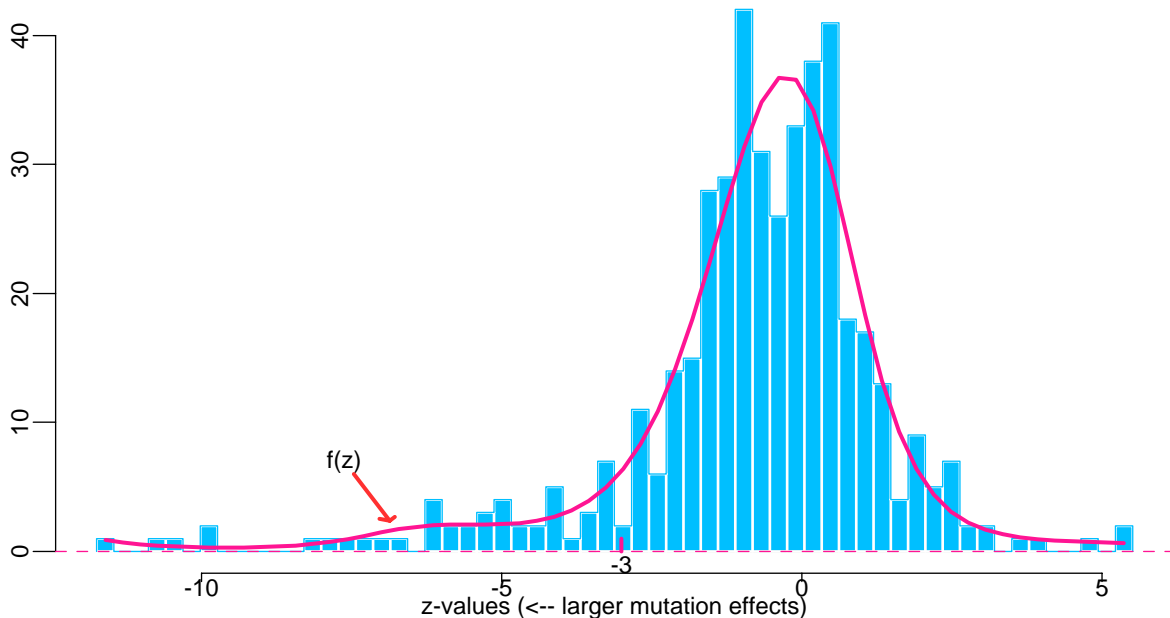


Figure 1: Histogram of 444 z -values from the Drug-mutation analysis; smooth curve $f(z)$ is natural spline fit to histogram counts. The central peak near $z = 0$ is approximately $N(-0.35, 1.20^2)$: the “empirical null hypothesis”. Simultaneous hypothesis tests for the 444 cases depend critically on the choice between the empirical or theoretical $N(0, 1)$ null.

hypothesis testing situations. A single hypothesis test is most often run with the expectation and hope of rejecting the null, “with 80% power” in a typical clinical trial. Nobody wants to reject 80% of $N = 5000$ null hypotheses. The usual point of large-scale testing is to identify a small percentage of interesting cases that deserve further investigation. While not exactly looking for a needle in a haystack, we don’t want the whole haystack either. An important assumption of what follows is that the proportion of interesting cases is small, perhaps 1%, or 5% of N , but not more than 10%. This is made explicit in Section 2, in our description of the local false discovery rate as an analytic tool for large-scale testing. There are situations where the 10% limit is irrelevant, for example, in constructing prediction models, but these lie outside our purpose here.

The terminology “Interesting/Uninteresting” used in this paper in preference to “Significant/Nonsignificant” is discussed near the end of Section 5. We conclude in Sections 7 and 8 with remarks, including most of the technical details, and a summary.

2. The Local False Discovery Rate It is convenient to discuss large-scale testing prob-

lems in terms of the local false discovery rate (fdr), an empirical Bayes version of Benjamini and Hochberg’s (1995) methodology focusing on densities rather than tail areas; see Efron et al. (2001) and Efron and Tibshirani (2002).

We begin with a simple Bayes model. Suppose that the N z -values fall into two classes, “Uninteresting” or “Interesting”, corresponding to whether or not z_i is generated according to the null hypothesis, with prior probabilities p_0 and $p_1 = 1 - p_0$, for the classes; and that z_i has density either $f_0(z)$ or $f_1(z)$ depending on its class,

$$\begin{aligned} p_0 &= \text{Prob}\{\text{Uninteresting}\}, & f_0(z) &\text{ density if Uninteresting (Null)} \\ p_1 &= \text{Prob}\{\text{Interesting}\}, & f_1(z) &\text{ density if Interesting (Non-Null)} . \end{aligned} \tag{2.1}$$

The smooth curve in Figure 1 estimates the *mixture density* $f(z)$,

$$f(z) = p_0 f_0(z) + p_1 f_1(z) . \tag{2.2}$$

According to Bayes theorem the *a posteriori* probability of being in the Uninteresting class given z is

$$\text{Prob}\{\text{Uninteresting}|z\} = p_0 f_0(z) / f(z) . \tag{2.3}$$

Here we define the *local false discovery rate* to be

$$\text{fdr}(z) \equiv f_0(z) / f(z) , \tag{2.4}$$

ignoring the factor p_0 in (2.3), so $\text{fdr}(z)$ is an upper bound on $\text{Prob}\{\text{Uninteresting}|z\}$. In fact p_0 can be roughly estimated, see Remark B, but we are assuming that p_0 is near 1, say $p_0 \geq 0.90$, so $\text{fdr}(z)$ is not a flagrant overestimator.

The local fdr provides a useful methodology for identifying Interesting cases in a situation like that of Figure 1: (1) estimate $f(z)$ from the observed ensemble of z -values, for example by the natural spline fit to the histogram counts; (2) assign a null density $f_0(z)$; (3) calculate $\text{fdr}(z) = f_0(z) / f(z)$; (4) report as Interesting those cases with $\text{fdr}(z_i)$ less than some threshold value, perhaps $\text{fdr}(z_i) \leq 0.10$. Remark B discusses the close connection between this algorithm and Benjamini and Hochberg’s (1995) method.

This paper concerns the choice of $f_0(z)$, the null hypothesis density. In the drug-mutation example it is crucial whether f_0 is taken to be the theoretical or empirical null, $N(0, 1)$ or $N(-0.35, 1.20^2)$. This is illustrated in Figure 2, a close-up view of Figure 1 focusing on the bin containing $z = -3$. The expected number of the 444 z_i values falling into this

bin is 6.37 for $f(z)$, and either 0.62 or 3.90 as $f_0(z)$ is $N(0, 1)$ or $N(-0.35, 1.20^2)$. Thus $\text{fdr}(z) = f_0(z)/f(z)$ at $z = -3$ is estimated to be either

$$\text{fdr}(-3) = \begin{cases} .097 \text{ using theoretical null } N(0, 1) \\ \text{or} \\ .612 \text{ using empirical null } N(-0.35, 1.20^2) . \end{cases} \quad (2.5)$$

In this bin, changing from the theoretical to empirical null changes our inferences from Interesting to definitely Uninteresting.

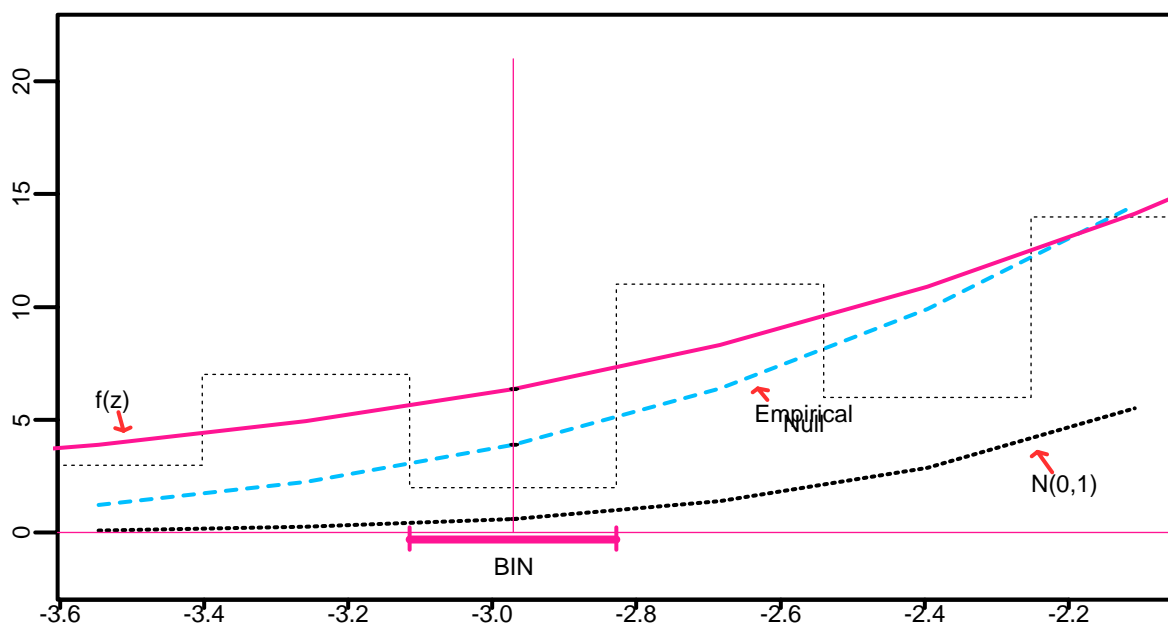


Figure 2: Close-up view of the bin containing $z = -3$ in Figure 1. Expected number in bin: 6.37 for $f(z)$, 0.62 for $f_0 = N(0, 1)$, 3.90 for $f_0 = N(-0.35, 1.20^2)$, the empirical null. Corresponding estimates of $\text{fdr}(-3)$: 0.097 for $N(0, 1)$ versus 0.612 for $N(-0.35, 1.20^2)$. Should we report the cases in this bin as Interesting?

Figure 3 compares the two estimates of $\log \text{fdr}(z)$ over most of the z scale. 18 of the 444 z -values have $\text{fdr}(z) < 0.10$ for $f_0 = N(0, 1)$ but > 0.10 for $f_0 = N(-0.35, 1.20^2)$, with 17 of these at the left end of the scale. All told the empirical null yields only two-thirds as many cases with $\text{fdr} < 0.10$ as the theoretical null, 35 compared to 53.

3. Estimating the Empirical Null Distribution Our estimate of the empirical null distribution for the Drug-mutation data was obtained in two steps: the curve $f(z)$ shown

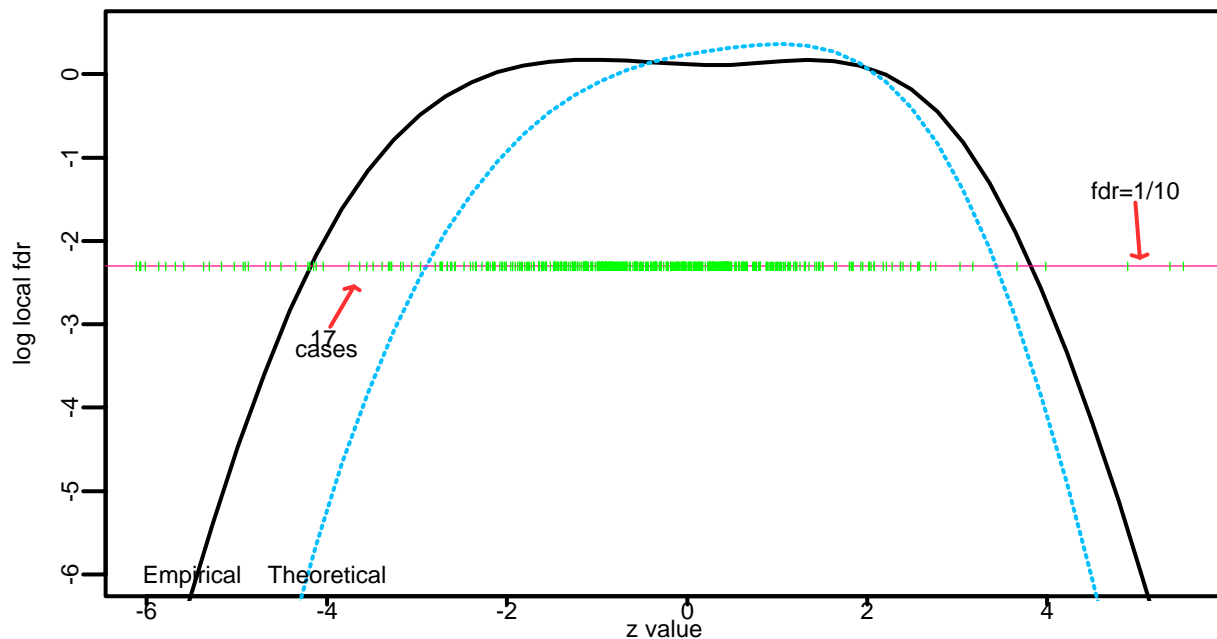


Figure 3: Comparison of estimates of $\log fdr(z)$ for the Drug-Mutation data; empirical null estimate (solid curve) declines more slowly than theoretical null estimate (dotted). Dashes indicate the 444 z -values. 17 cases on left have $fdr(z) < 1/10$ for theoretical but $> 1/10$ for empirical.

in Figure 1 was fit to the histogram counts by Poisson regression, and then the center and half-width of the central peak, say δ_0 and σ_0 , were obtained from $f(z)$,

$$\delta_0 = \arg \max\{f(z)\} \text{ and } \sigma_0 = \left[-\frac{d^2}{dz^2} \log f(z) \right]_{\delta_0}^{-\frac{1}{2}}, \quad (3.1)$$

yielding $(\delta_0, \sigma_0) = (-0.35, 1.20)$. Details are given in Remark D, where the possibility of a non-normal empirical null distribution is briefly discussed.

More direct estimation methods for f_0 seem possible, for example estimating δ_0 by the median of the z -values. Suppose though that 10% of the z -values came from the non-null distribution and all of these were located at the far left end of Figure 1. Then the median of all the z 's would be the 4/9 quartile of the actual null distribution, not its median, yielding a badly biased estimate of δ_0 . Similar comments apply to estimating σ_0 , Remark D. Method (3.1) does not require preliminary estimates of the proportion p_0 in the null population of (2.1), a considerable practical advantage.

How accurate are the estimates $(-0.35, 1.20)$? The usual standard error approximations for a Poisson regression fit are not appropriate here since the z_i 's are not independent of

each other. A nonparametric bootstrap analysis was performed instead, with the 1391 80-dimensional vectors (\mathbf{x}, \mathbf{v}) , (1.5-1.6) as the resampling units. This gave .09 and .08 for the bootstrap standard errors of δ_0 and σ_0 respectively, i.e.

$$(\delta_0, \sigma_0) = (-0.35, 1.20) \pm (.09, .08) . \quad (3.2)$$

It seems quite unlikely that estimation error alone accounts for the difference between the empirical null and the theoretical values $(\delta_0, \sigma_0) = (0, 1)$. (Notice that this type of bootstrap analysis, which requires independent sampling units, is not applicable to the microarray example of Section 6, where we expect correlations among the genes.)

The next two sections concern other possible causes for empirical/theoretical differences, diagnostics for these causes, and their interpretations. Our list is not exhaustive and in fact the microarray example of Section 6 demonstrates another form of pathology.

4. Permutation Tests and Unobserved Covariates The theoretical $N(0, 1)$ null hypothesis (1.4) is usually based on asymptotic approximations like those for the logistic regression coefficients in the Drug-mutation study. Permutation methods can be used to avoid these approximations, perhaps in the hope that an improved theoretical null will more closely match the empirical.

This was not the case for the Drug-mutation data. Permutation testing was implemented by randomly pairing the 1391 predictor vectors \mathbf{x} , (1.5), with the 1391 response vectors \mathbf{v} , (1.6), and recalculating the 444 z -values. This whole process was independently repeated 20 times, yielding a total of 20×444 permutation z 's. Their distribution was well approximated by a $N(0, .965^2)$ density (the ‘‘permutation null’’) except for a prominent spike near $z = 0.3$. In this case the permutation-improved theoretical null differs more rather than less from the empirical null $N(-0.35, 1.20^2)$.

Permutation methods are popular in the microarray literature as a way of avoiding assumptions and approximations, see Efron et al. (2001) or Dudoit et al. (2003), *but they do not automatically resolve the question of an appropriate null hypothesis*. This can be seen in the following hypothetical example, which is a stylized version of the two-sample microarray testing problem in Section 6: the data x_{ij} comes from N simultaneous two-sample experiments, each comparing $2n$ subjects,

$$x_{ij} \begin{cases} \text{Controls} & j = 1, 2, \dots, n \\ \text{Treatments} & j = n + 1, n + 2, \dots, 2n \end{cases} \quad (i = 1, \dots, N), \quad (4.1)$$

The i^{th} test statistic Y_i is the usual two-sample t -statistic, comparing Treatments versus Controls for the i^{th} experiment.

Suppose that, unknown to the statistician, the data was actually generated from

$$x_{ij} = u_{ij} + \frac{I_j}{2}\beta_i \begin{cases} u_{ij} \sim N(0, 1) \\ \beta_i \sim N(0, \sigma_\beta^2) \end{cases}, \quad (4.2)$$

with the u_{ij} and β_i mutually independent and

$$I_j = \begin{cases} -1 & j = 1, 2, \dots, n \\ +1 & j = n + 1, \dots, 2n \end{cases}. \quad (4.3)$$

Then it is easy to show that the statistics Y_i follow a dilated t -distribution with $2n - 2$ degrees of freedom,

$$Y_i \sim \left(1 + \frac{n}{2}\sigma_\beta^2\right)^{\frac{1}{2}} \cdot t_{2n-2}, \quad (4.4)$$

while the permutation distribution, permuting Treatments and Controls within each experiment, has nearly a standard t_{2n-2} null distribution. So for example if $\sigma_\beta^2 = 2/n$, the empirical density of the Y_i 's will be $\sqrt{2}$ times as wide as the permutation null.

The quantity β_i in (4.2)-(4.3) causes the only consistent differences between Treatments and Controls in experiment i . If β_i is a dependable feature of the i^{th} experiment, and would appear again with the same value in a replication of the study, then the permutation null t_{2n-2} is a reasonable basis for inference. With n large and $\sigma_\beta^2 = 2/n$, it results in $\text{fdr}(y_i) < 0.10$ for the most extreme 2% of the observed t -statistics, favoring those with the largest values of $|\beta_i|$.

Suppose though that β_i is not inherent to experiment i , but rather a purely random effect that would have a different value and perhaps a different sign if the study were repeated; that is, β_i is part of the noise and not part of the signal. In this case the appropriate choice is the empirical null (4.4). The equivalent of Figure 1 will be *all* central peak, with no interesting outliers, and there will be no cases having small values of $\text{fdr}(y_i)$. This is appropriate since now there is no real Treatment effect.

In this last context β_i acts as an *unobserved covariate*, a quantity which the statistician would use to correct the Treatment-Control comparison if it were observable. Unobserved covariates are ubiquitous in observational studies. There are several obvious ones in the Drug-mutation study: personal characteristics of the patients such as age and gender, prior use of AZT and other non-PI drugs, years since infection, geographical location, etc.

The effect of important unobserved covariates is to dilate the null hypothesis density $f_0(z)$, as happens in (4.4). Unobserved covariates will also dilate the “Interesting” density $f_1(z)$ in (2.1), and the mixture density $f(z)$, (2.2). However an empirical fitting method for estimating $f(z)$, such as the spline fit in Figure 1, automatically includes any dilation effects. In estimating $\text{fdr}(z) = f_0(z)/f(z)$ it is important to also allow for dilation of the numerator f_0 . *This is a strong argument for preferring the empirical null hypothesis in observational studies.*

5. A Structural Model for the z-values The Bayesian specifications (2.1) underlying our fdr results have the advantage of not requiring a structural model for the z -values; in particular it is not necessary to motivate, or even describe, the non-null density $f_1(z)$. There is however a simple structural model that helps elucidate the Interesting-Uninteresting distinction in (2.1).

The structural model assumes that z_i , the i^{th} z -value, is normally distributed around a “true value” μ_i , its expectation,

$$z_i \sim N(\mu_i, 1) \quad \text{for } i = 1, 2, \dots, N, \quad (5.1)$$

with μ_i having some prior distribution $g(\mu)$,

$$\mu_i \sim g(\mu) \quad \text{for } i = 1, 2, \dots, N. \quad (5.2)$$

Structure (5.1) is often a good approximation, see Section 4 of Efron (1988), and in fact proved reasonably accurate in the bootstrap experiment giving (3.2). Together (5.1)-(5.2) say that the mixture density $f(z)$, (2.2), is a convolution of $g(\mu)$ with the standard normal density $\varphi(z)$,

$$f(z) = \int_{-\infty}^{\infty} \varphi(z - \mu)g(\mu)d\mu \quad (5.3)$$

(with the understanding that $g(\mu)$ may include discrete probability atoms.)

As a first application of the structural model, suppose we insist that $g(\mu)$ put probability p_0 on $\mu = 0$,

$$\text{Prob}_g\{\mu = 0\} = p_0 \quad (5.4)$$

for some fixed value of p_0 between 0 and 1. This amounts to our original Bayes model (2.1) with $p_0 = \text{Prob}\{\text{Uninteresting}\}$, $f_0(z)$ the theoretical null hypothesis $N(0, 1)$, and

$$f_1(z) = \int_{\mu \neq 0} \varphi(z - \mu)g(\mu)d\mu / (1 - p_0). \quad (5.5)$$

In the context of this paper, p_0 should be 0.90 or greater.

For any $f(z)$ of the convolution form (5.3) let (δ_g, σ_g) be the center and width parameters (δ_0, σ_0) defined by (3.1). Figure 4 answers the following question: for a given choice of p_0 in constraint (5.4), what are the maximum possible values of $|\delta_g|$ and of σ_g ,

$$\delta_{\max} = \max\{|\delta_g| | p_0\} \quad \text{and} \quad \sigma_{\max} = \max\{\sigma_g | p_0\} . \quad (5.6)$$

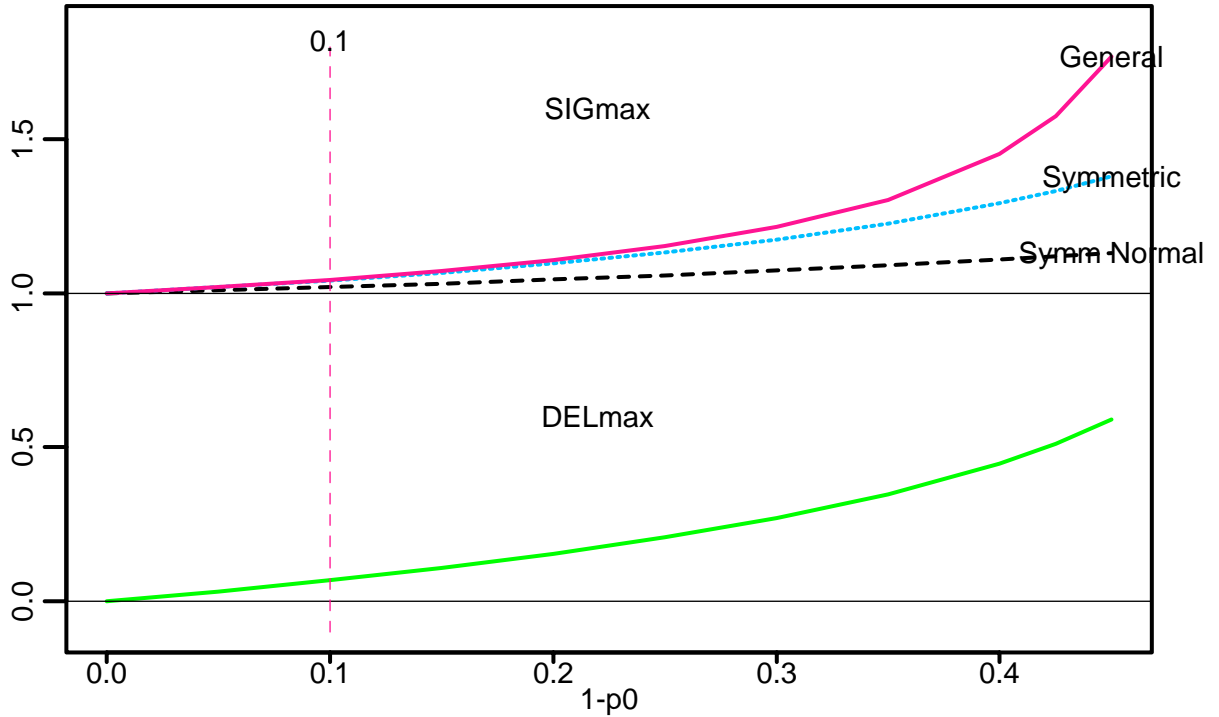


Figure 4: Maximum possible values of the center and width parameters (δ_0, σ_0) , (3.1), when the structural model (5.1)-(5.3) is constrained to put probability p_0 on $\mu = 0$. For $1-p_0 \leq 0.10$ the maxima are not much greater than the theoretical null values $(0, 1)$, as shown in Table 1.

Three curves appear for σ_{\max} , for the general case just described, for the case where the non-zero component of $g(\mu)$ is required to be symmetric around zero, and for the case where it is also required to be normal. Here we will only mention the general case. Remark F discusses the solution of (5.6), which turns out to have a simple “single-point” form.

The notable feature of Figure 4 is that for $p_0 \geq 0.90$, our preferred realm for large-scale hypothesis testing, $(\delta_{\max}, \sigma_{\max})$ must be quite near the theoretical null values $(0, 1)$:

$$\delta_{\max} \leq 0.07 \quad \text{and} \quad \sigma_{\max} \leq 1.04 . \quad (5.7)$$

Table 1 shows $(\delta_{\max}, \sigma_{\max})$ for various choices of p_0 . We see that the “Interesting” probability $1 - p_0$ would have to be nearly 0.30, very big by the standards of large-scale testing, in order to obtain the observed Drug-mutation values $(\delta_0, \sigma_0) = (-0.35, 1.20)$. The inference is that uninteresting effect, such as the unobserved covariates of Section 4, are dilating the null hypothesis.

Table 1: Value of σ_{\max} and δ_{\max} as a function of $1 - p_0$, (5.4).

$1 - p_0$:	0.05	0.10	0.20	0.30	(Drug-mutation)
σ_{\max} :	1.02	1.04	1.11	1.22	(1.20)
δ_{\max} :	0.03	0.07	0.15	0.27	(-0.35)

The main point here is that our measures (3.1) of center and width are quite robust to the arrangement of Interesting values μ_i as long as the Interesting percentage does not exceed 10%. If (δ_0, σ_0) for the central peak is much different than $(0, 1)$, as it is in Figure 1, then use of the theoretical null is bound to result in identifying an uncomfortably large percentage of supposedly Interesting cases.

We can pursue this last point for the Drug-mutation data by removing constraint (5.4). Figure 5 shows an unconstrained estimate of $g(\mu)$. For computational simplicity $g(\mu)$ was assumed to be discrete, with at most $J = 8$ support points $\mu_1, \mu_2, \dots, \mu_J$, so that (5.3) becomes

$$f(z) = \sum_{j=1}^J \pi_j \varphi(z - \mu_j), \quad (5.8)$$

π_j being the probability g puts on μ_j , with $\pi_j \geq 0$ and $\sum \pi_j = 1$. A non-linear minimization program was employed to find the best-fit curve of form (5.8) to the histogram counts in Figure 1, using Poisson deviance as the fitting criterion. The vertical bars in Figure 5 are located at the resulting 8 values μ_j , with the bar’s height proportional to π_j . For example the little bar at far left represents an atom of probability $\pi_1 = .015$ at $\mu_1 = -10.9$. The resulting $f(z)$ estimate (5.7) closely resembles the natural spline fit of Figure 1. Table 2 shows all 8 (π_j, μ_j) pairs.

Suppose for a moment that the estimated $g(\mu)$ is exactly correct, so 1.5% of the 444 cases have their μ_i ’s equal -10.9, 1.3% have -7.0, etc., and that an oracle has told us the eight (π_j, μ_j) values. Given an observed z_i we can now calculate $\text{Prob}\{\text{Uninteresting}|z\}$,

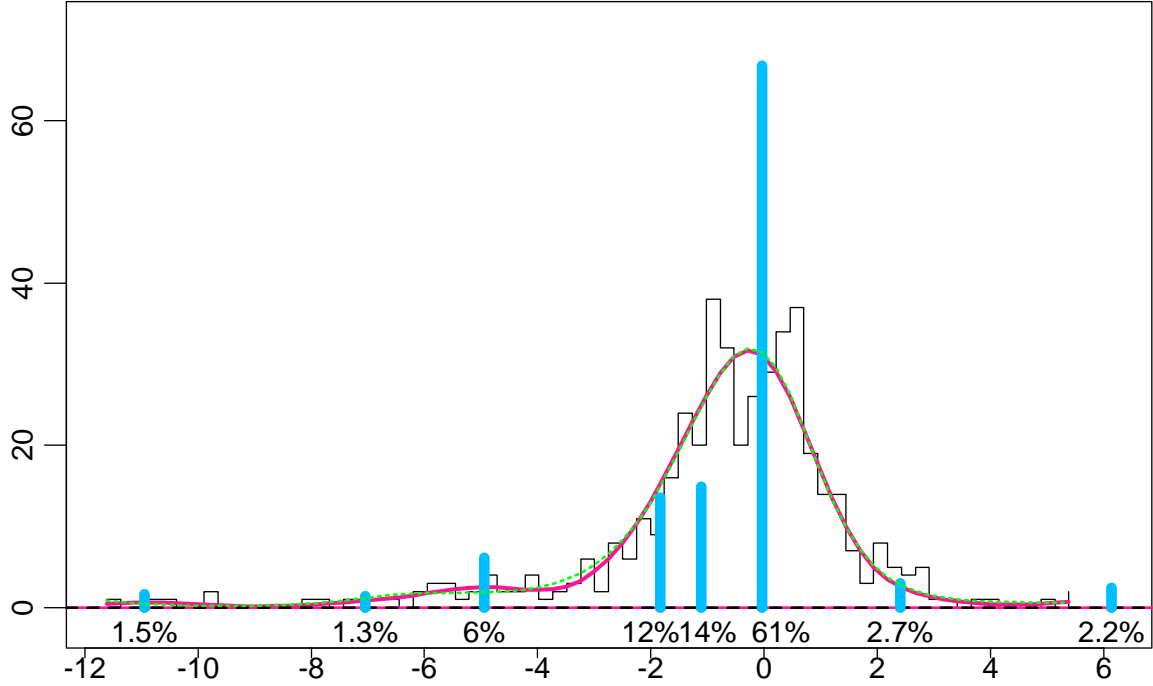


Figure 5: Best-fit discrete mixing function $g(\mu)$, (5.2) for Drug-mutation data; bars located at support points μ_j , heights proportional to weights π_j ; tall bar at $\mu_j = 0$ has weight $\pi_j = 0.61$. Solid curve is best-fit estimate $f(z) = \sum \pi_j \varphi(z - \mu_j)$; it closely matches natural spline fit from Figure 1 (dashed curve).

(2.3), exactly, once the scientist specifies the definition of Uninteresting versus Interesting. It seems obvious that the 60.8% at $\mu_j = 0$ are Uninteresting, and that the 10.6% at $\mu_j = -10.9, -7.0, -4.9$, and 6.1 deserve Interesting status. However the status of the 28.6% at $\mu_j = -1.8, -1.1$, and 2.4 is less clear.

If the 28.6% are deemed Interesting, this leaves only the 60.8% at $\mu_j = 0$ as Uninteresting. In terms of our Bayes model (2.1) we have $p_0 = .608$ and $f_0(z) \sim N(0, 1)$, the theoretical null. About 174 of the 444 cases will be identified as Interesting, too many for a typical screening exercise. Shifting the 28.6% to the Uninteresting classification increases p_0 to $.608 + .286 = .894$, a more manageable value, and changes $f_0(z)$ to the version of (5.7) supported on the four Uninteresting μ_j 's,

$$f_0(z) = \frac{\sum_{j=4}^7 \pi_j \varphi(z - \mu_j)}{\sum_{j=4}^7 \pi_j}; \quad (5.9)$$

this is approximately $N(-0.34, 1.19^2)$, almost the same as the empirical null (1.8).

In other words the definition of “Interesting” determines the relevant choice of the null hypothesis f_0 . If we want to keep the proportion of Interesting cases manageably small then $f_0(z)$ has to grow wider than $N(0, 1)$.

Use of the term “Interesting” rather than “Significant” reflects a difference in intent between large-scale and classical testing. In the hypothetical context of Figure 5 and Table 2, all of the 39.2% of the cases with non-zero μ_i ’s would eventually be declared as “significantly different from zero” if we vastly increased the sample size of patients. Section 4 suggests that minor deviations from $N(0, 1)$ might arise from scientifically uninteresting causes such as unobserved covariates. However even if a modestly non-zero μ_i is genuine in some sense, it may still be Uninteresting when viewed in comparison with an ensemble of more dramatic possibilities. Nonsignificant implies Uninteresting but not conversely.

6. A Microarray Example Microarrays have become a prime source of large-scale simultaneous testing problems. Figure 6 relates to a well-known microarray experiment concerning differences between two types of genetic mutations causing increased breast cancer risk, “BRCA1” and “BRCA2”; see Hedenfalk et al. (2001), also Efron and Tibshirani (2002), and Efron (2003).

The experiment included 15 breast cancer patients, seven with the BRCA1 mutation and eight with BRCA2. Each women’s tumor was analyzed on a separate microarray, each microarray reporting on the same set of $N = 3226$ genes. For each gene the two-sample t -statistic y_i comparing the 7 BRCA1 responses with the 8 BRCA2’s was computed. The y_i ’s were then converted to z -values.

$$z_i = \Phi^{-1}F_{13}(y_i) , \tag{6.1}$$

where F_{13} is the cdf of a standard t -distribution with 13 degrees of freedom. Figure 6 displays the histogram of the 3226 z -values.

Table 2:Weights π_j and locations μ_j for 8-point best-fit estimate $g(\mu)$ of Figure 8. Which locations we deem Interesting versus Uninteresting determines the choice between the theoretical or empirical null hypothesis. (Numerical results accurate to one decimal place.)

	<i>–Interesting–</i>		?	?	<i>Uninteresting</i>	?	<i>Interesting</i>	
$100 \cdot \pi_j$:	1.5%	1.3%	5.6%	12.3%	13.6%	60.8%	2.7%	2.2%
μ_j :	-10.9	-7.0	-4.9	-1.8	-1.1	0.0	2.4	6.1

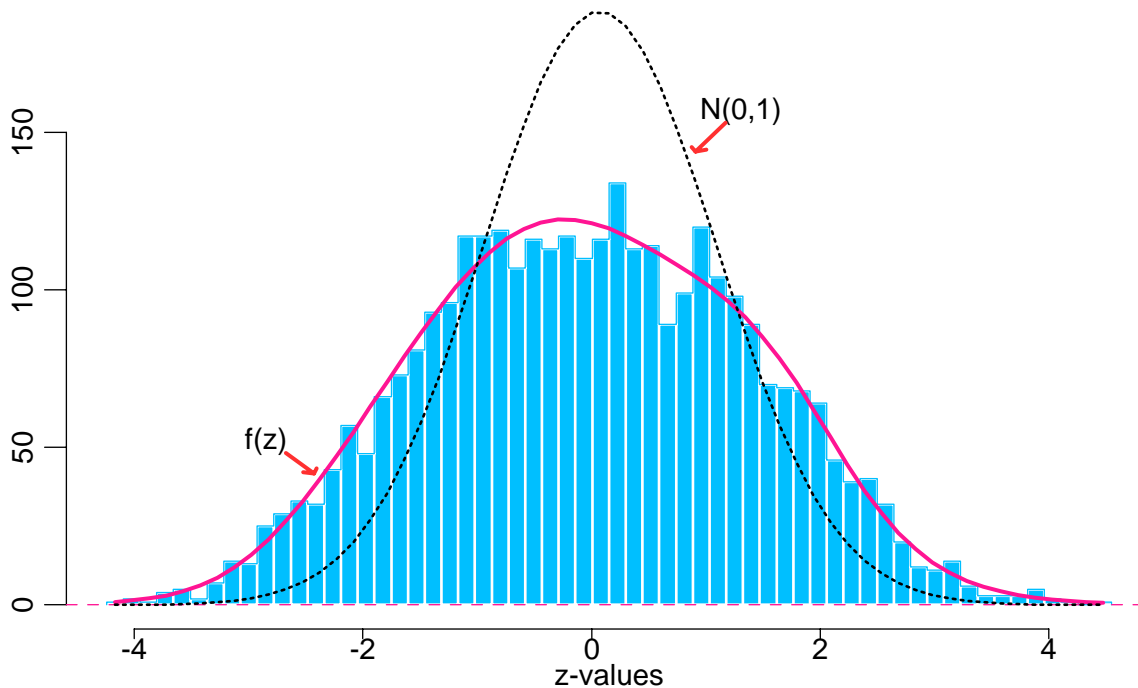


Figure 6: Histogram of $N = 3226$ z -values from breast cancer study. The theoretical $N(0, 1)$ null is much narrower than the central peak, which has $(\delta_0, \sigma_0) = (-0.02, 1.58)$. In this case the central peak seems to include the entire histogram.

The central peak is wider here than in Figure 1, with center-width estimates $(\delta_0, \sigma_0) = (-0.02, 1.58)$. More importantly, the histogram seems to be *all* central peak, with no interesting outliers such as those seen at the left of Figure 1. This was reflected in the local fdr calculations: using the theoretical $N(0, 1)$ null yielded 35 genes having $fdr(z_i) < 0.1$, those with $|z_i| > 3.35$; using the empirical $N(-0.02, 1.58^2)$ null, *no genes at all* had $fdr < 0.1$ (or for that matter $fdr < 0.9$, the histogram in fact being a little short-tailed compared to $N(-0.02, 1.58^2)$.)

There is ample reason to distrust the theoretical null in this case. The microarray experiment for all its impressive technology is still an observational study, with a wide range of unobserved covariates possibly distorting the BRCA1-BRCA2 comparison.

Another reason for doubt can be found in the data itself. The fdr methodology does not require independence of the y_i 's or z_i 's across genes. However it does require that the 15 measurements for *each* gene be independent across the microarrays. Otherwise the two-sample t -statistic y_i will not have an F_{13} null distribution, not even approximately.

Unfortunately the experimental methodology used in the breast cancer study seems to have induced substantial correlations among the various microarrays. In particular, as discussed in Remark G, the first four microarrays in the BRCA2 groups were mutually correlated, and likewise the last four. Correlations reduce the effective sample size for a two-sample t -statistic, just the type of effect that would induce overdispersion in (6.1).

This does not say that there are no BRCA1-BRCA2 differences, only that it is dangerous to compare the t -statistics with a standard t_{13} null distribution, even if simultaneous inference is accounted for.

7. Remarks

A. Drug-mutation Study The data base for the Drug-mutation study, Wu et al. (2002), included 2497 patients having HIV subtype B, of whom 1391 had received at least one of six popular Protease Inhibitor drugs. amprenavir, indinavir, lopinavir, nelfinavir, ritonavir, or saquinavir. Among the 1391, the mean number of PI drugs taken was 2.05 per patient. Amino acid sequences were obtained at all 99 positions on the HIV protease gene, and mutations from wild-type recorded; 25 positions showed 3 or fewer mutations among the 1391 patients, deemed too few for analysis, leaving 74 positions for the investigation here. Each of the 74 individual logistic regressions included an intercept term as well as the six PI main effects, but no other covariates.

B. The Local False Discovery Rate The local fdr, (2.3) or (2.4), is closely related to Benjamini and Hochberg's (1995) "tail-area" False Discovery Rate, as discussed in Efron et al. (2001) and Efron and Tibshirani (2002). Substituting cdf's F_0 and F for the densities f_0 and f , Bayes theorem gives a tail-area version of (2.3),

$$\text{Prob}\{\text{Uninteresting}|z \leq z_0\} = p_0 F_0(z_0)/F(z_0) \equiv \text{FDR}(z_0) . \quad (7.1)$$

$\text{FDR}(z_0)$ turns out to be the conditional expectation of $\text{fdr}(z) \equiv p_0 f_0(z)/f(z)$ given $z \leq z_0$,

$$\text{FDR}(z_0) = \int_{-\infty}^{z_0} \text{fdr}(z) f(z) dz / \int_{-\infty}^{z_0} f(z) dz . \quad (7.2)$$

Benjamini and Hochberg work in a frequentist framework but their False Discovery Rate control rule can be stated in empirical Bayes terms: given F_0 , which they usually take to be what we called the theoretical null, estimate $\text{FDR}(z_0)$ by

$$\widehat{\text{FDR}}(z_0) = p_0 F_0(z)/\widehat{F}(z_0) , \quad (7.3)$$

where \widehat{F} is the empirical cdf of the z_i 's; for a desired control level α , say $\alpha = .05$, define

$$z_0 = \arg \max_z \{\widehat{\text{FDR}}(z) \leq \alpha\} ; \quad (7.4)$$

then rejecting all cases with $z_i \leq z_0$ gives an expected (frequentist) rate of false discoveries no greater than α .

With z_0 as in (7.4), relation (7.2) (applied to the estimated versions of FDR, fdr, and f) says that the weighted average of $\text{fdr}(z_i)$ for the cases rejected by the FDR level- α rule is itself α . As an example take $\alpha = .05$ and f_0 equal the theoretical $N(0, 1)$ null. Applying the FDR control rule to the negative side of Figure 1's Drug-mutation data rejects the null hypothesis for the 56 cases having $z_i \leq -2.61$; the corresponding 56 values of $\text{fdr}(z_i)$ have weighted average $\alpha = .05$. They vary from nearly zero at the far left to .19 at the boundary value $z = -2.61$, justifying the name "local": z_i 's near the boundary are more likely to be false discoveries than the overall .05 rate suggests.

Our concern with a correct choice of null hypothesis applies to FDR just as well as fdr. In the microarray study, FDR with $F_0 = N(0, 1)$ gives 24 significant genes at $\alpha = .05$, while $F_0 = N(-.02, 1.58^2)$ gives none. In fact any simultaneous testing procedure, the popular Westfall-Young method (1993) for example, will depend on a correct assessment of p -values for the individual cases, i.e. on the choice of F_0 .

C. Estimating $f(z)$ The Poisson regression method used in Figure 1 to estimate the mixture density $f(z)$, (2.2), originates in an idea of Lindsey described in Section 2 of Efron and Tibshirani (1996): the range of the sample z_1, z_2, \dots, z_N is partitioned into K equal intervals, with interval k having midpoint x_k and containing count s_k of the N z -values; the expectation λ_k of s_k is nearly proportional to $f_k \equiv f(x_k)$, and if the z_i 's are independent the counts approximate independent Poisson variates,

$$s_k \overset{\text{ind}}{\sim} P_o(\lambda_k), \quad \lambda_k = cf_k \quad [k = 1, 2, \dots, K], \quad (7.5)$$

c a constant depending on N and the interval length.

Lindsey's method is to estimate the λ_k 's with a Poisson regression, which because of (7.5) amounts to estimating a scaled version of the f_k 's; in other words estimating $f(z)$. K equals 60 in Figure 1, with the regression model being a natural spline with 7 degrees of freedom, roughly equivalent to a sixth degree polynomial fit in z .

Poisson regression based on (7.5), is almost fully efficient for estimating $f(z)$ if the z_i 's are independent. Here we do not expect independence but we still have the expectation of

s_k proportional to f_k . The Poisson regression method will still tend to unbiasedly estimate $f(z)$, assuming the regression model is sufficiently flexible, though we may lose estimating efficiency.

The bootstrap analysis that gave the standard errors in (3.2) was also used to check (7.5). This turned out to be surprisingly accurate for the Drug-mutation data. If not we might have used the bootstrap estimate of covariance for the s_k 's to motivate a more efficient estimation procedure, though this is unlikely to be important for large values of N . In any case bootstrap analyses as in (3.2) will provide legitimate standard errors for the Poisson regression whether or not (7.5) is valid.

D. Estimating the Empirical Null Distribution The main tactic of this paper is to estimate the null distribution $f_0(x)$ in (2.1) from the central peak in the z -values' histogram. Assuming normality for f_0 gives

$$\log f(z) \doteq -\frac{1}{2} \left(\frac{z - \delta_0}{\sigma_0} \right)^2 + \text{constant} \quad (7.6)$$

for z near zero, so that δ_0 and σ_0 can be estimated by differentiating $\log f(z)$ as in (3.1). The constant depends on N and p_0 but the constant has no effect on the derivatives of (3.1).

Directly differentiating the spline estimate of $\log f(z)$ can give an overly variable estimate of σ_0 . One more smoothing step was employed here: a quadratic curve $a_0 + a_1x_k + a_2x_k^2$ was fit by ordinary least squares to the estimated values $\log f_k$, for x_k within 1.5 units of the maximum δ_0 , yielding $\sigma_0 = [-2a_2]^{-\frac{1}{2}}$ as in (3.1). This procedure gave the small bootstrap standard error estimate in (3.2).

None of this methodology is crucial, though it is important that the estimates δ_0 and σ_0 relate directly to $f_0(z)$, and are not much affected by the non-Null distribution $f_1(z)$ in (2.1). As an example of what can go wrong suppose we try to estimate σ_0 by a "robust" scale measure such as (84th quantile minus 16th quantile)/2. This gives $\sigma_0 = 1.47$ for the Drug-mutation data, reflecting long tails due to the Interesting cases in Figure 1. Similar difficulties arise using the central slope of a qq plot. Basically a density estimate of the central peak is required, and then some assessment of its center and width.

More ambitiously, we might try extending the estimation of $f_0(z)$ to third moments, permitting a skew null distribution. Expression (7.6) could be generalized to

$$-\log f(z) \doteq c_0 + c_1z + c_2z^2/2 + c_3z^3/6, \quad (7.7)$$

now requiring three derivatives to estimate the coefficients rather than the two of (3.1). This

is an unexplored path, and in particular Table 1 has not been extended to include skewness bounds.

Familiarity was the only reason for using z -values instead of t -values in Figures 1 and 6.

E. Estimating p_0 We can obtain reasonable upper bounds for p_0 in (2.1) from estimates of

$$\pi(c) \equiv \text{Prob}_f\{z_i \in \delta_0 \pm c\sigma_0\} . \quad (7.8)$$

Supposing $f_0(z) = N(\delta_0, \sigma_0^2)$, define

$$G_0(c) = 2\Phi(c) - 1 \quad \text{and} \quad G_1(c) = \int_{\delta_0 - c\sigma_0}^{\delta_0 + c\sigma_0} f_1(z) dz , \quad (7.9)$$

the probabilities that $z_i \in \delta_0 \pm c\sigma_0$ under f_0 and f_1 respectively. Then

$$p_0 = \frac{\pi(c) - G_1(c)}{G_0(c) - G_1(c)} \leq \frac{\pi(c)}{G_0(c)} , \quad (7.10)$$

the inequality following from the assumption that $G_1(c) \leq G_0(c)$, i.e. that the f_1 density is more dispersed than f_0 .

This leads to the estimated upper bound for p_0 ,

$$\hat{p}_0 = \frac{\hat{\pi}(c)}{G_0(c)} \quad \text{with} \quad \hat{\pi}(c) = \#\{z_i \in \delta_0 \pm c\sigma_0\}/N . \quad (7.11)$$

In particular if we assume $G_1(c) = 0$, in other words that the Interesting z_i 's always fall outside $\delta_0 \pm c\sigma_0$, then $\hat{p}_0 = \hat{\pi}(c)/G_0(c)$ is unbiased. (This is the same estimate suggested in Remark F of Efron et al. (2001).) Choosing $(\delta_0, \sigma_0) = (-0.35, 1.20)$ and $c = 1.5$ gave $\hat{p}_0 = 0.88$ for the Drug-mutation data, with bootstrap standard error 0.024.

F. Single-point Solutions for $(\delta_{\max}, \sigma_{\max})$ The distributions $g(\mu)$ providing $(\delta_{\max}, \sigma_{\max})$ in (5.6), as graphed in Figure 4, have their non-zero components supported at a single point μ_1 . For example, $g(\mu)$ for the entry giving $\sigma_{\max} = 1.04$ in Table 1 puts probability 0.90 at $\mu = 0$ and 0.10 at $\mu_1 = 1.47$. Single-point optimality was proved for three of the four cases in Figure 4, and verified by numerical maximization for the ‘‘General’’ case. Here is the proof for the σ_{\max} ‘‘Symmetric’’ case, the other two proofs being similar.

We consider symmetric distributions putting probability p_0 on $\mu = 0$ and probabilities p_j on symmetric pairs $(-\mu_j, \mu_j)$, $j = 1, 2, \dots, J$, so (5.3) becomes

$$f(z) = p_0\varphi(z) + \sum_{j=1}^J p_j[\varphi(z - \mu_j) + \varphi(z + \mu_j)]/2 . \quad (7.12)$$

Defining $c_0 = p_0/(1 - p_0)$, $r_j = p_j/p_0$, and $r_+ = \sum_1^J r_j = 1/c_0$, we can express σ_0 in (3.1) as

$$\sigma_0 = (1 - Q)^{-\frac{1}{2}} \quad \text{where} \quad Q = \frac{\sum_1^J r_j \mu_j^2 e^{-\mu_j^2/2}}{c_0 r_+ + \sum_1^J r_j e^{-\mu_j^2/2}} , \quad (7.13)$$

Here we have used $\delta_0 = 0$, which is true by symmetry assuming $p_0 \geq 1/2$. Then σ_{\max} in (5.6) can be found by maximizing Q .

We will show that with p_0 (and c_0) and $\mu_1, \mu_2, \dots, \mu_J$ held fixed in (7.12), Q is maximized by a choice of p_1, p_2, \dots, p_J having $J - 1$ zero values; this is a stronger version of the single-point result. Because Q is homogeneous in $\mathbf{r} = (r_1, r_2, \dots, r_J)$ in (7.13), we can consider the unconstrained maximization of $Q(\mathbf{r})$, subject only to $r_j \geq 0$ for $j = 1, 2, \dots, J$.

Differentiation gives

$$\partial Q / \partial r_j = \frac{1}{\text{den}} [\mu_j^2 e^{-\mu_j^2/2} - Q \cdot (c_0 + e^{-\mu_j^2/2})] , \quad (7.14)$$

“den” the denominator of Q . At a maximizing point \mathbf{r} we must have

$$\frac{\partial Q(\mathbf{r})}{\partial r_j} \leq 0 \quad \text{with equality if} \quad r_j > 0 . \quad (7.15)$$

Defining $R_j = \mu_j^2 / (1 + c_0 e^{\mu_j^2/2})$, (7.14)-(7.15) give

$$Q(\mathbf{r}) \geq R_j \quad \text{with equality if} \quad r_j > 0 . \quad (7.16)$$

Since $Q(\mathbf{r})$ is the maximum, this says that r_j , and p_j can only be non-zero if j maximizes R_j . In case of ties we can arbitrarily choose one of the maximizing j 's.

All of this shows that we need only consider $J = 1$ in (7.12). The global maximized value of r_0 in (7.12) is $\sigma_{\max} = (1 - R_{\max})^{-\frac{1}{2}}$ where

$$R_{\max} = \max_{\mu_1} \{ \mu_1^2 / (1 + c_0 e^{\mu_1^2/2}) \} . \quad (7.17)$$

The maximizing argument μ_1 ranges from 1.43 for $p_0 = .95$ to 1.51 for $p_0 = .70$. The corresponding result for δ_{\max} is simpler, $\mu_1 = \delta_{\max} + 1$.

G. Microarray Correlation in the Breast Cancer Study It is easy to spot an unwanted correlation structure among the eight BRCA2 microarrays. Let X be the 3226×8 matrix of BRCA2 data, with the columns of X standardized to have mean 0 and variance 1. A “de-gened” matrix \tilde{X} was formed by subtracting row-wise averages from each element of X ,

$$\tilde{x}_{ij} = x_{ij} - \sum_{k=1}^8 x_{ik} / 8 . \quad (7.18)$$

Table 3 shows the 8×8 correlation matrix of \tilde{X} . With genuine gene effects subtracted out, the correlations should vary around $-1/7 = -0.14$ if the columns of X are independent. Instead we see that the columns are correlated in blocks of four, with the off-diagonal block too negative and the on-diagonal blocks too positive.

Table 3:Correlation matrix for the BRCA2 data with row-wise means subtracted off, (7.17). It indicates positive correlations within the two blocks of four.

	1	2	3	4	5	6	7	8
1	1.00	0.02	0.02	0.23	-0.36	-0.35	-0.39	-0.34
2	0.02	1.00	0.10	-0.08	-0.30	-0.30	-0.23	-0.33
3	0.02	0.10	1.00	-0.17	-0.21	-0.26	-0.31	-0.27
4	0.23	-0.08	-0.17	1.00	-0.30	-0.23	-0.27	-0.32
5	-0.36	-0.30	-0.21	-0.30	1.00	-0.02	0.11	0.22
6	-0.35	-0.30	-0.26	-0.23	-0.02	1.00	0.15	0.13
7	-0.39	-0.23	-0.31	-0.27	0.11	0.15	1.00	0.07
8	-0.34	-0.33	-0.27	-0.32	0.22	0.13	0.07	1.00

8. Summary Large-scale simultaneous hypothesis testing, where the number of cases exceeds say 100, permits the empirical estimation of a null hypothesis distribution. The empirical null may be wider (more dispersed) than the theoretical null distribution that would ordinarily be used for a single hypothesis test. The choice between empirical and theoretical nulls can greatly influence which cases are identified as “Significant” or “Interesting”, as opposed to “Null” or “Uninteresting”, this being true no matter which simultaneous hypothesis testing method is used.

We present an analysis plan for large-scale testing situations:

- A density fitting technique is used to estimate the null hypothesis distribution f_0 , Figure 1 and Section 3.
- The local false discovery rate, an empirical Bayes version of standard FDR theory, provides inferences for the N cases, Figure 3 and Section 2.

There are many possible reasons for overdispersion of the empirical null distribution that would lead to the empirical null being preferred for simultaneous testing:

- Unobserved covariates in a observational study, Section 4.
- Hidden correlations, Section 6.
- A large proportion of genuine but uninterestingly small effects, Figure 5.

Large-scale testing differs in scientific intent from an individual hypothesis test. The latter is most often designed to reject the null hypothesis with high probability. Large-scale testing is usually more of a screening operation, intended to identify a *small* percentage of interesting cases, assumed to be on the order of 10% or less in this paper. Our estimation technique for the empirical null hypothesis is designed to be accurate under this constraint, Figure 4. More traditional estimation methods, involving permutations or quantiles, give incorrect f_0 estimates, Section 4 and Remark D.

Acknowledgment I am grateful to Robert Shafer, David Katzenstein, and Rami Kantor for bringing the Drug-mutation data to my attention, and to Robert Tibshirani for several helpful discussions.

References

- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J.R. Stat. Soc. Ser. B Stat. Methodol.* **57** 289-300.
- Dudoit, S., Shaffer J., and Boldrick J. (2003). “Multiple hypothesis testing in microarray experiments”. *Statistical Science* **18** 71-103.
- Efron, B. (2003). “Robbins, empirical Bayes, and microarrays”. *Annals Stat.* **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. and Tusher, V. (2001). Empirical Bayes analysis of a microarray experiment. *J. Amer. Statist. Assoc.* **96** 1151-1160.
- Efron, B. and Tibshirani, R. (1996). “Using specially designed exponential families for density estimation”. *Annals Stat.* **24** 2431-61.
- Efron, B. (1988). “Three examples of computer-intensive statistical inference”. *Sankhya* **50** 338-62.
- Hedenfalk, I., Duggen, D., Chen, Y., et al. (2001). “Gene expression profiles in hereditary breast cancer”. *New Engl. Jour. Medicine* **344** 539-48.
- Miller, R. (1981). *Simultaneous Statistical Inference*, Second Edition, Springer-Verlag, New York.
- Westfall, P. and Young, S. (1993). *Resampling-based multiple testing: examples and methods for p-value adjustments*. Wiley, New York.
- Wu, T., Schiffer, C., Shafer, R. et al. (2003). “Mutation patterns and structural correlates in Human Immunodeficiency Virus Type 1 Protease following different protease inhibitor treatments”. *Jour. Virology* **77(8)** 4836-47.