

Correlation and Large-Scale Simultaneous Significance Testing

Jinxi Liu

November 1, 2017

This article concerns the effect of correlation on multiple testing procedures, particularly false discovery rate techniques.

The experiments report two-sample t -statistics t_i comparing expression levels under two different conditions for N genes, $N = 3226$ for the breast cancer study, and $N = 7680$ for the HIV experiment; t_i tests the null hypothesis that gene i has the same expression distribution under both conditions.

The t_i 's have been converted to z -values for easy analysis.

$$Z_i = \Phi^{-1}(G_0(t_i)), \quad i = 1, 2, \dots, N, \quad (1)$$

where G_0 is a putative null cdf for the t -values. G_0 was taken to be a standard Student t cdf with appropriate degrees of freedom for the HIV study, whereas a permutation method provided G_0 for the breast cancer experiment (also nearly a Student t cdf).

Assuming that G_0 is the correct null distribution for t_i , transformation (1) yields for the null cases,

$$Z_i \sim N(0, 1) \tag{2}$$

called the theoretical null.

Microarray experiments involving genome-wide scans usually presuppose most of the genes to be null, the goal being to identify a small subset of interesting nonnull genes for future study, so we expect $N(0, 1)$ to fit the center of the z-value histogram.

This is not the case in Figure (1), where $N(0, 1)$ is too narrow for the breast cancer histogram and too wide for the HIV data.

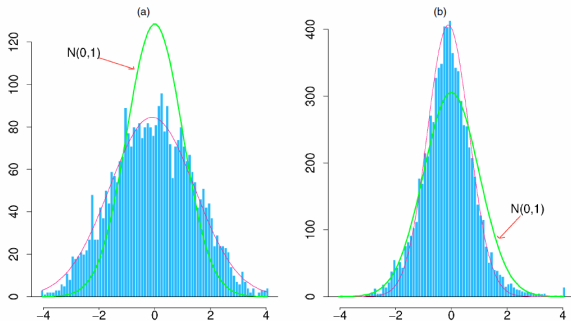


Figure 1: Histograms of z-Values From Two Microarray Experiments. (a) Breast cancer study, 3,226 genes. (b) HIV study, 7,680 genes. Heavy curves indicate $N(0, 1)$ theoretical null densities; Light curves indicate empirical null densities fit to central z-values, as done by Efron (2004). The theoretical null distributions are too narrow in (a) and too wide in (b). Both effects can be caused by correlations among the null z-values.

Correlation can cause effects like those shown in Figure 1, considerably widening or narrowing the distribution of the null z -values.

These effects have a substantial impact on simultaneous significance testing and must be accounted for in deciding which cases should be reported as nonnull.

Broadly speaking, a wide central histogram like that for the breast cancer data implies more null z -values in the tails, so that significance levels judged according to the theoretical $N(0, 1)$ null are too liberal. Conversely, the theoretical null is too conservative for the HIV data.

A surprisingly simple result emerges in which the main effect of all the pairwise correlations (several million of them for Fig. 1's examples) is summarized by a single dispersion variate, A :

a positive value of A widens the central peak of the z -value histogram, even assuming that the theoretical null (2) is individually correct for all the cases, whereas negative A narrows it, as in Figure 1(b).

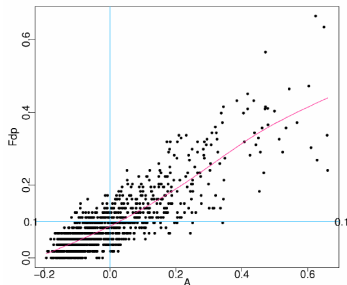


Figure 2: Benjamini-Hochberg FDR Controlling Procedure, $q = .10$, Run for 1,000 Simulation Trials; True False Discovery Proportion, FDP, Plotted versus Dispersion Variate A . Overall FDP averaged .096, close to q , but with a strong dependence on A , as shown by smooth regression curve.

We see a strong dependence on the dispersion variate A ; FDP averaged .34 for the upper 5% of A realizations but only .03 for the lower 5%.

A can be estimated from the width of the central peak of the z -value histogram; narrow peaks imply negative A 's and wide peaks imply positive A 's.

In terms of Figure 2, this enables the statistician to condition the FDR estimate on its approximate location along the A axis.

Correlation effects on the null distribution

For following calculations, it is assumed that all cases are null,

$$Z_i \sim N(0, 1), \quad \text{for } i = 1, 2, \dots, N \quad (3)$$

It is helpful to work with histogram counts rather than with the vector of z-values itself. Each histogram in Figure 1 has its z-axis partitioned into $K = 82$ bins of width $\Delta = .1$, running from -4.1 to 4.1 .

Correlation effects on the null distribution

We denote the count vector by y ,

$$y_k = \#\{z_i \text{ in } k\text{th bin}\}, \quad k = 1, 2, \dots, K \quad (4)$$

The histogram counts y_k arise from a partition of Z into K bins,

$$Z = \bigcup_{k=1}^K z_k, \quad (5)$$

each bin being of width Δ , with center point " $z[k]$ "

Correlation effects on the null distribution

The following definitions lead to useful representations for the mean and covariance of y .

$$\pi_k(i) = Pr\{z_i \in Z_k\}, \quad \pi_{k\cdot} = \frac{\sum_{i=1}^N \pi_k(i)}{N} \quad (6)$$

$$\gamma_{kl}(i, j) = Pr\{z_i \in Z_k \text{ and } z_j \in Z_l\}, \quad \gamma_{kl\cdot} = \frac{\sum_{i \neq j} \gamma_{kl}(i, j)}{N(N-1)}. \quad (7)$$

Correlation effects on the null distribution

Because of assumption (3), all of the $\pi_k(i)$ are determined by $\varphi(z)$, the standard normal density, with Taylor approximation around centerpoint $z[k]$,

$$\pi_{k\cdot} = \pi_k(i) \doteq \Delta \cdot \varphi(z[k]) \quad (8)$$

The expectation vector, $\mathbf{v} = (v_1, v_2, \dots, v_K)'$ of y is determined by (3),

$$\mathbf{v} = N\boldsymbol{\pi}_{\cdot} \doteq (\dots, N\Delta\varphi(z[k]), \dots)'. \quad (9)$$

Lemma 1

$$\text{cov}(\mathbf{y}) = C_0 + C_1 \quad (10)$$

$$C_0 = \text{diag}(\mathbf{v}) - \mathbf{v}\mathbf{v}'/N = N[\text{diag}(\boldsymbol{\pi}.) - \boldsymbol{\pi}.\boldsymbol{\pi}.'] \quad (11)$$

$$C_1 = (1 - \frac{1}{N})\text{diag}(\mathbf{v})\boldsymbol{\delta}\text{diag}(\mathbf{v}), \quad \text{with } \delta_{kl} = \gamma_{kl.}/\pi_{k.}\pi_{l.} - 1 \quad (12)$$

Correlation effects on the null distribution

If z_i and z_j are independent, then $\gamma_{kl}(i, j) = \pi_k(i)\pi_j(l)$.

Independence of all z -values implies that all of the elements of matrix δ in (12) are 0, leaving $\text{cov}(y) = C_0$.

Conversely, the amount of correlation between the z -values determines the size of δ and the increase of $\text{cov}(y)$ above C_0 .

Correlation effects on the null distribution

To approximate δ , we add the assumption of bivariate normality for any pair of z -values, $\text{cov}(z_i, z_j) \equiv \rho_{ij}$.

Let $g(\rho)$ indicate the empirical density of the $N(N - 1)$ correlations ρ_{ij}

Lemma 2 Under the bivariate normal approximation, the matrix δ has elements

$$\delta_{kl} \doteq \int_{-1}^1 \int_{-1}^1 \left[\frac{1}{\sqrt{1 - \rho^2}} \exp\left(\frac{\rho}{2(1 - \rho^2)} \times \{2z[k]z[l] - \rho(z[k]^2 + z[l]^2)\} \right) - 1 \right] g(\rho) d\rho \quad (13)$$

Correlation effects on the null distribution

Application of Lemmas 1 and 2 requires estimation of the correlation density $g(\rho)$, which we can obtain from observed correlations between the rows of X .

$$\text{breastcancer} : g(\rho) \sim N(0, .153^2) \quad (14)$$

$$\text{HIV} : g(\rho) \sim N(0, .42^2) \quad (15)$$

Approximation (14) indicates a substantial amount of global correlation among genes in the breast cancer study, and even more correlation for the HIV study.

Lemmas 1 and 2 decompose the covariance matrix of the count vector y into an independence term C_0 and an additional term C_1 that accounts for correlation among the z -values.

This section presents a simple approximation to C_1 in terms of its first eigenvector,

Lemma 3 Suppose that $g(\rho)$, the correlation density in (13) has mean 0 and standard deviation

$$\alpha = \left[\int_{-1}^1 \rho^2 g(\rho) d\rho \right]^{-1/2} \quad (16)$$

Then the matrix δ in (12) is approximated by the outer product

$$\delta \doteq \alpha^2 \mathbf{q} \mathbf{q}', \text{ where } q_k = \frac{z[k]^2 - 1}{\sqrt{2}} \quad (17)$$

First Eigenvector

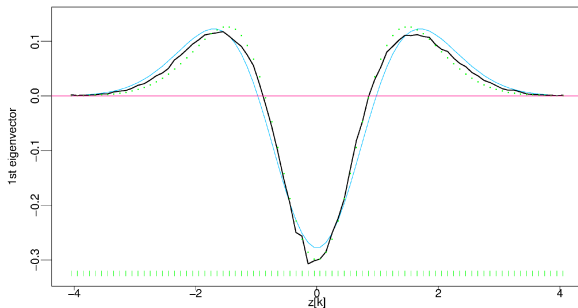


Figure 3: First Eigenvectors of Three Different Estimates of $\text{cov}(y)$: Normal-Theory Estimate for Breast Cancer Data (smooth curve); Normal-Theory Estimate for HIV Data (dots), and Permutation Estimate for Breast Cancer Data (jagged curve). The striking 'wing-shaped' form is proportional to the second derivative of the standard normal density.

First Eigenvector

Combining the three lemmas yields a useful approximation for the null covariance matrix of the count vector \mathbf{y} under the bivariate normal assumptions.

Theorem If $g(\rho)$ has mean 0 and standard deviation α , then

$$\text{cov}(\mathbf{y}) \doteq [\text{diag}(\mathbf{v}) - \mathbf{v}\mathbf{v}'/N] + \left(1 - \frac{1}{N}\right) (\alpha \mathbf{W})(\alpha \mathbf{W})', \quad (18)$$

Here $\mathbf{v} = E\{\mathbf{y}\}$, \mathbf{W} has components

$$W_k = N\Delta\varphi(z[k])\frac{z[k]^2 - 1}{\sqrt{2}} = N\Delta w(z[k]) \quad (19)$$

The Theorem summarizes the effect of z 's entire correlation structure in a single parameter α .

Poisson process considerations lead to a somewhat rough but evocative interpretation of the Theorem. Let $\mathbf{y} \sim Po(\mathbf{u})$ indicate a vector of independent Poisson variates, $y_k \overset{ind}{\sim} Po(u_k)$ for $k = 1, 2, \dots, K$, whereas $\mathbf{y} \sim (\mathbf{v}, \Gamma)$ denotes that vector \mathbf{v} has mean \mathbf{v} and covariance Γ .

Consider the number of cases N to be a Poisson variate, say

$$N \sim Po(N_0), \quad (20)$$

with $N_0 = 3226$ in the breast cancer study. This simplifies (18) slightly, to

$$\text{cov}(\mathbf{y}) \doteq \text{diag}(\mathbf{v}) + (\alpha \mathbf{W})(\alpha \mathbf{W})'. \quad (21)$$

If the z -values are independent, then (20) makes the counts, y_k , independent Poisson variates,

$$\mathbf{y} \sim Po(\mathbf{v}), \text{ agreeing with (21) at } \alpha = 0. \quad (22)$$

A hierarchical model generalizes (22) to incorporate dependence. We assume that \mathbf{y} depends on a mean vector \mathbf{u} , itself random,

$$\mathbf{y}|\mathbf{u} \sim Po(\mathbf{u}), \text{ where } \mathbf{u} \sim (\mathbf{v}, \Gamma), \quad (23)$$

so that the components of \mathbf{y} are conditionally independent given \mathbf{u} but marginally dependent, with mean and covariance

$$\mathbf{y} \sim (\mathbf{v}, \text{diag}(\mathbf{v}) + \Gamma). \quad (24)$$

To match (21), we set

$$\Gamma = (\alpha \mathbf{W})(\alpha \mathbf{W})'. \quad (25)$$

Formulas (24) and (25) suggest a hierarchical Poisson structure for the count vector \mathbf{y} ,

$$\mathbf{y} \sim \text{Po}(\mathbf{u}), \text{ where } \mathbf{u} = \mathbf{v} + A\mathbf{W}, \text{ with } A \sim (0, \alpha^2) \quad (26)$$

If $\alpha = 0$, then this reduces to the independence case (22); otherwise, the Poisson intensity vector \mathbf{v} is modified by the addition of an independent random multiple A of \mathbf{W} with standard deviation α .

Suppose for a moment that we knew which z_i 's among the full set of z -values correspond to null cases. For a given choice of x , define

$$Y(x) = \#\{\text{null } z_i \geq x\} \quad T(x) = \#\{z_i \geq x\} \quad (27)$$

$$FDR(x) = E\{Y(x)\}/T(x). \quad (28)$$

$FDR(x)$ is an empirical Bayes estimate of the a posteriori probability that case i is null given $z_i \geq x$, amounting to a version of Storey's "q-value".

The hierarchical structure gives conditional expectation

$$E\{\mathbf{y}|A\} = \mathbf{v} + A\mathbf{W} \quad (29)$$

Letting the bin width $\Delta \rightarrow 0$ produces a conditional version of (29)

$$E\{Y(x)|A\} = N\bar{\Phi}(x) \left[1 + A \frac{x\varphi(x)}{\sqrt{2}\bar{\Phi}(x)} \right] \quad (30)$$

Conditional FDR

$$FDR(Y(x)|A) = FDR_0(x) \left[1 + A \frac{x\varphi(x)}{\sqrt{2}\bar{\Phi}(x)} \right], \quad (31)$$

where $FDR_0(x)$ is the standard unconditional estimate $N\bar{\Phi}(x)/T(x)$ based on the theoretical $N(0, 1)$ null.

Large-scale Significance Testing

A can be estimated from the central spread of the histogram of z-values and then used to condition inferences, as in figure 2.

For $x_0 > 0$, let

$$Y_0 = \#\{z_i \in [-x_0, x_0]\}, \quad (32)$$

and define

$$P_0 = 2\Phi(x_0) - 1 \quad \text{and} \quad Q_0 = \sqrt{2}x_0\varphi(x_0) \quad (33)$$

Then (30) give $E\{Y_0|A\} \doteq N[P_0 - AQ_0]$, so

$$\hat{A} = \frac{P_0 - \hat{P}_0}{Q_0}, \quad \text{where } \hat{P}_0 = Y_0/N \quad (34)$$

We could substitute \hat{A} from (34) into (31) to obtain a conditional FDR estimate,

$$FDR(x|\hat{A}) = FDR_0(x) \left[1 + \hat{A} \frac{x\varphi(x)}{\sqrt{2}\bar{\Phi}(x)} \right], \quad (35)$$

In situations like that of Figure 2, the goal would be to accurately assess our position on the A axis, to better estimate FDP rather than estimating the unconditional average of FDP .

Large-scale Significance Testing

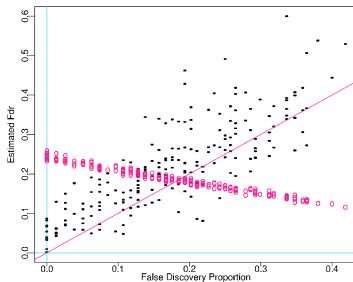


Figure 4: Simulation Experiment Comparing Conditional FDR Estimates (solid points), With Unconditional Estimates (open circles); $N = 3,000$, $p_0 = .95$. Null counts are generated as in (26), $\alpha = .15$; nonnull counts from $z \sim N(2.5, 1.25)$. The horizontal axis is the actual FDP, $FDP(x)$, $x = 2.5$, for each of the 200 trials. The unconditional estimate based on the theoretical null distribution declines as actual FDP increases.

Strikingly, the unconditional estimate goes in the wrong direction, declining as the actual FDP increases. This yields misleading inferences at both ends of the FDP scale. The conditional FDR estimate correctly tracks FDP , although with a considerable amount of noise.