

11-2006

# Quick Calculation for Sample Size while Controlling False Discovery Rate with Application to Microarray Analysis

Peng Liu

*Iowa State University*, [pliu@iastate.edu](mailto:pliu@iastate.edu)

J.T. Gene Hwang

*Cornell University*

Follow this and additional works at: [http://lib.dr.iastate.edu/stat\\_las\\_preprints](http://lib.dr.iastate.edu/stat_las_preprints)



Part of the [Statistics and Probability Commons](#)

---

## Recommended Citation

Liu, Peng and Hwang, J.T. Gene, "Quick Calculation for Sample Size while Controlling False Discovery Rate with Application to Microarray Analysis" (2006). *Statistics Preprints*. 55.

[http://lib.dr.iastate.edu/stat\\_las\\_preprints/55](http://lib.dr.iastate.edu/stat_las_preprints/55)

This Article is brought to you for free and open access by the Statistics at Iowa State University Digital Repository. It has been accepted for inclusion in Statistics Preprints by an authorized administrator of Iowa State University Digital Repository. For more information, please contact [digirep@iastate.edu](mailto:digirep@iastate.edu).

---

# Quick Calculation for Sample Size while Controlling False Discovery Rate with Application to Microarray Analysis

## Abstract

Sample size estimation is important in microarray or proteomic experiments since biologists can typically afford only a few repetitions. In the multiple testing problems involving these experiments, it is more powerful and more reasonable to control false discovery rate (FDR) or positive FDR (pFDR) instead of type I error, e.g., family-wise error rate (FWER) (Storey and Tibshirani, 2003). However, the traditional approach of estimating sample size is no longer applicable to controlling FDR, which has left most practitioners to rely on haphazard guessing. We propose a procedure to calculate sample size while controlling false discovery rate. Two major definitions of the false discovery rate (FDR in Benjamini and Hochberg, 1995, and pFDR in Storey, 2002) vary slightly. Our procedure applies to both definitions. The proposed method is straightforward to apply and requires minimal computation, as illustrated with two sample t-tests and F-tests. We have also demonstrated by simulation that, with the calculated sample size, the desired level of power is achievable by the q-value procedure (Storey, Taylor and Siegmund, 2004) when gene expressions are either independent or dependent.

## Disciplines

Statistics and Probability

## Comments

This preprint was published as Peng Liu and J. T. Gene Hwang, "Quick calculation for sample size while controlling false discovery rate with application to microarray analysis", *Bioinformatics* (2007): 739-746, doi: [10.1093/bioinformatics/btl664](https://doi.org/10.1093/bioinformatics/btl664)

# Quick Calculation for Sample Size while Controlling False Discovery Rate with Application to Microarray Analysis

Peng Liu<sup>1</sup> and J. T. Gene Hwang<sup>2</sup>

<sup>1</sup>Department of Statistics,

Iowa State University, Ames, IA 50011

<sup>2</sup>Department of Mathematics and Department of Statistical

Science,

Cornell University, Ithaca, NY 14853

November, 2006

## SUMMARY.

Sample size estimation is important in microarray or proteomic experiments since biologists can typically afford only a few repetitions. Classical procedures to calculate sample size are based on controlling type I error, e.g., family-wise error rate (FWER). In the context of microarray and other large-scale genomic data, it is more powerful and more reasonable to control false discovery rate (FDR) or positive FDR (pFDR)(Storey and Tibshirani, 2003). However, the traditional approach of estimating sample size is no longer applicable to controlling FDR, which has left most practitioners to rely on haphazard guessing.

We propose a procedure to calculate sample size while controlling false discovery rate. Two major definitions of the false discovery rate (FDR in Benjamini and Hochberg, 1995, and pFDR in Storey, 2002) vary slightly. Our procedure applies to both definitions. The proposed method is straightforward to apply and requires minimal computation, as illustrated with two sample  $t$ -tests and  $F$ -tests. We have also demonstrated by simulation that, with the calculated sample size, the desired level of power is achievable by the q-value procedure (Storey, Taylor and Siegmund, 2004) when gene expressions are either independent or dependent.

**KEY WORDS:** Power, Sample Size Calculation, pFDR, Experimental Design, Genomics

## 1. Introduction

Microarray and proteomic experiments are becoming popular and important in many biological disciplines, such as neuroscience (Mandel et al., 2003), pharmacogenomics, genetic disease and cancer diagnosis (Heller, 2002). These

**Table 1**  
*Outcomes when testing  $m$  hypothesis.*

Hypothesis	Accept	Reject	Total
Null true	$U$	$V$	$m_0$
Alternative true	$T$	$S$	$m_1$
Total	$W$	$R$	$m$

experiments are rather costly in terms of both materials (samples, reagents, equipments, etc.) and laboratory manpower. Many microarray experiments employ only a small number of replicates (2 to 8) (Yang and Speed, 2003). In many cases, the sample size is not adequate to achieve reliable statistical inference, resulting in waste of resources. Therefore, scientists often ask the following question. How big should the sample size be?

To answer this question, we will calculate sample size that controls some error rate and achieves a desired power. When calculating sample size for a single test, the error rate to control is traditionally the type I error rate, i.e., the probability of concluding a false positive by rejecting the true null hypothesis. However, we are simultaneously testing a huge number of hypotheses, each relating to a gene. Hence, multiple testing is commonly applied in the analysis of microarray data. There are several kinds of error rates to control in this context, such as family-wise error rate (FWER) or false discovery rate (FDR). Assume there are  $m$  genes on microarray chips and each gene is tested for the significance of differential expression. The test outcomes are summarized in Table 1, where, for example,  $V$  is the number of false positives and  $R$  is the number of rejections among the  $m$  tests (Benjamini and Hochberg, 1995).

The FWER is defined to be the probability of making at least one false positive error:  $FWER = Pr(V \geq 1)$ . Rejecting each individual test with a type I error rate of  $\alpha/m$  guarantees, by Bonferroni's type of argument, that FWER is controlled at level  $\alpha$  in the strong sense, *i.e.*,  $FWER \leq \alpha$  for any combinations of null and alternative hypotheses. Benjamini and Hochberg (1995) proposed another type of error to control – FDR, which is defined to be the expected proportion of false positives among the rejected hypotheses:

$$FDR = E \left[ \frac{V}{R} | R > 0 \right] Pr(R > 0) . \quad (1)$$

Storey (2002) proposed to control positive FDR (pFDR), *i.e.*,

$$pFDR = E \left[ \frac{V}{R} | R > 0 \right] = \frac{FDR}{Pr(R > 0)} . \quad (2)$$

In many cases of genomic data such as microarray, it was argued in Storey and Tibshirani (2003) that it is more reasonable and more powerful to control FDR or pFDR instead of FWER. However, the sample size has been traditionally calculated with a certain type I error rate and can not be directly applied with FDR control.

Recently, a few papers investigated the needs to calculate sample size while controlling FDR and proposed ways to pursue this goal. Yang et al. (2003) applied several inequalities to get a type I error rate that corresponds to the controlled level of FDR. Due to the inequalities applied, the sample size is likely overestimated. Pawitan et al. (2005) investigated several operating characteristic curves to visualize the relationship between FDR, sensitivity and sample size. Although their approach can be useful in calculating the sample size, no simple direct algorithm was provided. Jung (2005)

derived a formula which relates FDR and the type I error rate. Then FDR is controlled by an appropriate level of type I error rate. Pounds and Cheng (2005) proposed an algorithm to iteratively search for the sample size at which the desired power and controlled level of FDR can be achieved. Since FDR controlling procedure is gaining popularity in multiple testings for many problems including microarray analysis, it is important to be able to calculate sample size needed to control the FDR when designing the experiment.

Here we propose a procedure to calculate the sample size for multiple testing while controlling FDR. First, for any estimate of the proportion of non-differentially expressed genes and the level of FDR to control, we find a rejection region for each sample size. Then power is calculated for the selected rejection region for each sample size. According to the desired power, a sample size is finally decided.

Jung's approach (2005), which was known to us after we had finished our first draft, is more related to our proposed approach than others. Both Jung's and our approaches base on the same model assumptions which lead to the same FDR expression. The FDR expression is then controlled by studying its relationship to a quantity, which is the type I error rate for Jung and the critical value (the rejection region) for us. Our approach, however, is more graphical. This allows the visualization of the tradeoff between power and sample size and provides quick answer when the user-defined quantities such as power are modified.

In spite of the similarity, this paper extends the approach further to several different directions and we find our approach very satisfactory. First, we apply our approach to  $F$ -tests which are widely used in microarray data

analysis (Cui et al., 2005). Second, we study our approach carefully for the case when the means and variances for expression levels vary among genes, an important and practical setting for microarray. Third, we also show by simulation, that the q-value procedure for controlling FDR proposed by Storey, Taylor and Siegmund (2004) using our suggested sample size achieves the target power to a satisfactory degree. This answers the question positively as to whether there would be any statistical procedure that can realize the target power claimed by the proposed method. Finally, we also compare our approach with Yang et al. (2003) and Pounds and Cheng (2005) which provide more well-defined algorithms than other papers. Our simulation demonstrates that the proposed method is superior.

The paper is organized as follows. Section 2 describes our proposed method illustrated with two-sample  $t$ -tests and  $F$ -tests. In Section 3, we report the result of simulation studies that compare the power based on proposed method to the actual result from q-value procedure. Section 4 discusses our result.

Codes for the proposed method both in R and in Matlab are available to implement the method. The R code can be applied in conjunction with the `ssize` package from bioconductor.

## 2. Method

In this section, we first illustrate our idea and then show how to apply the proposed method for two designs of microarray experiment.

### 2.1 Proposed Method

The proposed method is derived from the definition of pFDR. Let  $H = 0$  if null hypothesis is true and  $H = 1$  if alternative hypothesis is true.

In a microarray experiment,  $H = 1$  represents differential expression for a gene whereas  $H = 0$  represents no differential expression. We assume as in Theorem 1 of Storey (2002) that all tests are identical, independent and Bernoulli distributed with  $Pr(H = 0) = \pi_0$ , where  $\pi_0$  is interpreted as the proportion of non-differentially expressed genes. By Storey's theorem,

$$pFDR(\Gamma) = Pr(H = 0|T \in \Gamma) , \quad (3)$$

where  $T$  denotes the test statistic and  $\Gamma$  denotes the rejection region. Because the number of genes is large, typically ranging from 5,000 to 30,000, the probability of no significant findings is close to zero (Storey and Tibshirani, 2003). Therefore our result also applies to controlling FDR because  $FDR = pFDR \cdot Pr(R > 0)$ . Suppose the level of FDR is chosen to be  $\alpha$ , the following relationship is derived via simple algebra (see Appendix A).

$$\frac{\alpha}{1 - \alpha} \frac{1 - \pi_0}{\pi_0} = \frac{Pr(T \in \Gamma|H = 0)}{Pr(T \in \Gamma|H = 1)} . \quad (4)$$

For simplicity in notation, we will denote

$$\Lambda = \frac{\alpha}{1 - \alpha} \frac{1 - \pi_0}{\pi_0} . \quad (5)$$

In order to achieve a FDR level to be  $\alpha$  (or less), we choose the rejection region  $\Gamma$  so that the right hand side of Equation (4) is equal to (or less than)  $\Lambda$  (see Appendix A).

## 2.2 Applications of Proposed Method

Microarray experiments are usually set up to find differentially expressed genes between different treatments. The data of scanned intensity for microarray usually go through quality control, transformation and normalization, as reviewed in Smyth et al. (2003) and Quackenbush (2002). We assume

that data first go through those steps before statistical tests are applied. Before the experiment, we have no observations to check the distribution. It seems reasonable to make a convenient assumption that the distribution of the pre-processed data is normal and hence two-sample  $t$ -tests and  $F$ -tests are applicable. The same assumption is also made by other proposed methods to calculate sample size (Dobbin and Simon, 2005; Hu et al., 2005; Hua et al., 2005; Jung, 2005).

*2.2.1 Two-Sample Comparison with  $t$ -test* Suppose we want to find differentially expressed genes between a treatment and a control group using two-sample  $t$ -tests. The tested hypothesis for each gene is  $H_0 : \mu_{T,g} = \mu_{C,g}$  versus  $H_1 : \mu_{T,g} \neq \mu_{C,g}$ , where  $\mu_{T,g}$  and  $\mu_{C,g}$  are mean expressions of  $g$ -th gene for treatment and control group respectively. Let  $x_{gj}$  and  $y_{gj}$  denote the observed gene expression levels for treatment and control group respectively for the  $g$ -th gene and  $j$ -th replicate. Assuming equal variance between treatment and control group, the test statistic for the  $g$ -th gene is:

$$T_g = \frac{\bar{x}_g - \bar{y}_g}{\sqrt{S_g^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} , \quad (6)$$

where  $S_g^2 = \frac{1}{n_1+n_2-2} [\sum_{j=1}^{n_1} (x_{gj} - \bar{x}_g)^2 + \sum_{j=1}^{n_2} (y_{gj} - \bar{y}_g)^2]$  is the pooled sample variance,  $\bar{x}_g$  and  $\bar{y}_g$  are the means of observed expression levels for gene  $g$  for the two groups respectively. The test statistic  $T_g$  has a central  $t$ -distribution under the null hypothesis and noncentral  $t$ -distribution under the alternative hypothesis. We reject the null hypothesis if  $|T_g| > c_g$ , for which  $c_g$  is to be

determined. Applying Equation (4), we find critical value  $c$  that satisfies:

$$\begin{aligned}\Lambda &= \frac{Pr(|T_g| > c_g | H = 0)}{Pr(|T_g| > c_g | H = 1)} \\ &= \frac{2 \cdot T_{n_1+n_2-2}(-c_g)}{1 - T_{n_1+n_2-2}(c_g|\theta_g) + T_{n_1+n_2-2}(-c_g|\theta_g)},\end{aligned}\quad (7)$$

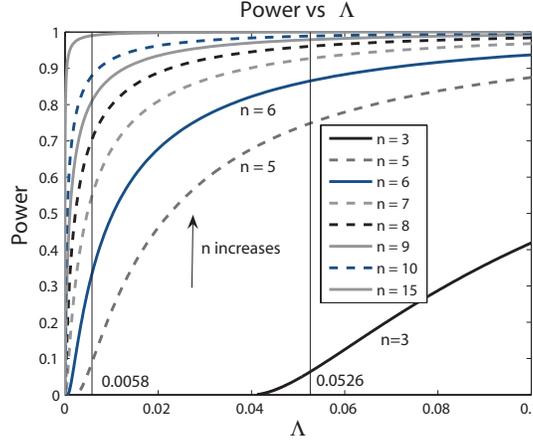
where  $T_d(\bullet|\theta)$  is the cumulative distribution function of a non-central  $t$ -distribution with  $d$  degrees of freedom and non-centrality parameter  $\theta$ . Moreover,  $T_d(\bullet)$  is  $T_d(\bullet|\theta)$  for  $\theta = 0$ . In (7),

$$\theta_g = \frac{\Delta_g}{\sigma_g \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}\quad (8)$$

where  $\Delta_g = \mu_{T,g} - \mu_{C,g}$  is the true difference between the mean expressions of treatment and control groups and  $\sigma_g$  is the standard deviation for gene  $g$ . In this section, we assume a simplified case that  $\Delta_g$  and  $\sigma_g$  are identical for all genes. Section 2.2.3 deals with the more realistic case when  $\Delta_g$  and  $\sigma_g$  vary among genes. So the subscript  $g$  is dropped in this section.

After we find critical values, power is calculated and sample size will be determined. A special and common case is the balanced design when the two groups have the same sample size, i.e.,  $n = n_1 = n_2$ .

Figure 1 plots the power against  $\Lambda$  for this case. For any selected values of  $\alpha$  (FDR level) and  $\pi_0$  (proportion of non-differentially expressed genes), a sample size can be determined based on this plot for a desired power. As an example, we want to determine the sample size when  $\pi_0 = 90\%$ . Suppose a two-fold change is desired (correspondingly,  $\Delta = \log_2(2) = 1$ ) and  $\sigma = 0.5$  from previous knowledge, then  $\frac{\Delta}{\sigma} = 2$ . Controlling FDR at 5% results in  $\Lambda = \frac{5\%}{1-5\%} \frac{100\%-90\%}{90\%} = 0.0058$ . The vertical line at  $\Lambda = 0.0058$  intersects



**Figure 1.** Summary of relationship between power and  $\Lambda = \frac{\alpha - \pi_0}{1 - \alpha - \pi_0}$  for  $\frac{\Delta}{\sigma} = 2$ . If  $\pi_0$  is estimated to be 90%, controlling FDR at 5% results in  $\Lambda = 0.0058$ . Along the vertical indicator line, we get power for each sample size. Another indicator line shows the position when  $\Lambda = 0.0526$  which is a result of FDR = 5% and  $\pi_0 = 50\%$ .

the power curves for different sample sizes. A desired power of 80% would determine a sample size of 9. Then 9 samples from each group are needed.

Figure 1 shows a flexible way to apply our method because we can get sample size for any  $\pi_0$  and controlled level of  $\alpha$ . A more straightforward way to view the result is presented in Figure 2, where we plot power versus sample size when FDR is controlled at 5% and various curves correspond to various  $\pi_0$ 's. For the same example, when  $\pi_0 = 90\%$  and  $\frac{\Delta}{\sigma} = 2$ , we determine  $n = 9$  to get at least 80% of power using the middle curve in Figure 2(a).

We shall take  $\sigma$  to be 0.2, which is approximately the 90th percentile of residual standard deviations for the granulosa cell tumor microarray data in Cui et al. (2005). Here 90th percentile is a conservative choice in that if we had used a percentage smaller than 90%, the sample size would be smaller. If still a 2-fold change ( $\Delta = \log_2(2) = 1$ ) is considered to be true effect

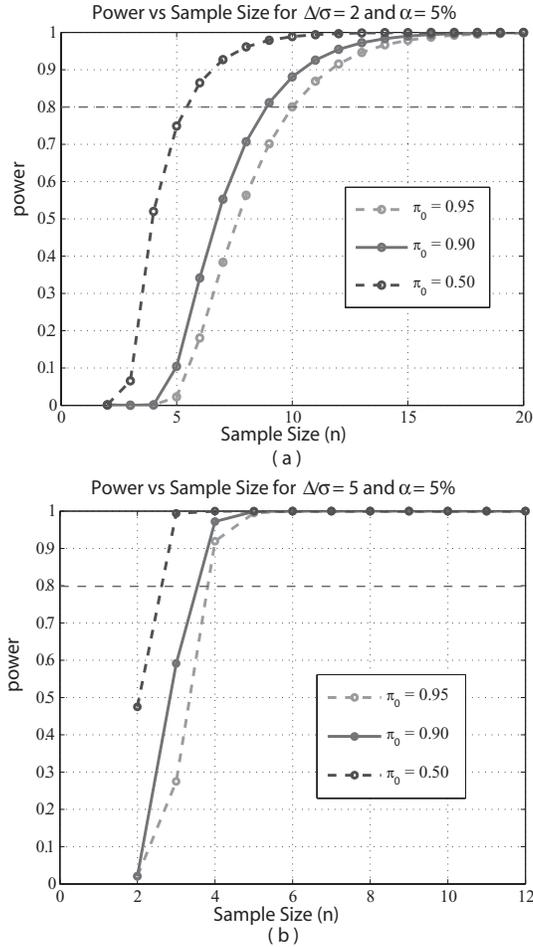
size, then  $\frac{\Delta}{\sigma} = 5$ . From the middle curve of Figure 2(b), corresponding to  $\pi_0 = 0.9$ , one can determine that a sample size of 4 is needed to obtain at least 80% of power.

*2.2.2 Multi-Sample Comparison with F-test* For microarray experiments comparing several treatments, there are different design schemes applied (Yang and Speed, 2003). Suppose without any replication, a design requires  $s$  slides. We call the  $s$  slides a *set* for this design. For example, we want to compare gene expressions among three treatments, such as livers from three genotypes of mice (Horton et al., 2003). If we apply a loop design, as shown in Figure 3, a “set” of three slides is needed for cDNA microarray experiment. Whether the replicates are different biological samples or different technical repetitions, our method is applicable as long as the appropriate parameter (means and variances) are used in the calculation. We recommend to use different biological samples in the experiment because this would provide more general conclusions. The question is how many sets of the slides is adequate to obtain a sufficient power and a controlled FDR.

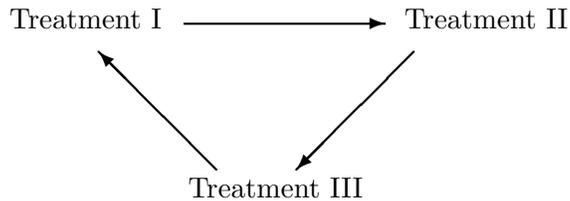
For each individual gene, the experimental design can be formulated with the same linear model for each set  $i$ ,  $i = 1, 2, \dots, n$ ,

$$Y_{g,i} = X\beta_g + \varepsilon_{g,i} , \quad (9)$$

where  $\beta_g(p \times 1)$  is the vector of parameters for gene  $g$ ,  $Y_{g,i}$  is the observed vector for  $g$ -th gene in the  $i$ -th set,  $X$  is the design matrix and  $\varepsilon_{g,i}$  is the error term. It is assumed that the errors are independent across genes and across sets in this section. For the design in Figure 3,  $Y_g$  would be the log-ratio of normalized gene expression levels for  $g$ -th gene, and two estimable



**Figure 2.** Plot of power versus sample size for  $t$ -test. Controlling FDR at 5%, we applied the proposed method to calculate power for each sample size. Panel (a) is for  $\frac{\Delta}{\sigma} = 2$  and panel (b) is for  $\frac{\Delta}{\sigma} = 5$ .



**Figure 3.** A design example for microarray experiment to compare gene expressions among three treatments. By convention, each arrow represents one two-color array with the green-labeled sample at the tail and the red-labeled sample at the head of the arrow. This design needs 3 arrays.

parameters can be the gene expression difference between treatment I and II, and difference between treatment I and III (Yang and Speed, 2003). Then the design matrix is

$$X = \begin{bmatrix} 1 & 0 \\ -1 & 1 \\ 0 & -1 \end{bmatrix}.$$

More complicated models can be constructed for more complex designs and corresponding terms should be added for effects that are not corrected during normalization, such as such array effects, dye effects and block effects. See, for example, Cui et al. (2005). For  $n$  sets of slides for a design, the least square estimate of  $\beta_g$  is:

$$\hat{\beta}_g = \sum_{i=1}^n (X'X)^{-1} X' Y_{g,i} / n = (X'X)^{-1} X' \sum_{i=1}^n Y_{g,i} / n. \quad (10)$$

With the assumption of normal distribution for the error,  $\hat{\beta}_g$  is also normally distributed,

$$\hat{\beta}_g \sim N(\beta_g, \sigma_g^2 (X'X)^{-1} / n).$$

We can apply this result and draw statistical inference for these parameters and their linear contrasts.

In general, assume that the question of interest is to test  $H_0 : L'\beta_g = 0$  versus  $H_1 : L'\beta_g \neq 0$ , where  $L$  is a  $p \times k$  coefficient matrix ( $k \leq p$ ) or  $p \times 1$  vector for the linear contrast(s) of interest. For simplicity, we omit the subscript  $g$  since we assume that the same test is applied for all genes separately. The  $F$ -tests based on  $n$  sets can be constructed with the following test statistic:

$$F_n = \frac{(L'\hat{\beta})' \cdot [L'(X'X)^{-1}L/n]^{-1} \cdot (L'\hat{\beta})/k}{\sum_{i=1}^n (Y_i - X\hat{\beta})'(Y_i - X\hat{\beta}) / (d(n))}. \quad (11)$$

Under  $H_0$ ,  $F_n$  follows a  $F$ -distribution with  $k$  and  $d(n)$  degrees of freedom where  $d(n)$  is a function of  $n$  and depends on the design. For example,  $d(n)$  for the design shown in Figure 3 is  $3n - 2$ . Under  $H_1$ ,  $F_n$  follows a non-central  $F$ -distribution with the same degrees of freedom and a non-centrality parameter  $\lambda$ :

$$\lambda = (L'\beta)'\Sigma^{-1}(L'\beta) , \quad (12)$$

where  $\Sigma = \sigma^2 L'(X'X)^{-1}L/n$ .

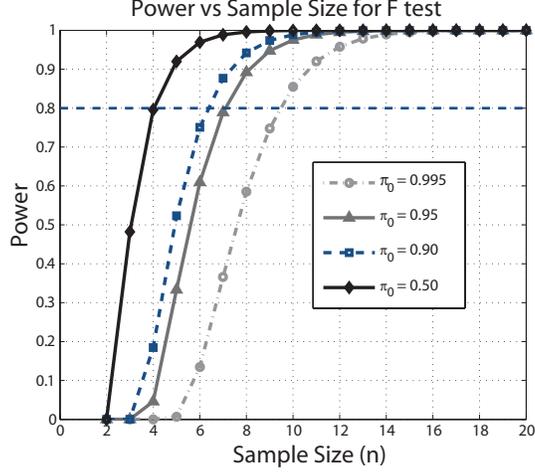
Applying Equation (4), we get

$$\Lambda = \frac{Pr(F_n > c|H = 0)}{Pr(F_n > c|H = 1)} = \frac{1 - F_{k,d(n)}(c)}{1 - F_{k,d(n)}(c|\lambda)} , \quad (13)$$

and the same procedure follows to calculate the sample size needed. Here, we choose  $c$  to satisfy Equation (13). Using such a  $c$ , we calculate the power  $Pr(F_n > c|H = 1)$  and then plot the power against  $n$ . Figure 4 shows the resulting curves that are similar as those in Figure 2.

*2.2.3 Case for unequal  $\Delta$ 's and  $\sigma$ 's* So far, we have proceeded as if all genes have the same set of parameters. So the average power across all genes would be the same as the power for individual genes. In reality, each gene may have different set of parameters. If we use the two-sample comparison as an example, the gene-specific parameters include  $\sigma_g$ , the standard deviation, and  $\Delta_g$ , the true difference between the mean expressions of the treatment and the control group.

To study the realistic case when  $\Delta_g$  and  $\sigma_g$  depend on  $g$ , we assume that they follow some distribution with the probability density function  $\pi(\Delta_g, \sigma_g)$ .



**Figure 4.** For the design as in Figure 3, if we test all treatment means are equal with FDR controlled at 5%, the power of F test is shown for each sample size (number of sets of 3 slides). Here power is calculated when the true difference of log of gene expression levels between treatment I and treatment II to be 1 and difference between treatment I and treatment III to be  $-0.5$ ,  $\sigma = 1$ .

The distribution can be a parametric or nonparametric one that has been estimated from data of similar experiments. For example, when designing an experiment, a pilot study could be available, based on which the distribution of parameters can be estimated. In this case, our procedure can be extended to calculate a sample size while obtaining an average power across all genes. Here by average power, we mean the power integrated with respect to  $\pi(\Delta_g, \sigma_g)$ ,

$$\begin{aligned}
 & Pr(T \in \Gamma | H = 1) \\
 &= \int \int Pr(T \in \Gamma | H = 1, \Delta_g, \sigma_g) \pi(\Delta_g, \sigma_g) d\Delta_g d\sigma_g. \quad (14)
 \end{aligned}$$

Using Equation (14) and the argument similar to what leads to Equation (4),

we conclude that the FDR is  $\alpha$  if

$$\Lambda = \frac{Pr(T \in \Gamma | H = 0)}{\int \int Pr(T \in \Gamma | H = 1, \Delta_g, \sigma_g) \pi(\Delta_g, \sigma_g) d\Delta_g d\sigma_g} . \quad (15)$$

where

$$\Lambda = \frac{\alpha}{1 - \alpha} \frac{1 - \pi_0}{\pi_0} .$$

When we apply this to the  $t$ -tests, similar to Equation (7), Equation (15) becomes

$$\Lambda = \frac{Pr(|T_g| > c | H = 0)}{\int \int Pr(|T_g| > c | H = 1, \Delta_g, \sigma_g) \pi(\Delta_g, \sigma_g) d\Delta_g d\sigma_g} , \quad (16)$$

where the numerator equals  $2 \cdot T_{n_1+n_2-2}(-c)$  and the denominator equals

$$\begin{aligned} & 1 - \int \int T_{n_1+n_2-2}(c|\theta) \pi(\Delta_g, \sigma_g) d\Delta_g d\sigma_g \\ & + \int \int T_{n_1+n_2-2}(-c|\theta) \pi(\Delta_g, \sigma_g) d\Delta_g d\sigma_g . \end{aligned} \quad (17)$$

Note that  $\theta$  is as defined in (8). As before,  $T_d(\bullet|\theta)$  denotes the cumulative distribution function (cdf) of  $t$ -distribution. We then solve for the critical value  $c$  and apply the same procedure to get the sample size needed. The same technique extends to the  $F$ -tests or other tests of interest.

To illustrate our idea in more details, we assume that the mean difference expression level of differentially expressed genes,  $\Delta_g$ , follows a normal distribution and variances of expression levels for all genes follow an inverse gamma distribution:

$$\begin{aligned} \Delta_g & \sim N(\mu_\Delta, \sigma_\Delta^2), \\ \sigma_g^2 & \sim \text{Inverse Gamma}(a, b), \end{aligned}$$

and we use  $\pi_1(\Delta_g)$  and  $\pi_2(\sigma_g)$  to denote the p.d.f. of  $\Delta_g$  and  $\sigma_g$  respectively. Then we solve for  $c$  based on Equations (16) and (17) for specified level ( $\alpha$ ) of

FDR and proportion of non-differentially expressed genes ( $\pi_0$ ). This involves integrations. To deal with the integration, say in (17), the inner integral equals (see the Appendix B for derivation)

$$\begin{aligned} & \int T_{n_1+n_2-2} \left( c|\Delta_g / \sigma_g \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right) \pi_1(\Delta_g) d\Delta_g \\ &= T_{n_1+n_2-2} \left( \frac{c}{\sqrt{\frac{\sigma_\Delta^2}{\sigma_g^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} + 1}} \middle| \frac{\mu_\Delta}{\sqrt{\sigma_g^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \right). \end{aligned} \quad (18)$$

For the integration with respect to  $\sigma_g$ , we can apply adaptive Lobatto quadrature for numerical integration which allows a stable calculation to get the root of  $c$ . The calculation with this numerical integration provides answers instantly. Once we get answers of  $c$  for each sample size, we calculate power accordingly and find the needed sample size based on power.

### 3. Simulation

How realistic is the calculated sample size proposed in this paper? More specifically, if the desired power is 80%, FDR = 5% and our approach results in a sample size of 9 for the two-sample comparison with  $t$ -test, is there a statistical test that would actually achieve all the operating characteristics with 9 slides? To find out, we simulate data with calculated sample size and perform multiple testing with a FDR controlling procedure. Then we checked:

- whether the multiple testing actually results in desired power for the calculated sample size, and

- whether the observed FDR is comparable with the level that we want to control.

If we can find a statistical procedure that achieves the desired FDR and power at the calculated sample size, our procedure is then demonstrated to be practical. This is indeed the case.

There are several procedures to control FDR, such as the q-value procedure proposed by Storey and Tibshirani (2003) and Storey, Taylor and Siegmund (2004), and the procedures proposed by Benjamini and Hochberg (1995) and Benjamini and Hochberg (2000). These procedures all have the FDR conservatively controlled (Storey et al., 2004). For the purpose of simulation study, we apply the q-value procedure as outlined in Storey et al. (2004) to control FDR. The earlier version of the manuscript applied the procedure in Storey and Tibshirani (2003) and the results were similar to the report here.

We first test the proposed method when observations (genes) are independent of each other. In a microarray setting, we suppose there are a total of 5000 genes and we have equal sample size for the treatment and the control groups ( $n_1 = n_2 = n$ ). Gene specific variances,  $\sigma_g^2$ , are simulated from an inverse gamma distribution. Same as in Wright and Simon (2003), we chose  $1/\sigma^2 \sim \Gamma(3, 1)$  because this distribution approximates well several microarray data sets that we have been analyzing. For the control group, gene expression values are simulated from  $N(0, \sigma_g^2)$ . For the treatment group, we set  $\Delta_g=0$  for non-differentially expressed genes and simulate  $\Delta_g$  from  $N(2, \sigma_\Delta^2)$  for differentially expressed genes, then gene expression values are simulated from  $N(\Delta_g, \sigma_g^2)$ .

**Table 2**  
*Parameter values in Simulation Study*

Parameter	Values in Simulation
$\pi_0$	0.995, 0.95, 0.9, 0.8
$\sigma_\Delta$	0.2, 1, 2
$\rho$	0, 0.2, 0.5, 0.8

There are several parameters involved for the simulation,  $\pi_0$  (the proportion of non-differentially expressed genes),  $\sigma_\Delta$  (the standard deviation of effect size) and for the dependent case, the correlation coefficient  $\rho$ . To evaluate the accuracy of our sample size calculation method, we perform the simulation with a factorial design and the levels(values) of each factor(parameter) are summarized in Table 2. For each of the 48 parameter settings, the FDR is controlled at 5% for multiple testing.

For each parameter setting, we calculate the anticipated power for each sample size and generate the power curve as described in Section 2. We also simulate 200 sets of data and perform  $t$ -tests for each data set with  $q$ -value procedure (Storey et al., 2004) to control FDR. The observed power is averaged over the 200 simulated data sets and observed proportion of false discoveries is also recorded. Comparing with the simulation results, the anticipated power curves based on our calculation are almost indistinguishable from the simulation results for all investigated parameter settings. An example is shown in Figure 5(a). Hence, our proposed method provides accurate estimate of sample sizes. The observed FDR is also close to the controlled level, 5%, as in Figure 5(b), justifying the validity of the procedure in Storey

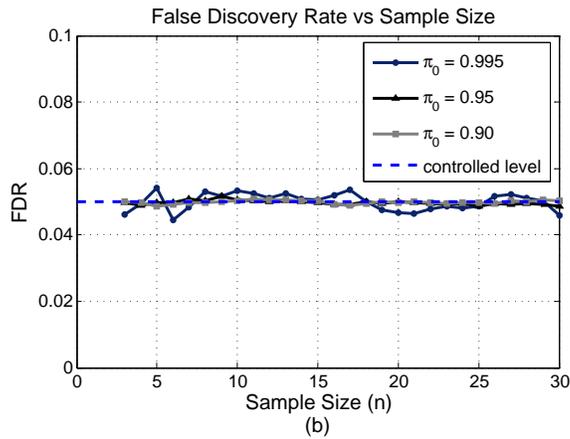
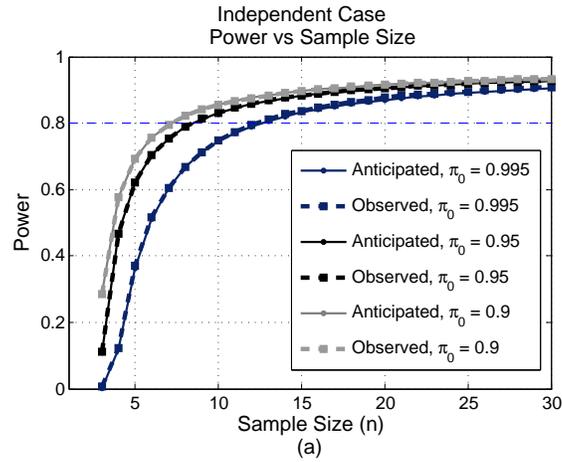
et al. (2004).

Since many genes may function as groups, it is very likely that dependencies exist in gene expression data. To check the performance of the proposed method when the assumption of independence is violated, gene expression levels are also simulated according to a dependence structure (Ibrahim et al., 2002). Then the same procedure of testing is applied and the resulting power curves are compared with our calculation.

More specifically, gene expression levels for differentially expressed genes are simulated in blocks of 25 according to the following hierarchical structure described in Section 4 of Ibrahim et al. (2002):

$$\begin{aligned}
\mu_X &\sim N(0, v_0^2) \\
\mu_Y &\sim N(2, v_0^2) \\
\mu_{Xg}|\mu_X &\sim N(\mu_X, \tau^2) \\
\mu_{Yg}|\mu_Y &\sim N(\mu_Y, \tau^2) \\
\sigma_g^2 &\sim \text{Inverse Gamma}(3, 1) \\
X_{gi}|\mu_{Xg} &\sim N(\mu_{Xg}, \sigma_g^2) \\
Y_{gi}|\mu_{Yg} &\sim N(\mu_{Yg}, \sigma_g^2).
\end{aligned}$$

where  $X_{gi}$  and  $Y_{gi}$  ( $g = 1, 2, \dots, G$ ,  $i = 1, 2, \dots, n$ ) are the gene expression levels for the control group (indexed with  $X$ ) and treatment group (indexed with  $Y$ ) respectively. For non-differentially expressed genes, we simulate  $\mu_{Xg}$  the same as above and set  $\mu_{Yg} = \mu_{Xg}$ , based on which we simulate the gene expression levels  $X_{gi}$  and  $Y_{gi}$ . Please note that the correlation coefficient,  $\rho$ , equals  $v_0^2/(v_0^2 + \tau^2)$  and  $\sigma_\Delta^2 = 2(v_0^2 + \tau^2)$  with  $\Delta_g = \mu_{Yg} - \mu_{Xg}$ . Examples of power curves are presented in Figure 6. For all 36 parameter settings of the



**Figure 5.** Simulation Results. (a) Observed power curves are plotted with dashed lines while the anticipated power curves based on our calculation are plotted with solid lines for different  $\pi_0$ 's. For all three  $\pi_0$ 's, the differences between the anticipated and observed power are almost indistinguishable. (b) Observed false discovery rates for the three parameter settings corresponding to (a) are plotted. The controlled level of 5% is indicated with the horizontal dashed line.

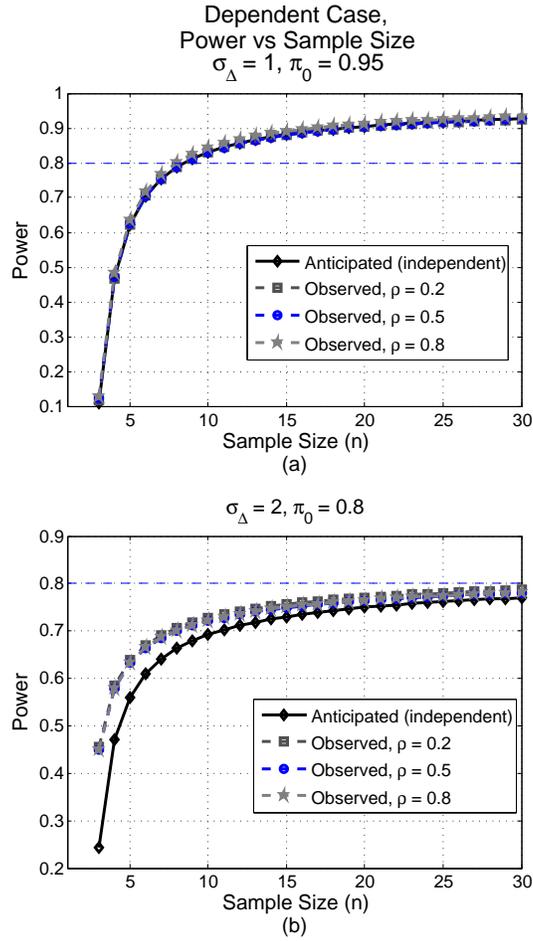
dependent case, 34 of them show results similar as in Figure 6 (a) , which demonstrate that the anticipated power approximates really well of the actual power. There are two settings that the discrepancy between anticipated power and calculation is relatively larger than the rest. Figure 6 (b) shows the worse one of the two. Still, the anticipated power based on our calculated sample size is close to the simulation result, and the difference shows that our method will provide a slightly conservative estimate of sample size.

When  $\Delta_g$  and  $\sigma_g^2$  are the same for all genes, simulation shows that our method can provide accurate sample size estimation both for independent genes and dependent data similarly as the simulation results shown above.

There are several papers addressing the question of calculating sample size while controlling FDR. Among these papers, Yang et al. (2003) and Pounds and Cheng (2005) provided clearly defined algorithms. We have compared our approach with these methods in the context of 2-sample  $t$ -test for fixed  $\Delta_g$  and  $\sigma_g^2$ . Table 3 shows that, the calculated sample size based on our proposed approach agrees with the actual sample size needed based on simulation results. Yang’s approach results in similar answers with ours except that in some case, it is a little conservative. Answers from Pounds and Cheng’s algorithm are too liberal in one situation (when  $\Delta/\sigma=1$ ) and deviate from the right answer a lot more than the other two methods.

#### 4. Discussion

The number of arrays included in microarray experiments directly affects the power of data analysis. It is critical to have a guideline to select a sample size. Because of the huge dimensionality associated with those data sets, controlling FWER is very conservative in many cases (Storey and Tibshirani,



**Figure 6.** Simulation Results. (a) Observed power curves are plotted with dashed lines while the anticipated power curve based on our calculation is plotted with a solid line. The anticipated power based on independence assumption approximates well the observed power for all three cases of different correlation coefficients. (a) Observed power curves are plotted with dashed lines while the anticipated power curve based on our calculation is plotted with a solid line. For these cases, the differences between the anticipated and observed power curves are relatively larger and our estimation for such cases is slightly conservative.

**Table 3**

*Comparison of sample size calculation methods including Yang's approach, Pounds and Cheng's approach (PC), the proposed method in this paper (LH) with the actual simulation result (Simu). The sample size is selected based on desired power of 80% and FDR at 5%.*

---



---

$\Delta/\sigma=2$	Yang's	PC	LH	Simu
$\pi_0 = 0.5$	8	7	6	6
$\pi_0 = 0.9$	10	10	9	10
$\pi_0 = 0.95$	11	11	11	11

---



---

$\Delta/\sigma=1$	Yang's	PC	LH	Simu
$\pi_0 = 0.5$	22	12	18	18
$\pi_0 = 0.9$	30	16	29	30
$\pi_0 = 0.95$	34	18	33	33

---



---

2003). Instead, FDR proposed by Benjamini and Hochberg (1995) and Storey (2002) seem to be a more appropriate error rate to control and has been widely applied to microarray analysis. Therefore, it is important to obtain a method to give the sample size that would control the FDR and guarantee a certain power.

The method is straightforward to apply as described in Section 2 for  $t$  and  $F$ -tests. The proposed method can be generalized to other tests, as long as there is an explicit form to calculate the type I error and power of an individual test. The method presented in this paper allows calculation for an accurate sample size with minimum effort when designing an experiment.

## Acknowledgement

The authors thank Dr. Gregory R. Warnes for insightful comments and suggestions. We also thank Dr. Chong Wang for pointing out the Lobatto Quadrature for numerical integration.

## REFERENCES

- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of Royal Statistical Society B* **57**, 289–300.
- Benjamini, Y. and Hochberg, Y. (2000). On adaptive control of the false discovery rate in multiple testing with independent statistics. *Journal of Educational and Behavioral Statistics* **25**, 60–83.
- Cui, X., Hwang, J., Qiu, J., Blades, N. J. and Churchill, . A. (2005). Improved statistical tests for differential gene expression by shrinking variance components estimates. *Biostatistics* **6**, 59–75.
- Dobbin, K. and Simon, R. (2005). Sample size determination in microarray experiments for class comparison and prognostic classification. *Biostatistics* **6**, 27–38.
- Heller, M. J. (2002). DNA microarray technology: devices, systems, and applications. *Annual Review in Biomedical Engineering* **4**, 129–153.
- Horton, J. D., Shah, N. A., Park, J. A. W. N. N. A. S. W., m. S. Brown and Goldstein, J. L. (2003). Combined analysis of oligonucleotide microarray data from transgenic and knockout mice identifies direct srebp target genes. *Proceedings of the National Academy of Sciences* **100**, 12027–12032.

- Hu, J., Zou, F. and Wright, F. A. (2005). Practical FDR-based sample size calculations in microarray experiments. *Bioinformatics* **21**, 3264–3272.
- Hua, J., Xiong, Z., Lowey, J., Suh, E. and Dougherty, E. R. (2005). Optimal number of features as a function of sample size for various classification rules. *Bioinformatics* **21**, 1509–1515.
- Ibrahim, J. G., Chen, M. and Gray, R. J. (2002). Bayesian models for gene expression with dna microarray data. *Journal of the American Statistical Association* **97**, 88–99.
- Jung, S.-H. (2005). Sample size for fdr-control in microarray data analysis. *Bioinformatics* **21**, 3097–3104.
- Mandel, S., Weinreb, O. and Youdim, M. B. H. (2003). Using cDNA microarray to assess parkinson’s disease models and the effects of neuroprotective drugs. *TRENDS in Pharmacological Sciences* **24**, 184–191.
- Pawitan, Y., Michiels, S., Koscielny, S., Gusnanto, A. and Ploner, A. (2005). False discovery rate, sensitivity and sample size for microarray studies. *Bioinformatics* **21**, 3017–3024.
- Pounds, S. and Cheng, C. (2005). Sample size determination for the false discovery rate. *Bioinformatics* **21**, 4263–4271.
- Quackenbush, J. (2002). Microarray data normalization and transformation. *Nature Genetics Supplement* **32**, 496–501.
- Smyth, G. K., Yang, Y. H. and Speed, T. (2003). Statistical issues in cDNA microarray data analysis. *Methods in Molecular Biology* **224**, 111–136.
- Storey, J. D. (2002). A direct approach to false discovery rates. *Journal of Royal Statistical Society B* **64**, 479–498.
- Storey, J. D., Taylor, J. E. and Siegmund, D. (2004). Strong control, con-

- servative point estimation and simultaneous rates: a unified approach. *Journal of Royal Statistical Society B* **66**, 187–205.
- Storey, J. D. and Tibshirani, R. (2003). Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences* **100**, 9440–9445.
- Wright, G. W. and Simon, R. M. (2003). A random variance model for detection of differential gene expression in small microarray experiments. *Bioinformatics* **19**, 2448–2455.
- Yang, M. C. K., Yang, J. J., McIndoe, R. A. and She, J. X. (2003). Microarray experimental design: power and sample size considerations. *Physiological Genomics* **16**, 24–28.
- Yang, Y. H. and Speed, T. (2003). *Design and analysis of comparative microarray experiments*. Chapman and Hall.

## Appendices

### *Appendix A: Derivation of Equation (4)*

Suppose that the test statistic is  $T$  and the rejection region is  $\Gamma$ . Let  $H = 0$  if null hypotheses is true and  $H = 1$  if alternative hypothesis is true. We make the same assumptions as Theorem 1 in Storey (2002), that is, all tests are identical, independent and the probability that the hypothesis is null is  $Pr(H_0) = Pr(H = 0) = \pi_0$ . Based on the Bayes rule, the FDR is

$$\begin{aligned}
 & Pr(H_0|T \in \Gamma) \\
 = & \frac{Pr(T \in \Gamma|H_0) \cdot \pi_0}{Pr(T \in \Gamma|H_0) \cdot \pi_0 + Pr(T \in \Gamma|H_1) \cdot (1 - \pi_0)}
 \end{aligned}
 \tag{19}$$

where  $H_0$  and  $H_1$  denote  $H = 0$  and  $H = 0$  respectively. To control the level of FDR to be  $\alpha$  or less, we set (19) to be less than or equal to  $\alpha$ , which is equivalent to:

$$\frac{\alpha}{1 - \alpha} \frac{1 - \pi_0}{\pi_0} \leq \frac{Pr(T \in \Gamma | H = 0)}{Pr(T \in \Gamma | H = 1)}. \quad (20)$$

*Appendix B: Derivation of Equation (18)*

If we let  $Z$  stands for a random variable with the standard normal distribution, and let  $t_v(\theta)$  denotes the non-central t distribution with  $v$  degrees of freedom and non-centrality parameter  $\theta$ , then a random variable  $X$  which follows distribution  $t_v(\theta)$  can be viewed as:

$$X \stackrel{d}{=} \frac{Z + \theta}{\sqrt{\frac{\chi_v^2}{v}}} \quad (21)$$

where “ $\stackrel{d}{=}$ ” denotes that the two random variables have the same distribution. In the case of two-sample t-test as in Section 2.2.1, for given  $\Delta_g$  and  $\sigma_g$ , we know that

$$T_{n,g} = \frac{\bar{x}_g - \bar{y}_g}{S_g * \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2} \left( \frac{\Delta_g}{\sigma_g \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right)$$

Then

$$T_{n,g} \stackrel{d}{=} U \Big/ \sqrt{\frac{\chi_{n_1+n_2-2}^2}{n_1 + n_2 - 2}}$$

where

$$U = Z + \frac{\Delta_g}{\sigma_g \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

For given  $\sigma_g$ , we assume that

$$\Delta_g \sim N(\mu_\Delta, \sigma_\Delta^2),$$

then

$$U \sim N\left(\frac{\mu_\Delta}{\sigma_g \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}, \frac{\sigma_\Delta^2}{\sigma_g^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} + 1\right).$$

If we scale  $T_{n,g}$  to obtain another non-central t-distribution, we get

$$\begin{aligned} \frac{T_{n,g}}{\sqrt{\frac{\sigma_\Delta^2}{\sigma_g^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} + 1}} &\stackrel{d}{=} \frac{U}{\sqrt{\frac{\sigma_\Delta^2}{\sigma_g^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} + 1}} \bigg/ \sqrt{\frac{\chi_{n_1+n_2-2}^2}{n_1+n_2-2}} \\ &\stackrel{d}{=} \frac{Z+p}{\sqrt{\frac{\chi_{n_1+n_2-2}^2}{n_1+n_2-2}}} \end{aligned}$$

where the non-centrality parameter

$$\begin{aligned} p &= \frac{\mu_\Delta}{\sigma_g \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \bigg/ \sqrt{\frac{\sigma_\Delta^2}{\sigma_g^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} + 1} \\ &= \frac{\mu_\Delta}{\sigma_\Delta^2 + \sigma_g^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}. \end{aligned}$$

Based on the above result, we can avoid integration with respect to  $\Delta_g$ , instead, we will have

$$\begin{aligned} &\int Pr(T_{n,g} < c | \Delta_g) \pi_1(\Delta_g) d\Delta_g \\ &= T_{n_1+n_2-2} \left( \frac{c}{\sqrt{\frac{\sigma_\Delta^2}{\sigma_g^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} + 1}} \mid p \right) \end{aligned}$$

In the case of  $n_1 = n_2 = n$ ,

$$\begin{aligned} & \int \Pr (T_{n,g} < c | \Delta_g) \pi_1(\Delta_g) d\Delta_g \\ &= T_{2n-2} \left( \frac{c}{\sqrt{\frac{n \cdot \sigma_\Delta^2}{2 \cdot \sigma_g^2} + 1}} \mid \frac{\mu_\Delta}{\sigma_\Delta^2 + \frac{2\sigma_g^2}{n}} \right) \end{aligned}$$

where  $T_d(\bullet|\theta)$  denote the cumulative distribution function (c.d.f.) of  $X$  in (21).