

QVALUE: The Manual

Version 1.0

Alan Dabney and John D. Storey
Department of Biostatistics
University of Washington
Email: jstorey@u.washington.edu

March 2003; Updated June 2003; Updated January 2004

Table of Contents

| | |
|--|---------|
| 1. Introduction | Page 1 |
| 2. Citing this software..... | Page 1 |
| 3. How to install QVALUE..... | Page 2 |
| 4. How to use QVALUE in point-and-click mode..... | Page 2 |
| 5. How to use QVALUE in command line mode..... | Page 4 |
| 6. How to use the q-value plots to make decisions..... | Page 7 |
| 7. What is a q-value? (A primer) | Page 8 |
| 8. Frequently asked questions..... | Page 10 |
| 9. References..... | Page 12 |

1. Introduction

This document provides instructions for how to use the QVALUE software package, as well as a short tutorial on false discovery rates and q-values. If you are unfamiliar with false discovery rates and q-values, then it may be helpful to read Section 7 on page 8 first.

2. Citing this software

The QVALUE software involves research done by David Siegmund, John Storey, Jonathan Taylor, and Rob Tibshirani. The software was written by Alan Dabney and John Storey. Please cite at least one of the following articles when reporting results from the software.

Storey JD. (2002) A direct approach to false discovery rates. *Journal of the Royal Statistical Society, Series B*, **64**: 479-498.

Storey JD and Tibshirani R. (2003) Statistical significance for genome-wide experiments. *Proceeding of the National Academy of Sciences*, **100**: 9440-9445.

Storey JD, Taylor JE, and Siegmund D. (2004) Strong control, conservative point estimation, and simultaneous conservative consistency of false discovery rates: A unified approach. *Journal of the Royal Statistical Society, Series B*, **66**: 187-205.

Note that Benjamini & Hochberg (1995) propose the false discovery rate concept and provide the first step-wise p-value method to control it. Storey (2003) generally defines the q-value and initially studies its properties in a Bayesian context.

3. How to install QVALUE

QVALUE runs on top of the free statistical software program R. Please go to <http://faculty.washington.edu/~jstorey/qvalue/> for specific details on how to install R and QVALUE for Linux/Unix, Macintosh, or Windows.

4. How to use QVALUE in point-and-click mode

Step 1. Save your p-values into a text file, preferably with one p-value per line. This can be done in Excel, for example, by creating a worksheet with the p-values listed in a single column and saving the worksheet as a tab-delimited text file.

Step 2. Start R. This can be done at the command line, or by clicking on the R icon on your desktop. At the prompt in R, type:

```
> library(qvalue)
> qvalue.gui()
```

Something like Figure 1 will appear on your screen.

Note: For Windows users, we have provided an installation method where one only has to click on a desktop icon in order to start the software package. See <http://faculty.washington.edu/~jstorey/qvalue/> for more on this.

Step 3. In order to load the p-values, press the Browse button and select the file containing the p-values. Then press Load.

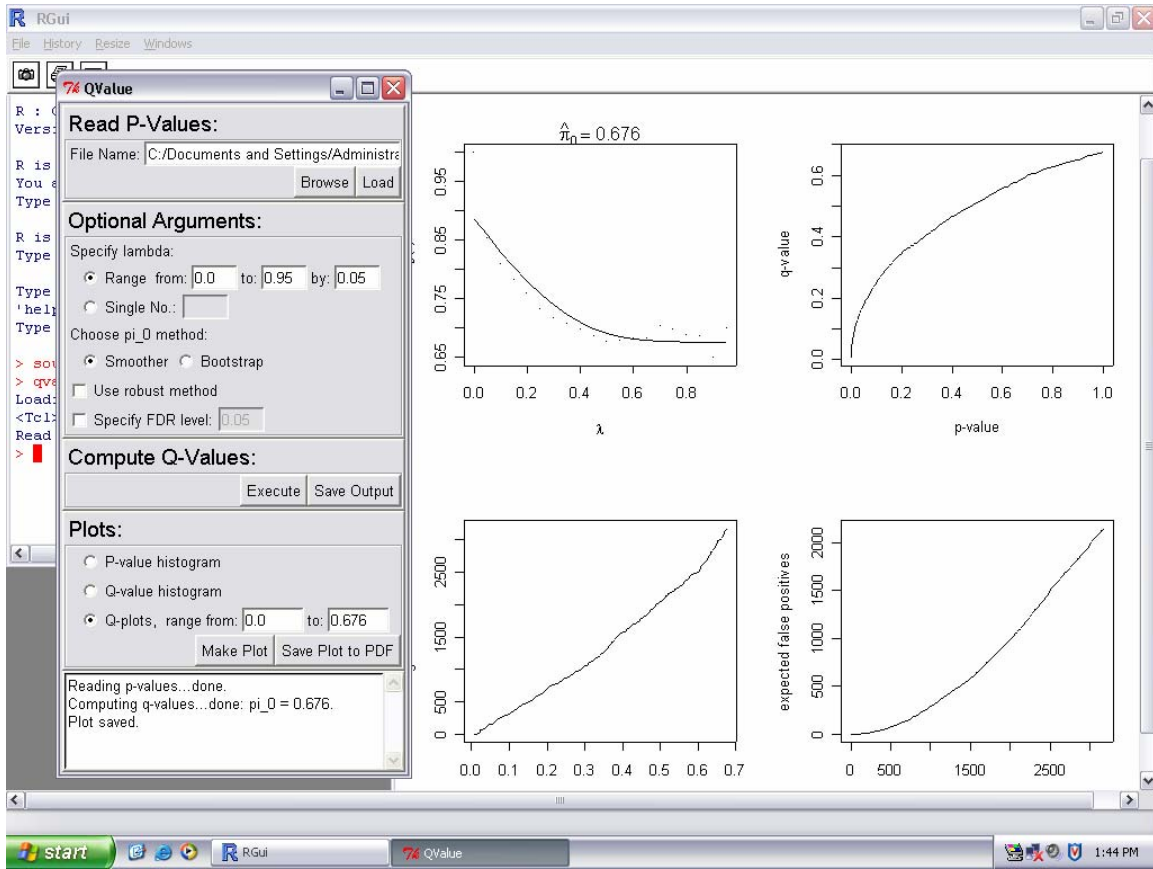


Figure 1. A screenshot of the QVALUE point-and-click interface in Windows.

Step 4. Several optional parameters can be set that affect how the q-values are estimated. The most delicate aspect of this process is estimating π_0 , where π_0 is the *overall* proportion of true null hypotheses.

- The *Specify lambda* option allows one to choose the range of the tuning parameter used in estimating π_0 (see Storey 2002 and Section 5 below).
- The *Choose pi_0 method* allows one to decide between using the smoother method (Storey & Tibshirani 2003) or the bootstrap method (Storey, Taylor & Siegmund 2004) when estimating π_0 .
- If the *Use robust method* is selected, then a q-value estimate is used that is more robust for small p-values and a direct finite sample estimate of pFDR (Storey 2002).
- The *Specify FDR level* option writes a true/false indicator into the results file that indicates for each test whether its estimated q-value is less than or equal to that level or not.

Step 5. Given that p-values have been successfully loaded (Step 3), click **Execute** under the *Compute Q-values* section to calculate the q-value estimates under the current optional settings. If the optional settings are changed, then simply click **Execute** again to update the estimates.

Step 6. The *Plots* section allows one to view three types of plots. Once the p-values are successfully loaded, then selecting `P-value histogram` and clicking `Make Plot` will create a histogram of the p-values. Similarly, once the q-values have been calculated, then selecting `Q-value histogram` and clicking `Make Plot` will create a histogram of the q-values. Selecting `Q-plots` and clicking `Make Plot` will produce a plot of the four “Q-plots” discussed in Storey & Tibshirani (2003). [See also Section 6 of this manual for more on the Q-plots.] The range of q-values considered in these plots can be set by the user; the range is automatically set to be between zero and the estimate of π_0 otherwise. Any plot can be saved as a PDF file by clicking `Save Plot to PDF` when the plot is currently displayed.

Step 7. Once you are satisfied with your q-value estimates (perhaps after looking at some plots and adjusting the optional settings), then click `Save Output` in the *Compute Q-values* section to save the q-values, p-values, and π_0 estimate into a text file. The order will be the same as the original file containing the p-values.

The most delicate aspect of this software is choosing how to estimate π_0 in the optional arguments. If no options are selected, then by default the smoother method proposed in Storey & Tibshirani (2003) is used. Our experience indicates that this often works better than the bootstrap method, but can backfire for a small number of p-values or in pathological situations. An overall safer option is the bootstrap method proposed in Storey, Taylor & Siegmund (2004). If one selects the single value `lambda=0`, then this produces the estimate implicit in the Benjamini & Hochberg (1995) methodology. In particular, setting `lambda=0` estimates π_0 to be 1. This can be viewed as a special conservative case of the Storey (2002) methodology, so at the very least we recommend using the bootstrap method rather than setting `lambda` to a single predetermined number, such as `lambda=0`. Very importantly, one can use the first plot of the Q-plots (a plot of the π_0 estimate versus its tuning parameter λ) in order to gauge the quality of the final π_0 estimate.

5. How to Use QVALUE in command line mode

Step 1. Save your p-values into a text file, preferably with one p-value per line. In this section, we will call this file `pvalues.txt`. This can be done in Excel, for example, by creating a worksheet with the p-values listed in a single column and saving the worksheet as a tab-delimited text file.

Step 2. Start R. This can be done at the command line, or by clicking on the R icon on your desktop. At the prompt in R, type the following in order to load the library of functions:

```
> library(qvalue)
```

Step 3. Select `File` → `Change directory`. Select the directory where you have stored `pvalues.txt`. Now type the following commands:

```

> p <- scan("pvalues.txt")
> qobj <- qvalue(p)
> qplot(qobj)
> qwrite(qobj, filename="myresults.txt")

```

The first line saves the p-values into the R object `p`. The second command saves the output from the main q-value function into the R object `qobj`. The third command makes four useful plots that can be used to assess which significance cut-offs make sense for your study. These plots are labeled and self-explanatory, and they are also discussed in the next section. The fourth command writes the results to a file called `myresults.txt`, which will be written in the same directory as `pvalues.txt`. The file contains the function call used and the estimate of π_0 , where π_0 is the *overall* proportion of true null hypotheses. (The false discovery rate is the proportion of true null hypotheses among those called significant, and π_0 is the proportion of true null hypotheses among all tests. See Section 5 on page 6.) Starting on the third line, the file lists each p-value and corresponding estimated q-value, one per line in the same order as `pvalues.txt`.

Several other arguments can be used in the function `qvalue`. The following lists all the possible arguments, with a description of each:

- `p`: A vector of p-values. *This is the only necessary input.*
- `lambda`: The values of the tuning parameter to be considered in estimating π_0 . These must be in $[0,1]$ and are set to `lambda=seq(0, 0.95, 0.05)` by default. Optional; see Storey (2002) for more information.
- `pi0.meth`: Either "smoother" or "bootstrap"; the method for automatically choosing tuning parameter `lambda` in the estimate of π_0 . If the `lambda` argument above is only given one value, then this option is ignored. Optional; the choice "smoother" is the default choice.
- `fdr.level`: The level at which to control the false discovery rate. Optional; if this is selected, a vector of TRUE and FALSE is returned that specifies whether each q-value is less than `fdr.level` or not.
- `robust`: An indicator of whether it is desired to make the estimate more robust for small p-values. This uses the point estimate of the "positive false discovery rate" (pFDR). Optional; see Storey (2002) for more information.

The most delicate aspect of this software is choosing how to estimate π_0 via `lambda` and `pi0.meth`. If no options are selected, then by default the smoother method (`pi0.meth="smoother"`) proposed in Storey and Tibshirani (2003) is used. My experience indicates that this often works better than the bootstrap method, but can backfire for a small number of p-values or in pathological situations. An overall safer option is the bootstrap method (`pi0.meth="bootstrap"`) proposed in Storey, Taylor & Siegmund (2002). If one selects `lambda=0`, then this produces the estimate implicit in the Benjamini and Hochberg (1995) methodology. In particular, setting

$\lambda=0$ estimates π_0 to be 1. This can be viewed as a special conservative case of the Storey (2002) methodology, so at the very least we recommend using `pi0.meth="bootstrap"` rather than setting `lambda` to some predetermined number, such as `lambda=0`. Here are three examples of using `qvalue` with non-default options:

```
> qobj <- qvalue(p, lambda=seq(0.2,0.8,0.01), robust=TRUE)
> qobj <- qvalue(p, lambda=0, fdr.level=0.05)
> qobj <- qvalue(p, pi0.meth="bootstrap")
```

The function `qplot` has an option to change the range of q-values for which the plots can be made. If one wants to view a range of q-values in, say 0 to 0.3, then type:

```
> qplot(qobj, rng=0.3)
```

The function `qwrite` currently has no options other than designating the file name, as was done above.

We advocate reporting the estimated q-value for each test. However, sometimes one wants to estimate the false discovery rate incurred for a given p-value cut-off, or estimate the p-value cut-off to control the false discovery rate at a certain level. Below are instructions on how to do this in `QVALUE`.

Estimating the false discovery rate for a given p-value cut-off. If one wants to estimate the false discovery rate when calling all p-values less than or equal to 0.01 significant, then type:

```
> max(qobj$qvalues[qobj$pvalues <= 0.01])
```

This calculates the maximum estimated q-value among all p-values less than or equal to 0.01, which is equivalent to estimating the false discovery rate when calling all p-values less than or equal to 0.01 significant. Clearly, if a cut-point different than 0.01 is desired, then replace 0.01 in the above command with that number.

Estimating a p-value cut-off for a given false discovery rate level. If one wants to estimate the p-value cut-off for controlling the false discovery rate at level 0.05, then type:

```
> max(qobj$pvalues[qobj$qvalues <= 0.05])
```

This calculates the largest p-value with estimated q-value less than or equal to 0.05. If we set `lambda=0`, then this is equivalent to the Benjamini & Hochberg (1995) step-wise p-value method. If π_0 is estimated (rather than set to 1), then this is equivalent to the false discovery rate controlling procedure proposed in Storey, Taylor & Siegmund (2002). Clearly, if a level of false discovery rate control other than 0.05 is desired, then replace 0.05 in the above command with the desired number.

6. How to use the q-values to make decisions

Here, we give some concise guidelines for interpreting the output of the software. A more thorough discussion in the context of genomics can be found in Storey & Tibshirani (2003).

One very important number that is obtained with the software is an estimate of the overall proportion of true null hypotheses π_0 . In the point-and-click interface, this estimate is printed in the dialogue box when the q-values are estimated. This estimate can be accessed at the command line by:

```
> qobj$pi0
```

Clearly, an estimate of *the proportion of significant tests* is one minus this number. This is a useful number to know, even if all the truly significant tests cannot all be explicitly identified. The p-values and q-values can be accessed in the file where the results are written.

If one wants to “control” the false discovery rate at a pre-determined level α , then calling all tests significant with estimated q-values $\leq \alpha$ accomplishes this under certain mathematical assumptions, including some cases where the p-values are dependent (Storey, Taylor & Siegmund 2002). In other words, we guarantee in some sense that

$$\frac{\text{\# of false positives}}{\text{\# of significant } t \text{ tests}} \leq \alpha$$

by calling all tests significant with estimated q-values $\leq \alpha$. One can automatically denote whether the estimated q-values are less than or equal to some α by using the FDR level option.

The more likely case is that one will want to investigate the overall behavior of the estimated q-values before making such a decision. The `Q-plots` command in point-and-click mode and the `qplot` function in the command line mode allow one to view several useful plots:

1. The estimated π_0 versus the tuning parameter λ
2. The q-values versus the p-values
3. The number of significant tests versus each q-value cut-off
4. The number of expected false positives versus the number of significant tests

The main purpose of the first plot is to gauge how reliable the estimate of π_0 is. Basically, a tuning parameter λ has to be chosen to estimate π_0 . The variable λ is called `lambda` in the software; as stated above it can be fixed or automatically chosen. The estimated π_0 is plotted versus the tuning parameter λ . As λ gets larger, the bias of the estimate decreases, yet the variance increases. See Storey (2002) for more on this. Comparing your final estimate of π_0 to this plot gives a good sense as to its quality. A smoother is fit to the plot

in order to elucidate the trend of the estimates. The remaining plots show how many tests are significant, as well as how many false positives to expect for each q-value cut-off. A thorough discussion of these plots can be found in Storey & Tibshirani (2003).

Finally, note that the most informative approach is to report the estimated q-value with each test, rather than making potentially arbitrary decisions about cut-offs for significance.

7. What is a q-value? (A primer)

The q-value is similar to the well known p-value. It gives each hypothesis test a measure of significance in terms of a certain error rate. The p-value of a test measures the minimum *false positive rate* that is incurred when calling that test significant. Likewise, the q-value of a test measures the minimum *false discovery rate* that is incurred when calling that test significant.

Whereas the p-value is commonly used for performing a single significance test, the q-value is useful for assigning a measure of significance to each of *many* tests performed *simultaneously*. (An example is testing thousands of genes for differential expression using DNA microarray data.) For each of these tests, there is a *null hypothesis* tested against an *alternative hypothesis*. A measure of significance therefore roughly measures how much a single test deviates from the null. The false positive rate and false discovery rate accomplish this quite differently.

A *false positive* is the term used to describe rejecting the null hypothesis (i.e., calling the test significant) when it is really true. Suppose we have defined a rule for calling tests significant. The false positive rate of the rule can then be loosely described by:

$$\text{false positive rate} \approx \frac{\# \text{ of false positives}}{\# \text{ of true null tests}} .$$

Therefore, the false positive rate measures the proportion of true null hypotheses that were (incorrectly) called significant by this rule.

A *false discovery* is also a false positive, however, the different terminology stresses the fact that we are concerned with false positives among the significant tests (i.e., the discoveries). The false discovery rate is the expected proportion of false positives among the tests found to be significant. The false discovery rate can then be loosely described by:

$$\text{false discovery rate} \approx \frac{\# \text{ of false positives}}{\# \text{ of significant } t \text{ tests}} .$$

The false positive rate and the false discovery rate therefore tell us two very different things about a method for calling tests significant. For a single test, the false positive rate can be useful for measuring how likely it is for a truly null case to be as significant as what has been observed. However, for many tests this is not as useful. For example,

suppose we decide that we can live with a false positive rate of 5%. Then about 5% of the time, we will call a truly null hypothesis significant. If we perform 1000 tests at a 5% false positive rate, then we can expect up to 50 false positives. This will typically be too many in practical situations.

When performing many significance tests, the false discovery rate gives more useful information. If we are willing to incur a false discovery rate of 5%, then this means that among all tests we call significant, about 5% of them will be false positives. If there are 100 significant tests, then this results in about 5 false positives; 500 significant tests results in about 25 false positives, etc.

If all tests are called significant then the false positive rate = 1 since all tests are called significant, and therefore all true null hypotheses are called significant. On the other hand, the false discovery rate is

$$\pi_0 \equiv \frac{\# \text{ of true null tests}}{\# \text{ of total tests}}$$

when all tests are called significant. The quantity π_0 is the overall proportion of true null hypotheses in the study. This is a useful number to consider as well as $\pi_1 \equiv 1 - \pi_0$, which is the proportion of significant results in the study. An estimate π_0 of is provided in the software.

In most significance testing situations, the null hypothesis is defined in such a way that either the null distribution of the test statistic is known (e.g., the null distribution of a t-test is the t distribution when the data are normal) or the null distribution can be simulated (e.g., via permutations or the bootstrap). Regardless of the method used, the false positive rate is easily measured, making it straightforward to obtain p-values. The false discovery rate, on the other hand, involves information about the false null hypotheses. Therefore to make a precise false discovery rate calculation, we would have to know which tests are truly significant and what their alternative distributions are.

False discovery rates methods can be described without loss of generality in terms of p-value. Specifically, Storey (2001) has shown that the q-value is the same whether we estimate it from the original statistics or from their corresponding p-values. Therefore, the software available here calculates q-values based on p-values. Because of the ease at which p-values can be obtained, this is a useful way to make the methods widely available.

Storey (2002) has developed methods for estimating false discovery rates that can be applied in a variety of ways. Rather than using only information from the null distribution, it utilizes information from all the p-values at once. For a given p-value threshold, say 5%, the false discovery rate is estimated in such a way that on average this estimate will exceed the true false discovery rate. This is a good property – we don't want to report a smaller false discovery rate than truly exists. Recently, it has been shown that this same estimate can be used to pick a false discovery rate beforehand, say 1%, and find the p-value threshold that guarantees on average that the true false discovery rate will be

less than or equal to the desired level (Storey, Taylor & Siegmund 2002). This is also a desirable property.

Fixing a significance threshold beforehand or fixing the false discovery rate beforehand may be useful in some situations. What is most general and useful however, is a test-specific false discovery rate measure. This essentially allows us to look at all possible thresholds at once, as well as providing each test with a measure of significance that can be easily interpreted. *This is exactly what the q-value accomplishes.* For a given test, we estimate the q-value by calculating the minimum estimated false discovery rate among all thresholds at which the false discovery rate is called significant. Conditions under which the q-values are simultaneously conservative have been given in Storey, Taylor & Siegmund (2002). If this property holds, then one can consider all q-values simultaneously without worrying about incurring bias. A neat Bayesian posterior probability view of q-values has been shown in Storey (2001), which gives the origin of its name.

See the recent talk at <http://faculty.washington.edu/~jstorey/qvalue/talk.pdf> for another short introduction to false discovery rates and q-values, including a brief summary of the formulas used in this software.

8. Frequently asked questions

A. I have a study with the following design ... and I formed p-values by Is your software appropriate for my study?

Unfortunately, we only have time to answer questions about the software. As long as your p-values are correctly calculated, then this software should provide a decent guide for significance in terms of false discovery rates. If one can show that the null p-values are independent or that the weak dependence criteria are satisfied, then it is possible to claim “strong control” of the false discovery rate for a given significance cut-off.

B. I cannot get R installed properly. What am I doing wrong?

There is extensive online help available at <http://www.r-project.org/> for installing R.

C. I am using Windows, and the QVALUE window keeps disappearing behind the R window. How can I stop this?

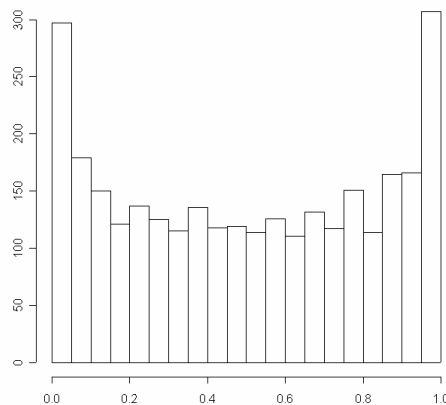
Our point-and-click interface is written using a library in R called `tcltk`. On some Windows machines, there appears to be a bug in the way that a `tcltk` based window interacts with the R window. Hopefully, this will be fixed in future releases of R.

D. Why am I getting an error that says my π_0 estimate is less than zero?

Usually this indicates that the p-values were incorrectly calculated. Sometimes, π_0 is very small, in which case this error can be avoided by using the bootstrap option rather than the smoother option.

E. Why does the histogram of my p-values look so different from the one in the Storey & Tibshirani (2003) *PNAS* paper?

There are many possibilities. There could be a difference in power between your data and the data used in Storey & Tibshirani (2003). Your p-values could have been calculated incorrectly. One interesting case is when the p-value histogram looks like the following:



Usually, this means that a one-sided test was performed when there is actually a two-sided signal. (Recall, that a flat portion of the histogram indicates the p-values are uniformly distributed there.) For example, one could test for over-expression while there is differential expression in both the over-expressed and under-expressed directions. The estimation methodology assumes that the null p-values are uniformly distributed or are a more conservative version of the uniform distribution. It can be seen in the above histogram that null p-values are actually pulled towards 1. Therefore, the methodology is still valid, however, it is very conservative. This is not due to carelessness in developing the methodology. Rather, it is a drawback of frequentist hypothesis testing. One can avoid this problem by more aggressively estimating the null distribution, taking into account that the signal in one direction is null. When your p-value histogram looks as above, the *bootstrap* method for estimating π_0 is the appropriate method to use – do not use the *smoother* method.

F. Can I modify your software and re-distribute it?

QVALUE has a LGPL license. Please see <http://www.gnu.org/> for the specific details of such a license. In most situations the software can be modified and redistributed as long as it stays open source and Alan Dabney and John Storey are given credit.

9. References

Benjamini Y and Hochberg Y. (1995) Controlling the false discovery rate: A practical and powerful approach to multiple testing. . *Journal of the Royal Statistical Society, Series B*, **57**: 289-300.

Ihaka R and Gentleman R. (1996) R: A language for data analysis and graphics. *Journal of Computational and Graphical Statistics*, **5**: 299-314.

Storey JD. (2002) A direct approach to false discovery rates. *Journal of the Royal Statistical Society, Series B*, **64**: 479-498.

Storey JD and Tibshirani R. (2003) Statistical significance for genome-wide experiments. *Proceeding of the National Academy of Sciences*, **100**: 9440-9445.

Storey JD. (2003) The positive false discovery rate: A Bayesian interpretation and the q-value. *Annals of Statistics*, **31**: 2013-2035.

Storey JD, Taylor JE, and Siegmund D. (2004) Strong control, conservative point estimation, and simultaneous conservative consistency of false discovery rates: A unified approach. *Journal of the Royal Statistical Society, Series B*, **66**:187-205.