

edge:
Extraction of Differential Gene Expression
Version 0.99.0

John D. Storey^{*1}, Jeffrey T. Leek² and Andrew J. Bass¹

¹Princeton University
²John Hopkins University

March 25, 2015

Contents

1	Introduction	3
2	Citing this package	4
3	Getting help	5
4	Quick start guide	5
5	Case study: static experiment	6
5.1	Importing the data	7
5.2	Creating the alternative and null models	7
5.3	The <code>edgeSet</code> object	8
5.4	Fitting the data	9
5.5	Significance analysis	11
5.5.1	Likelihood ratio test	11
5.5.2	Optimal discovery procedure	11
5.6	Significance results	11
6	Case study: independent time course experiment	14
6.1	Importing the data	14
6.2	Creating the alternative and null models	15
6.3	The <code>edgeSet</code> object	16
6.4	Fitting the data	17
6.5	Significance analysis	18
6.5.1	Likelihood ratio test	19
6.5.2	Optimal discovery procedure	19
6.6	Significance results	19

^{*}<http://genomine.org/contact.html>

7	Case study: longitudinal time course experiment	22
7.1	Importing the data	22
7.2	Creating the alternative and null models	22
7.3	The <code>edgeSet</code> object	24
7.4	Fitting the data	25
7.5	Significance analysis	26
7.5.1	Likelihood ratio test	26
7.5.2	Optimal discovery procedure	26
7.6	Significance results	27
8	sva: Surrogate variable analysis	29
9	snm: Supervised normalization of microarray data	31
10	qvalue: Estimate the q-values	33
11	Advanced topic: Using the <code>ExpressionSet</code> object	34

1 Introduction

edge is a package for significance analysis of DNA micro-array experiments and is able to identify genes that are differentially expressed between two or more different biological conditions (e.g., healthy versus diseased tissue). **edge** performs significance analysis by using a new method developed by Storey (2007) called the optimal discovery procedure (ODP). Whereas previously existing methods employ statistics that are essentially designed for testing one gene at a time (e.g., t-statistics and F-statistics), the ODP-statistic uses information across all genes to test for differential expression. Storey et al. (2007) shows that the ODP is a more intuitive, often times more powerful, approach to multiple hypothesis testing when compared to traditional methods. The improvements in power from using the optimal discovery procedure are substantial; Figure 1 shows a comparison between **edge** and five leading software packages based on the Hedenfalk et al. (2001) breast cancer expression study.

edge also implements strategies that have been specifically designed for time course experiments. Many things can go wrong when using methods that have been designed for static experiments, and even though some significance analysis packages allow for users to enter information about time points, Storey et al. (2005) developed a procedure that simplifies the modelling process for time course experiments. In addition to identifying differentially expressed genes in both static and time course studies, **edge** includes implementations of popular packages such as **snn**, **sva** and **qvalue** to help simplify the analysis process for researchers.

The rest of the document details how to use **edge** in three different case studies: static, independent time course and longitudinal time course. For additional information regarding the optimal discovery procedure or the Storey et al. (2005) methodology for time course experiments, see section 2.

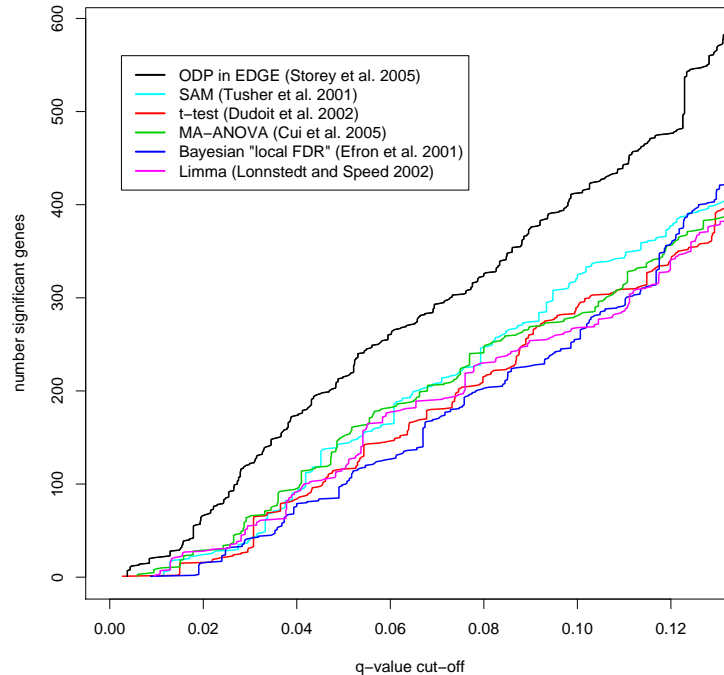


Figure 1: Comparison of EDGE to various other leading methods for identifying differential expressed genes in the [Hedenfalk et al. \(2001\)](#) study. Figure is from [Leek et al. \(2006\)](#).

2 Citing this package

[1] John D. Storey. The optimal discovery procedure: a new approach to simultaneous significance testing. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(3):347–368, 2007. ISSN 1467-9868. doi: 10.1111/j.1467-9868.2007.005592.x. URL <http://dx.doi.org/10.1111/j.1467-9868.2007.005592.x>

Theory paper that introduces the optimal discovery procedure and shows that it maximizes the expected true positive results for each number of fixed false positive results. The optimality is closely related to the false discovery rate.

[2] John D. Storey, James Y. Dai, and Jeffrey T. Leek. The optimal discovery procedure for large-scale significance testing, with applications to comparative microarray experiments. *Biostatistics*, 8(2):414–432, 2007. doi: 10.1093/biostatistics/kxl019. URL <http://biostatistics.oxfordjournals.org/content/8/2/414.abstract>

Discusses various ways of estimating the ODP statistic with applications to microarray experiments.

[3] Sangsoo Woo, Jeffrey T. Leek, and John D. Storey. A computationally efficient modular optimal discovery procedure. *Bioinformatics*, 27(4):509–515, 2011. doi: 10.1093/bioinformatics/btq701. URL <http://bioinformatics.oxfordjournals.org/content/27/4/509.abstract>

Previous implementations of the ODP are computationally infeasible for a large number of hypothesis tests. This paper introduces a computationally efficient implementation of ODP that this package is based on.

[4] John D. Storey, Wenzhong Xiao, Jeffrey T. Leek, Ronald G. Tompkins, and Ronald W. Davis. Significance analysis of time course microarray experiments. *Proceedings of the National Academy of Sciences of the United States of America*, 102(36):12837–12842, 2005. doi: 10.1073/pnas.0504609102. URL <http://www.pnas.org/content/102/36/12837.abstract>

A methodology for analyzing time course microarray data is introduced and applied to two time course studies on humans.

3 Getting help

Hopefully, most questions relating to the package will be answered in the vignette but to get a more detailed account of how to use the functions simply type within R:

```
help(package = "edge")
```

Please contact the authors directly with any issues regarding bugs. Otherwise, any questions or problems implementing **edge** will most efficiently be addressed on the Bioconductor mailing list, <http://stat.ethz.ch/mailman/listinfo/bioconductor>.

4 Quick start guide

To get started, first load the **kidney** dataset included in the package:

```
library(edge)
data(kidney)
names(kidney)

## [1] "age"      "sex"      "kidexpr"
```

The **kidney** study is interested in determining differentially expressed genes in the kidney as it ages. The **age** variable is the age of the subjects, the **sex** variable is whether the subjects were male or female and the **tissue** variable is whether the tissue sample is from the cortex or medula. The expression values for the genes are contained in the **kidexpr** variable. In this example, we are only interested in the cortex samples:

```
kidexpr <- kidney$kidexpr
age <- kidney$age
sex <- kidney$sex
```

Once the data has been loaded, the user has two options to create the experimental models: **edgeModel** or **edgeStudy**. If the experiment models are unknown to the user, **edgeStudy** can be used to create the models:

```
edgeObj <- edgeStudy(data = kidexpr, adj.var = sex,
  tme = age, sampling = "timecourse")
fullMod <- fullModel(edgeObj)
nullMod <- nullModel(edgeObj)
```

The variable **sampling** describes the type of experiment performed, **adj.var** is the adjustment variable and **tme** is the time variable in the study. If the experiment is more complex then type **?edgeStudy** for additional

arguments.

If the alternative and null models are known to the user then `edgeModel` can be used to make an `edgeSet` object:

```
library(splines)
cov <- data.frame(sex = sex, age = age)
null.model <- ~sex
full.model <- ~sex + ns(age, df = 4)
edgeObj <- edgeModel(data = kidexpr, cov = cov, nullMod = null.model,
  altMod = full.model)
```

The `cov` is a data frame of covariates, the `nullMod` is the null model and the `altMod` is the alternative model. The input `cov` is a data frame with the column names the same as the variables in the alternative and null models.

The `odp` or `lrt` function can be used on `edgeObj` to implement either the optimal discovery procedure or the likelihood ratio test, respectively:

```
# optimal discovery procedure
edgeODP <- odp(edgeObj, verbose = FALSE)
# likelihood ratio test
edgeLRT <- lrt(edgeObj)
```

To access the π_0 estimate, p-values, q-values and local false discovery rates for each gene, use the function `qvalueObj`:

```
qvalObj <- qvalueObj(edgeODP)
qvals <- qvalObj$qvalues
pvals <- qvalObj$pvalues
lfdr <- qvalObj$lfdr
pi0 <- qvalObj$pi0
```

The following sections of the manual go through various case studies for a more comprehensive overview of the `edge` package.

5 Case study: static experiment

In the static sampling experiment, the arrays have been collected from distinct biological groups without respect to time. The goal is to identify genes that have a statistically significant difference in average expression across these distinct biological groups. The example data set that will be used in this section is the `gibson` data set and it is a random subset of the data from [Idaghdour et al.](#)

The `gibson` data set provides gene expression measurements in peripheral blood leukocyte samples from three Moroccan Amazigh groups leading distinct ways of life: desert nomadic (DESERT), mountain agrarian (VILLAGE), and coastal urban (AGADIR). We are interested in finding the genes that differentiate the Moroccan Amazigh groups the most. See [Idaghdour et al.](#) for additional information regarding the data.

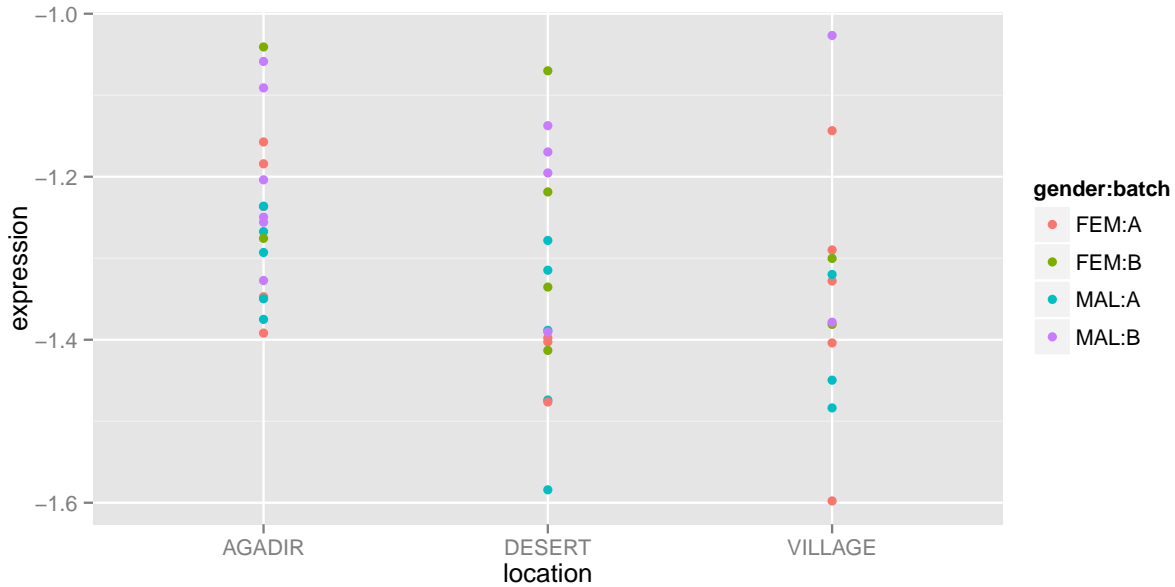


Figure 2: Plot of gene 1 in the `gibson` study.

5.1 Importing the data

To import the `gibson` data use the `data` function:

```
data(gibson)
names(gibson)

## [1] "batch"    "gibexpr"  "gender"   "location"
```

There are a few variables in the data set: `batch`, `gibexpr`, `gender`, and `location`. The three covariates of interest are `gender`, `batch` and `location`. The biological variable is the `location` variable, which contains information on where individuals are sampled: “VILLAGE”, “DESERT” or “AGADIR”. The `gender` variable specifies whether the individual is a male or a female and there are two different `batches` in the study. The `gibexpr` variable contains the gene expression measurements.

As an example, the expression values of the first gene are shown in Figure 2. In the figure, it appears that the individuals from “VILLAGE” are more expressed when compared to other lifestyles. We should stop short of that observation because the data needs to be adjusted with the experimental models. Before that, the alternative and null model of the study needs to be carefully formulated which is discussed in the next section.

5.2 Creating the alternative and null models

In order to find differentially expressed genes, there first needs to be an alternative and null model for the study. There are two ways to input the experimental models in `edge`: `edgeModel` and `edgeStudy`. `edgeStudy` should be used by users unfamiliar with formulating the alternative and null models but are familiar with the covariates in the study:

```
edgeObj <- edgeStudy(data = gibson$gibexpr, adj.var = cbind(gibson$gender,
  gibson$batch), grp = gibson$location, sampling = "static")
```

`adj.var` is for the adjustment variables, `grp` is the variable containing the group assignments for each individual in the study and `sampling` describes the type of experiment. Since `gibson` is a static study, the `sampling` argument will be “static”. The `grp` variable will be the `location` variable and the adjustment variables are `gender` and `batch`.

Alternatively, if the user is familiar with their alternative and null models in the study then `edgeModel` can be used to input the models directly:

```
cov <- data.frame(Gender = gibson$gender, Batch = gibson$batch,
  Location = gibson$location)
null.model <- ~Gender + Batch
alt.model <- ~Gender + Batch + Location
edgeObj <- edgeModel(data = gibson$gibexpr, cov = cov,
  altMod = alt.model, nullMod = null.model)
```

The `cov` argument is a data frame of all the relevant covariates, `altMod` and `nullMod` are the alternative and null models of the experiment, respectively. Notice that the models must be formatted as a formula and contain the same variable names as in the `cov` data frame. The null model contains the `gender` and `batch` covariates and the alternative model includes the `location` variable. Therefore, we are interested in testing whether the alternative model improves the model fit of a gene significantly when compared to the null model. If it does not, then we can conclude that there is no significant difference between Moroccan Amazigh groups for this particular gene.

The variable `edgeObj` is an `edgeSet` object that stores all the relevant experimental data. The `edgeSet` object is discussed further in the next section.

5.3 The edgeSet object

Once either `edgeModel` or `edgeStudy` is used, an `edgeSet` object is created. To view the slots contained in the object:

```
slotNames(edgeObj)

## [1] "null.model"      "full.model"
## [3] "null.matrix"     "full.matrix"
## [5] "individual"      "qvalueObj"
## [7] "experimentData"  "assayData"
## [9] "phenoData"       "featureData"
## [11] "annotation"      "protocolData"
## [13] ".__classVersion__"
```

A description of each slot is listed below:

- `full.model`: the alternative model of the experiment
- `null.model`: the null model of the experiment
- `full.matrix`: the alternative model in matrix form

- `null.matrix`: the null model in matrix form
- `individual`: variable that keeps track of individuals (same individuals are sampled multiple times)
- `qvalueObj`: `qvalue` list. Contains p-values, q-values and local false discovery rates of the significance analysis. See the [qvalue package](#) for more details.
- `ExpressionSet`: inherits the slots from `ExpressionSet` object

`ExpressionSet` contains the expression measurements and other information from the experiment. The `edgeSet` object inherits all the functions from an `ExpressionSet` object. As an example, to access the expression values, one can use the function `exprs` or to access the covariates, `pData`:

```
gibexpr <- exprs(edgeObj)
cov <- pData(edgeObj)
```

The `ExpressionSet` class is a widely used object in Bioconductor and more information can be found [here](#). See the section 11 on `ExpressionSet` to get a better understanding of how it integrates into the `edge` framework.

As an example of how to access the slots of `edgeObj` suppose we are interested in viewing the alternative and null models. The models can be accessed by:

```
fullModel(edgeObj)

## ~Gender + Batch + Location

nullModel(edgeObj)

## ~Gender + Batch
```

Next, we can extract the models in matrix form for computational analysis:

```
full.matrix <- fullMatrix(edgeObj)
null.matrix <- nullMatrix(edgeObj)
```

See `?edgeSet` for additional functions to access different slots of the `edgeSet` object.

5.4 Fitting the data

The `edgeFit` function is an implementation of least squares using the alternative and null models:

```
efObj <- edgeFit(edgeObj, stat.type = "lrt")
```

The `stat.type` argument specifies whether you want the `odp` or `lrt` fitted values. The difference between choosing “odp” and “lrt” is that “odp” centers the data by the null model fit which is necessary for downstream analysis in the optimal discovery procedure. `edgeFit` creates another object with the following slots:

- `fit.full`: fitted values from the alternative model
- `fit.null`: fitted values from null model
- `res.full`: residuals from the alternative model



Figure 3: Plot of gene 1 in the `gibson` study after applying the alternative and null model fit. The “raw” column is the expression values of the original data.

- `res.null`: residuals from the null model
- `dH.full`: diagonal elements in the projection matrix for the full model
- `beta.coef`: the coefficients for the full model
- `stat.type`: statistic type used, either “odp” or “lrt”

To access the fitted coefficients of the alternative model in `efObj`:

```
betaCoef(efObj)
```

To access the alternative and null residuals:

```
alt.res <- resFull(efObj)
null.res <- resNull(efObj)
```

To access the fitted values:

```
alt.fitted <- fitFull(efObj)
null.fitted <- fitNull(efObj)
```

See `?edgeFit` for more details on accessing the slots in an `edgeFit` object. The fitted values of the first gene is shown in Figure 3. The null model fit is the average expression value across the interaction of `batch` and `sex`. The alternative model fit seems to pick up some differences relative to the null model. Next, we have to test whether the observed differences between the model fits are significant.

5.5 Significance analysis

Interpreting the models in a hypothesis test is very intuitive: Does the alternative model better fit the data when compared to the null model? For the fitted values of the first gene plotted in Figure 3, it seems that the alternative model fits the data better than the null model. In order to conclude that it is significant, we need to calculate the p-value. The user can use either the optimal discovery procedure or likelihood ratio test.

5.5.1 Likelihood ratio test

The `lrt` function performs a likelihood ratio test to determine p-values:

```
edgeLRT <- lrt(edgeObj, nullDistn = "normal")
```

If the null distribution, `nullDistn`, is calculated using “bootstrap” then residuals from the alternative model are re-sampled and added to the null model to simulate a distribution where there is no differential expression. Otherwise, the default input is “normal” and the assumption is that the null statistics follow a F-distribution. See `?lrt` for additional arguments.

5.5.2 Optimal discovery procedure

`odp` performs the optimal discovery procedure, which is a new approach developed by Storey et al. (2005) for optimally performing many hypothesis tests in a high-dimensional study. When testing a feature, information from all the features is utilized when testing for significance of a feature. It guarantees to maximize the number of expected true positive results for each fixed number of expected false positive results which is related to the false discovery rate. The optimal discovery procedure can be implemented on `edgeObj` by the `odp` function:

```
edgeODP <- odp(edgeObj, bs.its = 30, verbose = FALSE,
  n.mods = 50)
```

The number of bootstrap iterations is controlled by `bs.its`, `verbose` prints each bootstrap iteration number and `n.mods` is the number of clusters in the k-means algorithm. A k-means algorithm is used to assign genes to groups in order to speed up the computational time of the algorithm. If `n.mods` is equal to the number of genes then the original optimal discovery procedure is used. Depending on the number of genes, this setting can take a very long time. Therefore, it is recommended to use a small `n.mods` value to substantially decrease the computational time. In Woo et al. (2011), it is shown that assigning `n.mods` to about 50 will cause a negligible loss in power. Type `?odp` for more details on the algorithm.

5.6 Significance results

The `summary` function can be used on an `edgeSet` object to give an overview of the analysis:

```
summary(edgeODP)
##
## ExpressionSet Summary
##
```

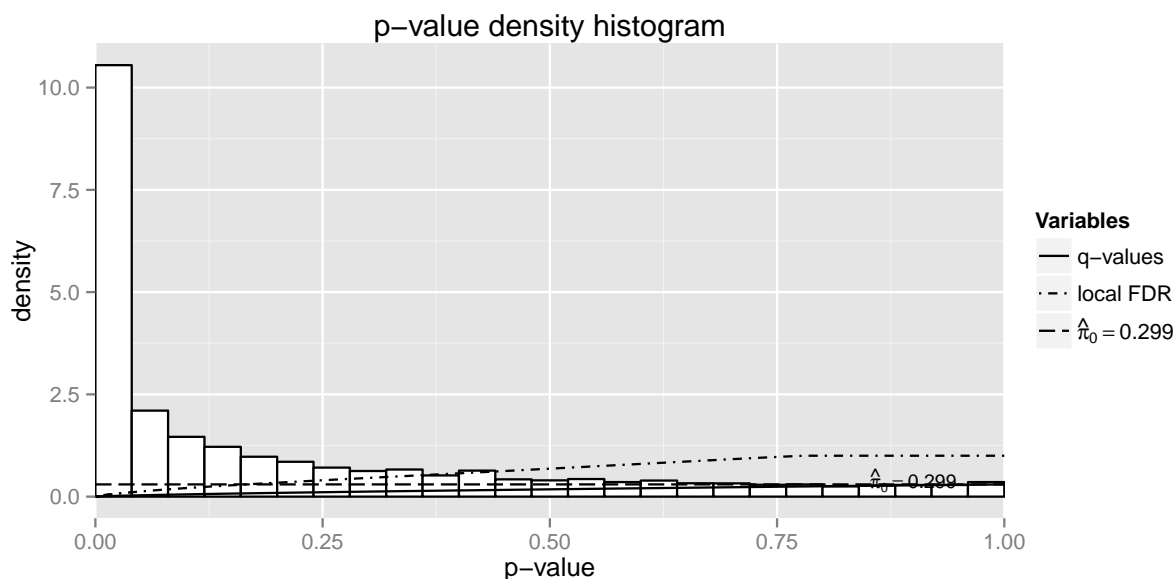


Figure 4: Applying the function `hist` to the slot `qvalueObj` in the `gibson` data set. Function is derived from the `qvalue` package.

```
## ExpressionSet (storageMode: lockedEnvironment)
## assayData: 10177 features, 46 samples
##   element names: exprs
## protocolData: none
## phenoData
##   sampleNames: 1 2 ... 46 (46 total)
##   varLabels: Gender Batch Location
##   varMetadata: labelDescription
## featureData: none
## experimentData: use 'experimentData(object)'
## Annotation:
##
## edge Analysis Summary
##
## Total number of arrays: 46
## Total number of probes: 10177
##
## Biological variables:
##   Null Model: ~Gender + Batch
##
##   Full Model: ~Gender + Batch + Location
##
## .....
##
## Statistical significance summary:
## pi0: 0.2991622
```

```
##
## Cumulative number of significant calls:
##
##          <1e-04 <0.001 <0.01 <0.025 <0.05 <0.1
## p-value      932   1679   2934   3783   4571  5487
## q-value      744   1362   2903   4117   5251  6863
## local fdr    492    847   1869   2500   3185  4194
##
##          <1
## p-value    10177
## q-value    10177
## local fdr   9530
```

There are three core summaries: **ExpressionSet** summary, **edge** analysis and statistical significance summary. The **ExpressionSet** summary shows a summary of the **ExpressionSet** object. **edge** analysis shows an overview of the models used and other information about the data set. The significance analysis shows the proportion of null genes, π_0 , and significant genes at various cutoffs in terms of p-values, q-values and local false discovery rates.

The function `qvalueObj` can be used on `edgeODP` to extract the significance results:

```
sig.results <- qvalueObj(edgeODP)
```

The object `sig.results` is a list with the following slots:

```
names(sig.results)

## [1] "call"      "pi0"       "qvalues"
## [4] "pvalues"   "lfdr"      "pi0.lambda"
## [7] "lambda"    "pi0.smooth"
```

The key variables are `pi0`, `pvalues`, `lfdr` and `qvalues`. The `pi0` variable provides an estimate of the proportion of null p-values, `pvalues` are the p-values, `qvalues` are the estimated q-values and `lfdr` are the local false discovery rates. Using the function `hist` on `sig.results` will produce a p-value histogram along with the density curves of q-values and local false discovery rate values:

```
hist(sig.results)
```

The plot is shown in Figure 4. To extract the p-values, q-values, local false discovery rates and the π_0 estimate:

```
pvalues <- sig.results$pvalues
qvalues <- sig.results$qvalues
lfdr <- sig.results$lfdr
pi0 <- sig.results$pi0
```

Making significance decisions based on p-values in multiple hypothesis testings problems can lead to accepting a lot of false positives in the study. Instead, using q-values to determine significant genes is recommended because it controls the false discovery rate at a level `alpha`. Q-values measure the proportion of false positives incurred when calling a particular test significant. For example, to complete our analysis of gene 1 in this example, let's view the q-value estimate:

```
qvalues[1]

## [1] 3.659477e-05
```

So for this particular gene, the q-value is 3.6594766×10^{-5} . If we consider a false discovery rate cutoff of 0.1 then this gene is significant. Therefore, the observed differences observed in Figure 3 are significant so this particular gene is differentially expressed between locations.

To get a list of all the significant genes at a false discovery rate cutoff of 0.01:

```
fdr.level <- 0.01
sigGenes <- qvalues < fdr.level
```

View the [qvalue vignette](#) to get a more thorough discussion in how to use p-values, q-values, π_0 estimate and local false discovery rates to determine significant genes.

6 Case study: independent time course experiment

In the independent time course study, the arrays have been sampled with respect to time from one biological group and the goal is to identify genes that show “within-class temporal differential expression”, i.e., genes that show statistically significant changes in expression over time. The example data set used in this section is a kidney data set by [Rodwell et al. \(2004\)](#). Gene expression measurements, from cortex and medulla samples in the kidney, were obtained from 72 human subjects ranging in age from 27 to 92 years. Only one array was obtained per sample and the age and tissue type of each subject was recorded. See [Rodwell et al. \(2004\)](#) for additional information regarding the data set.

6.1 Importing the data

To import the `kidney` data use the `data` function:

```
data(kidney)
names(kidney)

## [1] "age"      "sex"      "kidexpr"
```

There are a few covariates in this data set: `sex`, `age`, `tissue`, `kidexpr` and `kidcov`. In this example, we will focus on the cortex `tissue` samples:

```
sex <- kidney$sex
age <- kidney$age
kidexpr <- kidney$kidexpr
```

The two main covariates of interest for this example are the `sex` and `age` covariates. The `sex` variable is whether the subject was male or female and the `age` variable is the age of the patients. `kidexpr` contains the gene expression values for the study.

As an example of a gene in the study, the expression values of the fifth gene are shown in Figure 5. It is very difficult to find a trend for this particular gene. Instead, we need to adjust the data with the models in the study which is discussed in the next section.



Figure 5: Plot of gene 5 in the kidney study.

6.2 Creating the alternative and null models

In order to find differentially expressed genes, the alternative and null model for the study need to be formulated. There are two ways to input the experimental models in `edge`: `edgeModel` and `edgeStudy`. `edgeStudy` should be used by users unfamiliar with formulating the alternative and null models but are familiar with the covariates in the study:

```
edgeObj <- edgeStudy(data = kidexpr, adj.var = sex,
  tme = age, sampling = "timecourse", basis.df = 4)
```

`adj.var` is for the adjustment variables, `tme` is the time variable, `basis.df` is the degrees of freedom for the spline fit, and `sampling` describes the type of experiment. Since `kidney` is a time course study, the `sampling` argument will be “timecourse”. The `tme` variable will be the `age` variable, `basis.df` will be 4 based on previous work by [Storey et al. \(2005\)](#) and the adjustment variable is `sex`. To view the models generated by `edgeStudy`:

```
fullModel(edgeObj)

## ~adj.var + ns(tme, df = 4, intercept = FALSE)
## <environment: 0xd6bc1b8>

nullModel(edgeObj)

## ~adj.var
## <environment: 0xd6bc1b8>
```

Notice that the difference between the alternative and null model is the natural spline fit of the `age` variable. If we look at [Figure 5](#), it becomes evident that a spline curve can be used to approximate the fit of the data, and 4 degrees of freedom is chosen based on previous analysis of the expression patterns. See [Storey et al.](#)

(2005) for a detailed discussion on modelling in time course studies.

Alternatively, if the user is familiar with their alternative and null models in the study then `edgeModel` can be used to input the models directly:

```
library(splines)
cov <- data.frame(sex = sex, age = age)
null.model <- ~sex
alt.model <- ~sex + ns(age, df = 4)
edgeObj <- edgeModel(data = kidexpr, cov = cov, altMod = alt.model,
  nullMod = null.model)
```

The `cov` argument is a data frame of all the relevant covariates, `altMod` and `nullMod` are the alternative and null models of the experiment, respectively. Notice that the models must be formatted as a formula and contain the same variable names as in the `cov` data frame. The null model contains the `sex` covariate and the alternative model includes the `age` variable. Therefore, we are interested in testing whether the alternative model improves the model fit of a gene significantly when compared to the null model. If it does not, then we can conclude that there is no significant difference in the gene as it ages in the cortex.

The variable `edgeObj` is an `edgeSet` object that stores all the relevant experimental data. The `edgeSet` object is discussed further in the next section.

6.3 The edgeSet object

Once either `edgeModel` or `edgeStudy` is used, an `edgeSet` object is created. To view the slots contained in the object:

```
slotNames(edgeObj)

## [1] "null.model"      "full.model"
## [3] "null.matrix"     "full.matrix"
## [5] "individual"      "qvalueObj"
## [7] "experimentData"  "assayData"
## [9] "phenoData"       "featureData"
## [11] "annotation"      "protocolData"
## [13] ".__classVersion__"
```

A description of each slot is listed below:

- `full.model`: the alternative model of the experiment
- `null.model`: the null model of the experiment
- `full.matrix`: the alternative model in matrix form
- `null.matrix`: the null model in matrix form
- `individual`: variable that keeps track of individuals (same individuals are sampled multiple times)
- `qvalueObj`: `qvalue` list. Contains p-values, q-values and local false discovery rates of the significance analysis. See the [qvalue package](#) for more details.
- `ExpressionSet`: inherits the slots from `ExpressionSet` object

`ExpressionSet` contains the expression measurements and other information from the experiment. The `edgeSet` object inherits all the functions from an `ExpressionSet` object. As an example, to access the expression values, one can use the function `exprs` or to access the covariates, `pData`:

```
gibexpr <- exprs(edgeObj)
cov <- pData(edgeObj)
```

The `ExpressionSet` class is a widely used object in Bioconductor and more information can be found [here](#). See the section 11 on `ExpressionSet` to get a better understanding of how it integrates into the `edge` framework.

As an example of how to access the slots of `edgeObj` suppose we are interested in viewing the alternative and null models. The models can be accessed by:

```
fullModel(edgeObj)

## ~sex + ns(age, df = 4)

nullModel(edgeObj)

## ~sex
```

Next, we can extract the models in matrix form for computational analysis:

```
full.matrix <- fullMatrix(edgeObj)
null.matrix <- nullMatrix(edgeObj)
```

See `?edgeSet` for additional functions to access different slots of the `edgeSet` object.

6.4 Fitting the data

The `edgeFit` function is an implementation of least squares using the alternative and null models:

```
efObj <- edgeFit(edgeObj, stat.type = "lrt")
```

The `stat.type` argument specifies whether you want the `odp` or `lrt` fitted values. The difference between choosing “`odp`” and “`lrt`” is that “`odp`” centers the data by the null model fit which is necessary for downstream analysis in the optimal discovery procedure. `edgeFit` creates another object with the following slots:

- `fit.full`: fitted values from the alternative model
- `fit.null`: fitted values from null model
- `res.full`: residuals from the alternative model
- `res.null`: residuals from the null model
- `dH.full`: diagonal elements in the projection matrix for the full model
- `beta.coef`: the coefficients for the full model
- `stat.type`: statistic type used, either “`odp`” or “`lrt`”

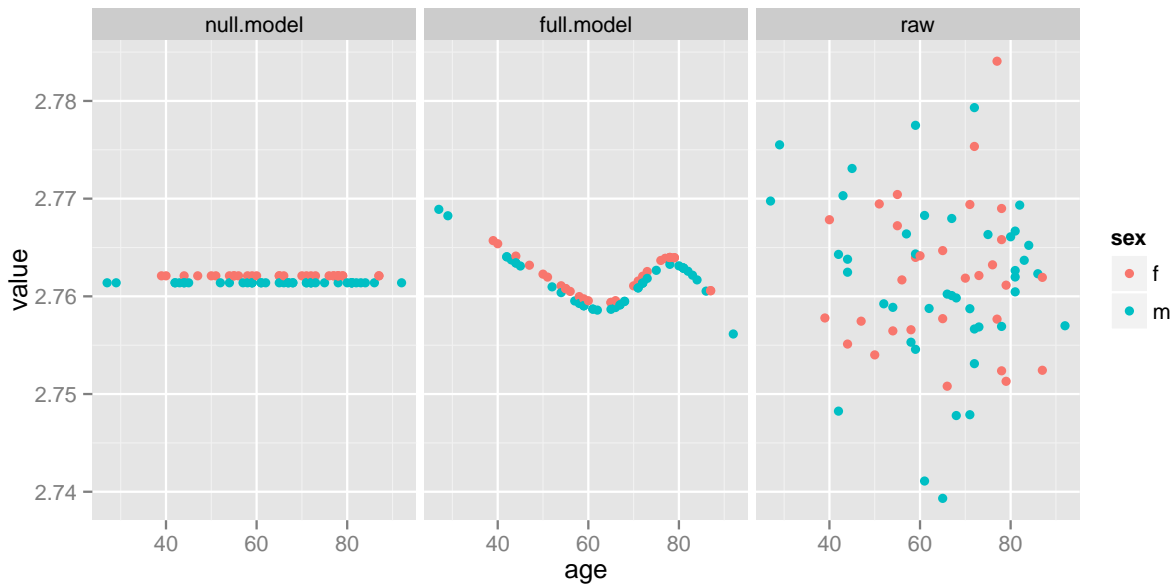


Figure 6: Plot of gene 5 in the `kidney` study after applying the alternative and null model fit. The “raw” column is the expression values of the original data.

To access the fitted coefficients of the alternative model in `efObj`:

```
betaCoef(efObj)
```

To access the alternative and null residuals:

```
alt.res <- resFull(efObj)
null.res <- resNull(efObj)
```

To access the fitted values:

```
alt.fitted <- fitFull(efObj)
null.fitted <- fitNull(efObj)
```

See `?edgeFit` for more details on accessing the slots in an `edgeFit` object. The fitted values of the fifth gene is shown in Figure 6. The null model fit is the average expression. It appears that the alternative model fits the raw data better than the null model. Next, we have to test whether the observed differences between the model fits are significant.

6.5 Significance analysis

Interpreting the models in a hypothesis test is very intuitive: Does the alternative model better fit the data when compared to the null model? For the fitted values of the fifth gene plotted in Figure 6, it seems that the alternative model fits the data better than the null model. In order to conclude it is significant, we need to calculate the p-value. The user can use either the optimal discovery procedure or likelihood ratio test.

6.5.1 Likelihood ratio test

The `lrt` function performs a likelihood ratio test to determine p-values:

```
edgeLRT <- lrt(edgeObj, nullDistn = "normal")
```

If the null distribution, `nullDistn`, is calculated using “bootstrap” then residuals from the alternative model are re-sampled and added to the null model to simulate a distribution where there is no differential expression. Otherwise, the default input is “normal” and the assumption is that the null statistics follow a F-distribution. See `?lrt` for additional arguments.

6.5.2 Optimal discovery procedure

`odp` performs the optimal discovery procedure, which is a new approach developed by [Storey et al. \(2005\)](#) for optimally performing many hypothesis tests in a high-dimensional study. When testing a feature, information from all the features is utilized when testing for significance of a feature. It guarantees to maximize the number of expected true positive results for each fixed number of expected false positive results which is related to the false discovery rate. The optimal discovery procedure can be implemented on `edgeObj` by the `odp` function:

```
edgeODP <- odp(edgeObj, bs.its = 30, verbose = FALSE,
  n.mods = 50)
```

The number of bootstrap iterations is controlled by `bs.its`, `verbose` prints each bootstrap iteration number and `n.mods` is the number of clusters in the k-means algorithm. A k-means algorithm is used to assign genes to groups in order to speed up the computational time of the algorithm. If `n.mods` is equal to the number of genes then the original optimal discovery procedure is used. Depending on the number of genes, this setting can take a very long time. Therefore, it is recommended to use a small `n.mods` value to substantially decrease the computational time. In [Woo et al. \(2011\)](#), it is shown that assigning `n.mods` to about 50 will cause a negligible loss in power. Type `?odp` for more details on the algorithm.

6.6 Significance results

The `summary` function can be used on an `edgeSet` object to give an overview of the analysis:

```
summary(edgeODP)

##
## ExpressionSet Summary
##
## ExpressionSet (storageMode: lockedEnvironment)
## assayData: 5000 features, 72 samples
##   element names: exprs
## protocolData: none
## phenoData
##   sampleNames: 1 2 ... 72 (72 total)
##   varLabels: sex age
##   varMetadata: labelDescription
## featureData: none
## experimentData: use 'experimentData(object)'
```

```
## Annotation:
##
## edge Analysis Summary
##
## Total number of arrays: 72
## Total number of probes: 5000
##
## Biological variables:
## Null Model:~sex
##
## Full Model:~sex + ns(age, df = 4)
##
## .....
##
## Statistical significance summary:
## pi0: 0.4202189
##
## Cumulative number of significant calls:
##
##          <1e-04 <0.001 <0.01 <0.025 <0.05 <0.1
## p-value      33      134      370      605      888 1311
## q-value       0       16       80      184      327 720
## local fdr     0        0       36      113      174 328
##
##          <1
## p-value    5000
## q-value    5000
## local fdr 4967
```

There are three core summaries: **ExpressionSet** summary, **edge** analysis and statistical significance summary. The **ExpressionSet** summary shows a summary of the **ExpressionSet** object. **edge** analysis shows an overview of the models used and other information about the data set. The significance analysis shows the proportion of null genes, π_0 , and significant genes at various cutoffs in terms of p-values, q-values and local false discovery rates.

The function `qvalueObj` can be used on `edgeODP` to extract the significance results:

```
sig.results <- qvalueObj(edgeODP)
```

The object `sig.results` is a list with the following slots:

```
names(sig.results)

## [1] "call"      "pi0"       "qvalues"
## [4] "pvalues"   "lfdr"      "pi0.lambda"
## [7] "lambda"    "pi0.smooth"
```

The key variables are `pi0`, `pvalues`, `lfdr` and `qvalues`. The `pi0` variable provides an estimate of the proportion of null p-values, `pvalues` are the p-values, `qvalues` are the estimated q-values and `lfdr` are the local false discovery rates. Using the function `hist` on `sig.results` will produce a p-value histogram along with the density curves of q-values and local false discovery rate values:

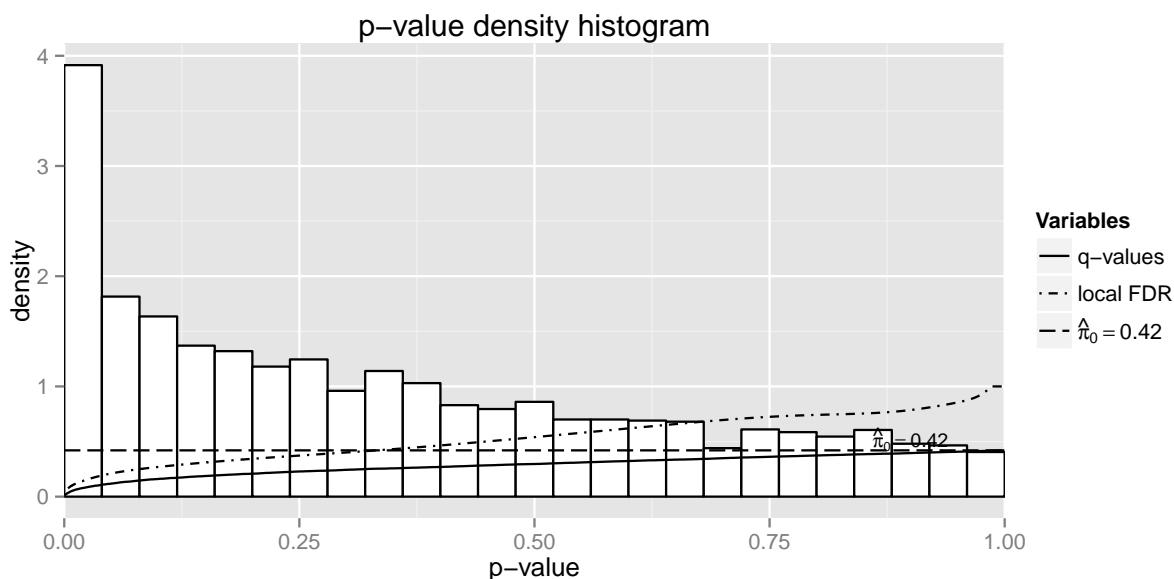


Figure 7: Applying the function `hist` to the slot `qvalueObj` in the `kidney` data set. Function is derived from the `qvalue` package.

```
hist(sig.results)
```

The plot is shown in Figure 7. To extract the p-values, q-values, local false discovery rates and the π_0 estimate:

```
pvalues <- sig.results$pvalues
qvalues <- sig.results$qvalues
lfdr <- sig.results$lfdr
pi0 <- sig.results$pi0
```

Making significance decisions based on p-values in multiple hypothesis testings problems can lead to accepting a lot of false positives in the study. Instead, using q-values to determine significant genes is recommended because it controls the false discovery rate at a level `alpha`. Q-values measure the proportion of false positives incurred when calling a particular test significant. For example, to complete our analysis of gene 5 in this example, let's view the q-value estimate:

```
qvalues[5]
## [1] 0.1586294
```

So for this particular gene, the q-value is 0.1586294. If we consider a false discovery rate cutoff of 0.1 then this gene is not significant. Therefore, the observed differences observed in Figure 6 are not significant so this particular gene is not differentially expressed as the kidney ages.

To get a list of all the significant genes at a false discovery rate cutoff of 0.1:

```
fdr.level <- 0.1
sigGenes <- qvalues < fdr.level
```

View the [qvalue vignette](#) to get a more thorough discussion in how to use p-values, q-values, π_0 estimate and local false discovery rates to determine significant genes.

7 Case study: longitudinal time course experiment

In the longitudinal time course study, the goal is to identify genes that show “between-class temporal differential expression”, i.e., genes that show statistically significant differences in expression over time between the various groups. The `endotoxin` data set provides gene expression measurements in an endotoxin study where four subjects were given endotoxin and four subjects were given a placebo. Blood samples were collected and leukocytes were isolated from the samples before infusion. Measurements were recorded at times 2, 4, 6, 9, 24 hours. We are interested in identifying genes that vary over time between the endotoxin and control groups. See [Calvano et al. \(2005\)](#) for more details regarding the `endotoxin` dataset.

7.1 Importing the data

To import the `endotoxin` data use the `data` function:

```
data(endotoxin)
names(endotoxin)

## [1] "class"      "endoexpr"  "ind"       "time"
```

There are a few covariates in this data set: `expr`, `class`, `individual`, and `time`. There are 8 individuals in the experiment (`ind`) that were sampled at multiple time points (`time`) that were either “endotoxin” or “control” (`class`). The `expr` variable contains the expression values of the experiment:

To show an example gene, the expression values of the second gene are shown in Figure 8. It is very difficult to find a trend for this particular gene. Instead, we need to adjust the data with the models in the study.

7.2 Creating the alternative and null models

In order to find differentially expressed genes, there first needs to be an alternative and null model for the study. There are two ways to input the experimental models in `edge`: `edgeModel` and `edgeStudy`. `edgeStudy` should be used by users unfamiliar with formulating the alternative and null models but are familiar with the covariates in the study:

```
edgeObj <- edgeStudy(data = endotoxin$endoexpr, grp = endotoxin$class,
  tme = endotoxin$time, ind = endotoxin$ind, sampling = "timecourse")
```

`grp` is for the variable which group each individual belongs to, `tme` is the time variable, `ind` is used when individuals are sampling multiple times and `sampling` describes the type of experiment. Since `endotoxin` is a time course study, the `sampling` argument will be “timecourse”. The `tme` variable will be the time variable, `ind` is the individuals variable and the `grp` variable is `class`. To view the models created by `edgeStudy`:

```
fullModel(edgeObj)
```

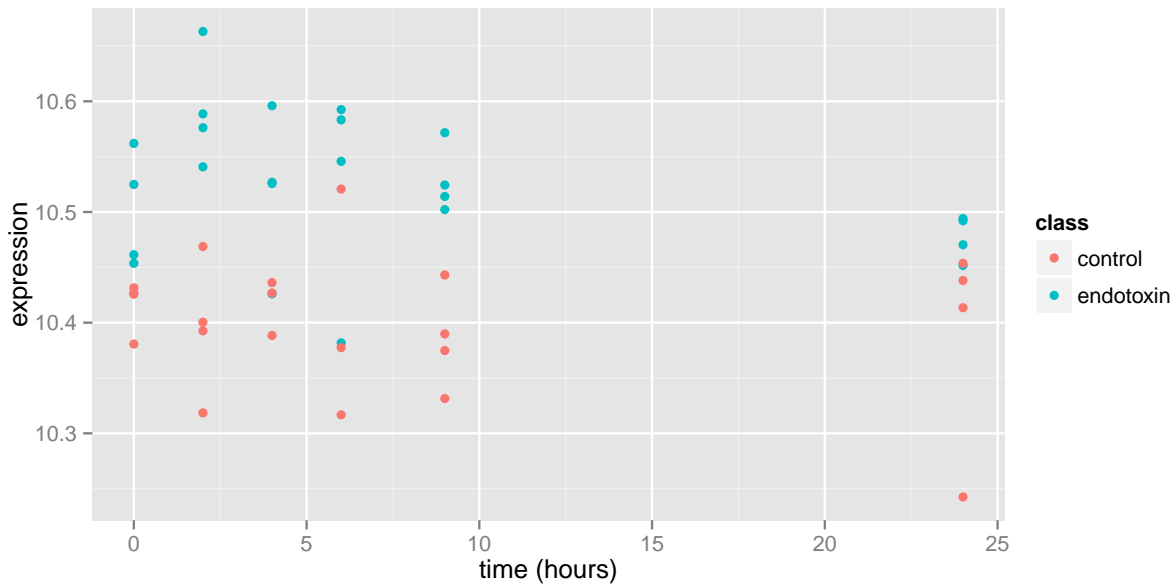


Figure 8: Plot of gene 2 in the **endotoxin** study.

```
## ~grp + ns(tme, df = 2, intercept = FALSE) + (grp):ns(tme, df = 2,
##       intercept = FALSE)
## <environment: 0xdb9ce28>

nullModel(edgeObj)

## ~grp + ns(tme, df = 2, intercept = FALSE)
## <environment: 0xdb9ce28>
```

See [Storey et al. \(2005\)](#) for how the models in the **endotoxin** experiment are formed. Alternatively, if the user is familiar with their alternative and null models in the study then `edgeModel` can be used to input the models directly:

```
cov <- data.frame(ind = endotoxin$ind, tme = endotoxin$time,
  grp = endotoxin$class)
null.model <- ~grp + ns(tme, df = 2, intercept = FALSE)
alt.model <- ~grp + ns(tme, df = 2, intercept = FALSE) +
  (grp):ns(tme, df = 2, intercept = FALSE)
edgeObj <- edgeModel(data = endotoxin$endoexpr, cov = cov,
  altMod = alt.model, nullMod = null.model)
```

The `cov` argument is a data frame of all the relevant covariates, `altMod` and `nullMod` are the alternative and null models of the experiment, respectively. Notice that the models must be formatted as a formula and contain the same variable names as in the `cov` data frame. We are interested in testing whether the alternative model improves the model fit of a gene significantly when compared to the null model. If it does not, then we can conclude that there is no significant difference in this gene between the endotoxin and the control as time goes on.

The variable `edgeObj` is an `edgeSet` object that stores all the relevant experimental data. The `edgeSet`

object is discussed further in the next section.

7.3 The edgeSet object

Once either `edgeModel` or `edgeStudy` is used, an `edgeSet` object is created. To view the slots contained in the object:

```
slotNames(edgeObj)

## [1] "null.model"      "full.model"
## [3] "null.matrix"     "full.matrix"
## [5] "individual"      "qvalueObj"
## [7] "experimentData"  "assayData"
## [9] "phenoData"       "featureData"
## [11] "annotation"      "protocolData"
## [13] ".__classVersion__"
```

A description of each slot is listed below:

- `full.model`: the alternative model of the experiment
- `null.model`: the null model of the experiment
- `full.matrix`: the alternative model in matrix form
- `null.matrix`: the null model in matrix form
- `individual`: variable that keeps track of individuals (same individuals are sampled multiple times)
- `qvalueObj`: `qvalue` list. Contains p-values, q-values and local false discovery rates of the significance analysis. See the [qvalue package](#) for more details.
- `ExpressionSet`: inherits the slots from `ExpressionSet` object

`ExpressionSet` contains the expression measurements and other information from the experiment. The `edgeSet` object inherits all the functions from an `ExpressionSet` object. As an example, to access the expression values, one can use the function `exprs` or to access the covariates, `pData`:

```
gibexpr <- exprs(edgeObj)
cov <- pData(edgeObj)
```

The `ExpressionSet` class is a widely used object in Bioconductor and more information can be found [here](#). See the section 11 on `ExpressionSet` to get a better understanding of how it integrates into the `edge` framework.

As an example of how to access the slots of `edgeObj` suppose we are interested in viewing the alternative and null models. The models can be accessed by:

```
fullModel(edgeObj)

## ~grp + ns(tme, df = 2, intercept = FALSE) + (grp):ns(tme, df = 2,
##      intercept = FALSE)
```

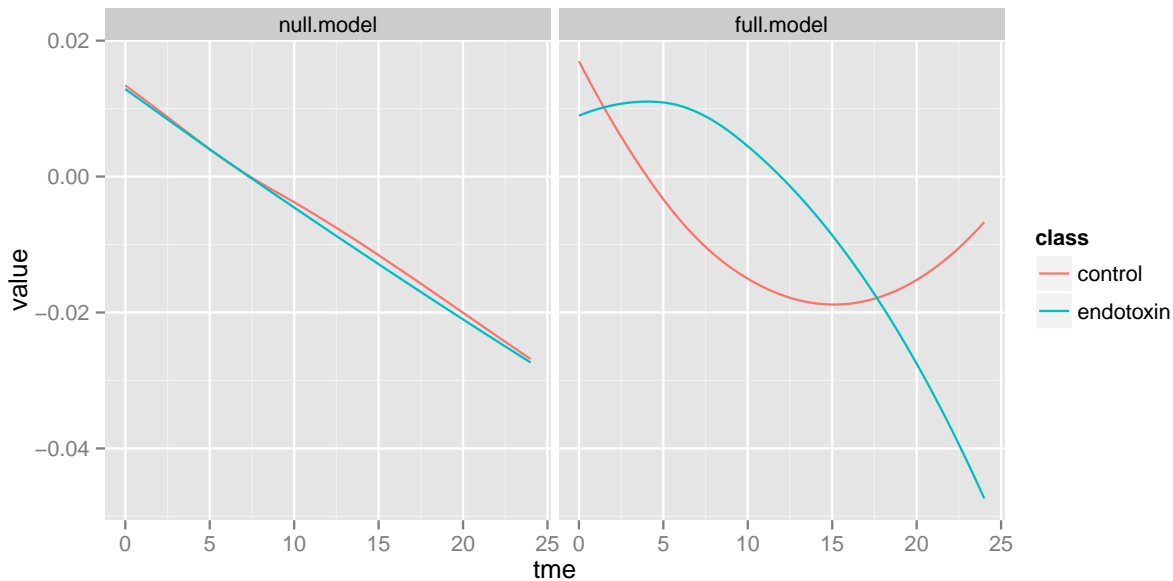



Figure 9: Plot of gene 2 in the **endotoxin** study after applying the alternative and null model fit. The “raw” column is the expression values of the original data.

```
nullModel(edgeObj)

## ~grp + ns(tme, df = 2, intercept = FALSE)
```

Next, we can extract the models in matrix form for computational analysis:

```
full.matrix <- fullMatrix(edgeObj)
null.matrix <- nullMatrix(edgeObj)
```

See `?edgeSet` for additional functions to access different slots of the `edgeSet` object.

7.4 Fitting the data

The `edgeFit` function is an implementation of least squares using the alternative and null models:

```
efObj <- edgeFit(edgeObj, stat.type = "lrt")
```

The `stat.type` argument specifies whether you want the `odp` or `lrt` fitted values. The difference between choosing “`odp`” and “`lrt`” is that “`odp`” centers the data by the null model fit which is necessary for downstream analysis in the optimal discovery procedure. `edgeFit` creates another object with the following slots:

- `fit.full`: fitted values from the alternative model
- `fit.null`: fitted values from null model
- `res.full`: residuals from the alternative model

- `res.null`: residuals from the null model
- `dH.full`: diagonal elements in the projection matrix for the full model
- `beta.coef`: the coefficients for the full model
- `stat.type`: statistic type used, either “odp” or “lrt”

To access the fitted coefficients of the alternative model in `efObj`:

```
betaCoef(efObj)
```

To access the alternative and null residuals:

```
alt.res <- resFull(efObj)
null.res <- resNull(efObj)
```

To access the fitted values:

```
alt.fitted <- fitFull(efObj)
null.fitted <- fitNull(efObj)
```

See `?edgeFit` for more details on accessing the slots in an `edgeFit` object. The fitted values of the second gene is shown in Figure 9. The null model fit is the average expression. It appears that the alternative model fits a pattern that might be observed in the raw data. Next, we have to test whether the observed differences between the model fits are significant.

7.5 Significance analysis

Interpreting the models in a hypothesis test is very intuitive: Does the alternative model better fit the data when compared to the null model? For the fitted values of the second gene plotted in Figure 9, it seems that the alternative model fits the data better than the null model. In order to conclude it is significant, we need to calculate the p-value. The user can use either the optimal discovery procedure or likelihood ratio test.

7.5.1 Likelihood ratio test

The `lrt` function performs a likelihood ratio test to determine p-values:

```
edgeLRT <- lrt(edgeObj, nullDistn = "normal")
```

If the null distribution, `nullDistn`, is calculated using “bootstrap” then residuals from the alternative model are re-sampled and added to the null model to simulate a distribution where there is no differential expression. Otherwise, the default input is “normal” and the assumption is that the null statistics follow a F-distribution. See `?lrt` for additional arguments.

7.5.2 Optimal discovery procedure

odp performs the optimal discovery procedure, which is a new approach developed by Storey et al. (2005) for optimally performing many hypothesis tests in a high-dimensional study. When testing a feature, information

from all the features is utilized when testing for significance of a feature. It guarantees to maximize the number of expected true positive results for each fixed number of expected false positive results which is related to the false discovery rate. The optimal discovery procedure can be implemented on `edgeObj` by the `odp` function:

```
edgeODP <- odp(edgeObj, bs.its = 30, verbose = FALSE,
  n.mods = 50)
```

The number of bootstrap iterations is controlled by `bs.its`, `verbose` prints each bootstrap iteration number and `n.mods` is the number of clusters in the k-means algorithm. A k-means algorithm is used to assign genes to groups in order to speed up the computational time of the algorithm. If `n.mods` is equal to the number of genes then the original optimal discovery procedure is used. Depending on the number of genes, this setting can take a very long time. Therefore, it is recommended to use a small `n.mods` value to substantially decrease the computational time. In [Woo et al. \(2011\)](#), it is shown that assigning `n.mods` to about 50 will cause a negligible loss in power. Type `?odp` for more details on the algorithm.

7.6 Significance results

The `summary` function can be used on an `edgeSet` object to give an overview of the analysis:

```
summary(edgeODP)

##
## ExpressionSet Summary
##
## ExpressionSet (storageMode: lockedEnvironment)
## assayData: 5000 features, 46 samples
##   element names: exprs
## protocolData: none
## phenoData
##   sampleNames: 1 2 ... 46 (46 total)
##   varLabels: ind tme grp
##   varMetadata: labelDescription
## featureData: none
## experimentData: use 'experimentData(object)'
## Annotation:
##
## edge Analysis Summary
##
## Total number of arrays: 46
## Total number of probes: 5000
##
## Biological variables:
##   Null Model:~grp + ns(tme, df = 2, intercept = FALSE)
##
##   Full Model:~grp + ns(tme, df = 2, intercept = FALSE) + (grp):ns(tme, df = 2,
##     intercept = FALSE)
##
## .....
##
##
```

```
## Statistical significance summary:
## pi0: 0.6843337
##
## Cumulative number of significant calls:
##
##          <1e-04 <0.001 <0.01 <0.025 <0.05 <0.1
## p-value      107     190    380     549    750 1056
## q-value       0      82    142     213    277 406
## local fdr     0      70     94     121    169 221
##
##          <1
## p-value     5000
## q-value     5000
## local fdr   5000
```

There are three core summaries: **ExpressionSet** summary, **edge** analysis and statistical significance summary. The **ExpressionSet** summary shows a summary of the **ExpressionSet** object. **edge** analysis shows an overview of the models used and other information about the data set. The significance analysis shows the proportion of null genes, π_0 , and significant genes at various cutoffs in terms of p-values, q-values and local false discovery rates.

The function `qvalueObj` can be used on `edgeODP` to extract the significance results:

```
sig.results <- qvalueObj(edgeODP)
```

The object `sig.results` is a list with the following slots:

```
names(sig.results)

## [1] "call"      "pi0"       "qvalues"
## [4] "pvalues"   "lfdr"      "pi0.lambda"
## [7] "lambda"    "pi0.smooth"
```

The key variables are `pi0`, `pvalues`, `lfdr` and `qvalues`. The `pi0` variable provides an estimate of the proportion of null p-values, `pvalues` are the p-values, `qvalues` are the estimated q-values and `lfdr` are the local false discovery rates. Using the function `hist` on `sig.results` will produce a p-value histogram along with the density curves of q-values and local false discovery rate values:

```
hist(sig.results)
```

The plot is shown in Figure 10. To extract the p-values, q-values, local false discovery rates and the π_0 estimate:

```
pvalues <- sig.results$pvalues
qvalues <- sig.results$qvalues
lfdr <- sig.results$lfdr
pi0 <- sig.results$pi0
```

Making significance decisions based on p-values in multiple hypothesis testings problems can lead to accepting a lot of false positives in the study. Instead, using q-values to determine significant genes is recommended because it controls the false discovery rate at a level `alpha`. Q-values measure the proportion of false positives incurred when calling a particular test significant. For example, to complete our analysis of gene 2 in this example, let's view the q-value estimate:

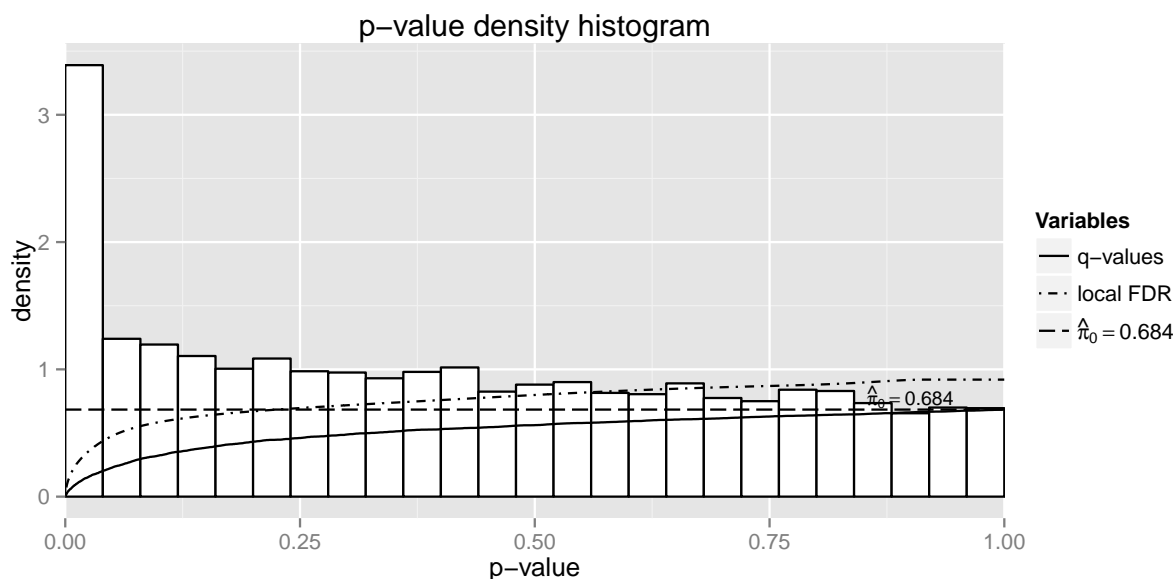


Figure 10: Applying the function `hist` to the slot `qvalueObj` in the `endotoxin` data set. Function is derived from the `qvalue` package.

```
qvalues[2]
## [1] 0.5336279
```

So for this particular gene, the q-value is 0.5336279. If we consider a false discovery rate cutoff of 0.1 then this gene is not significant. Therefore, the observed differences observed in Figure 9 are not significant so this particular gene is not differentially expressed between `class` as time varies.

To get a list of all the significant genes at a false discovery rate cutoff of 0.1:

```
fdr.level <- 0.1
sigGenes <- qvalues < fdr.level
```

View the [qvalue vignette](#) to get a more thorough discussion in how to use p-values, q-values, π_0 estimate and local false discovery rates to determine significant genes.

8 sva: Surrogate variable analysis

The `sva` package is useful for removing batch effects or any unwanted variation in an experiment. It does this by forming surrogate variables to adjust for sources of unknown variation. Details on the algorithm can be found in [Leek and Storey \(2007\)](#). `edge` uses the `sva` package in the function `edgeSVA`. Suppose we are working with the `kidney` data in 6, then the first step is to create an `edgeSet` object by either using `edgeModel` or `edgeStudy`:

```
library(splines)
cov <- data.frame(sex = sex, age = age)
```

```

null.model <- ~sex
alt.model <- ~sex + ns(age, df = 4)
edgeObj <- edgeModel(data = kidexpr, cov = cov, altMod = alt.model,
  nullMod = null.model)

```

To find the surrogate variables and add them to the experimental models in `edgeObj`, use the function `edgeSVA`:

```

edgeObj <- edgeSVA(edgeObj, n.sv = 3, B = 10)

## Number of significant surrogate variables is: 3
## Iteration (out of 10 ):1 2 3 4 5 6 7 8 9 10

```

`n.sv` is the number of surrogate variables and `B` is the number of bootstraps. See `?edgeSVA` for additional arguments. To see the terms that have been added to the models:

```

fullModel(edgeObj)

## ~sex + ns(age, df = 4) + SV1 + SV2 + SV3
## <environment: 0x10a8f2b8>

nullModel(edgeObj)

## ~sex + SV1 + SV2 + SV3
## <environment: 0x10a8f2b8>

```

The variables `SV1`, `SV2` and `SV3` are the surrogate variables formed by `sva`. To access the surrogate variables:

```

cov <- pData(edgeObj)
names(cov)

## [1] "sex" "age" "SV1" "SV2" "SV3"

surrogate.vars <- cov[, 3:ncol(cov)]

```

Now `odp` or `lrt` can be used as in previous examples:

```

edgeODP <- odp(edgeObj, verbose = FALSE)
edgeLRT <- lrt(edgeObj, verbose = FALSE)
summary(edgeODP)

##
## ExpressionSet Summary
##
## ExpressionSet (storageMode: lockedEnvironment)
## assayData: 500 features, 72 samples
##   element names: exprs
## protocolData: none
## phenoData
##   sampleNames: 1 2 ... 72 (72 total)
##   varLabels: sex age ... SV3 (5 total)
##   varMetadata: labelDescription
## featureData: none

```

```
## experimentData: use 'experimentData(object)'
## Annotation:
##
## edge Analysis Summary
##
## Total number of arrays: 72
## Total number of probes: 500
##
## Biological variables:
## Null Model:~sex + SV1 + SV2 + SV3
## <environment: 0x10a8f2b8>
##
## Full Model:~sex + ns(age, df = 4) + SV1 + SV2 + SV3
## <environment: 0x10a8f2b8>
##
## .....
##
##
## Statistical significance summary:
## pi0: 0.2503495
##
## Cumulative number of significant calls:
##
##          <1e-04 <0.001 <0.01 <0.025 <0.05 <0.1
## p-value      2      7     30     59     94    146
## q-value      0      0      3     21     55    175
## local fdr    0      0      2     10     27     81
##
##          <1
## p-value    500
## q-value    500
## local fdr  500
```

And to extract the π_0 estimate, q-values, local false discovery rates and p-values:

```
qvalObj <- qvalueObj(edgeODP)
qvals <- qvalObj$qvalues
lfdr <- qvalObj$lfdr
pvals <- qvalObj$pvalues
pi0 <- qvalObj$pi0
```

9 snm: Supervised normalization of microarray data

There has been a lot of work done on separating signal from confounding factors, but a lot of algorithms fail to consider both the models of the study and the technical factors such as batch or array processing date. The `snm` package allows for supervised normalization of microarrays on a gene expression matrix. It takes into account both the experimental models and other technical factors in the experiments. Details on the algorithm can be found in [Mecham et al. \(2010\)](#). The `snm` package is implemented in the `edgeSNM` function. Continuing the analysis on the kidney study in 6:

```
# create models
edgeObj <- edgeStudy(data = kidexpr, adj.var = sex,
  tme = age, basis.df = 4, sampling = "timecourse")
```

Now that we have created `edgeObj`, we can adjust for additional array effects, dye effects and other intensity-dependent effects. In this example, we created array effects that are not existent in the real data set in order to show how to use the function:

```
int.var <- data.frame(array.effects = as.factor(1:72))
edgeObj <- edgeSNM(edgeObj, int.var = int.var, diagnose = FALSE,
  verbose = FALSE)
```

The `int.var` is where the data frame of intensity-dependent effects are inputted, `diagnose` is a flag to let the software know whether to produce diagnostic plots. Additional arguments can be found by typing `?edgeSNM`.

Once the data has been normalized, we can access the normalized matrix by using `exprs`:

```
norm.matrix <- exprs(edgeObj)
```

To run the significance analysis, `odp` or `lrt` can be used:

```
edgeODP <- odp(edgeObj, verbose = FALSE)
summary(edgeODP)

##
## ExpressionSet Summary
##
## ExpressionSet (storageMode: lockedEnvironment)
## assayData: 500 features, 72 samples
##   element names: exprs
## protocolData: none
## phenoData
##   sampleNames: 1 2 ... 72 (72 total)
##   varLabels: adj.var tme
##   varMetadata: labelDescription
## featureData: none
## experimentData: use 'experimentData(object)'
## Annotation:
##
## edge Analysis Summary
##
## Total number of arrays: 72
## Total number of probes: 500
##
## Biological variables:
##   Null Model:~adj.var
##   <environment: 0xab402e0>
##
##   Full Model:~adj.var + ns(tme, df = 4, intercept = FALSE)
##   <environment: 0xab402e0>
##
## .....
```



```
##
##
## Statistical significance summary:
## pi0: 0.3805793
##
## Cumulative number of significant calls:
##
##          <1e-04 <0.001 <0.01 <0.025 <0.05 <0.1
## p-value      1      7    30     61    98   157
## q-value      0      0     1     13    18   103
## local fdr     0      0     0      6    15    52
##
##          <1
## p-value     500
## q-value     500
## local fdr   500
```

And to extract the π_0 estimate, q-values, local false discovery rates and p-values:

```
qvalObj <- qvalueObj(edgeODP)
qvals <- qvalObj$qvalues
lfdr <- qvalObj$lfdr
pvals <- qvalObj$pvalues
pi0 <- qvalObj$pi0
```

10 qvalue: Estimate the q-values

After `odp` or `lrt` is used, the user may wish to change some parameters used when calculating the q-values. This can be done by using the `edgeQvalue` function. Lets review the analysis process for the `kidney` dataset in 6: create the alternative and null models and then run `odp` or `lrt` to get significance results. Applying these steps in the `kidney` dataset:

```
# create models
edgeObj <- edgeStudy(data = kidexpr, adj.var = sex,
  bio.var = age, sampling = "timecourse")
# run significance analysis
edgeObj <- odp(edgeObj, verbose = FALSE)
```

Suppose we wanted to estimate π_0 using the “bootstrap” method in `qvalue` (see [qvalue vignette](#) for more details):

```
old_pi0est <- qvalueObj(edgeObj)$pi0
edgeObj <- edgeQvalue(edgeObj, pi0.method = "bootstrap")
new_pi0est <- qvalueObj(edgeObj)$pi0

##   old_pi0est new_pi0est
## 1  0.7885024  0.7928632
```

In this case, there is a small difference between using the “smoother” method and “bootstrap” method but the point is that the arguments from the `qvalue` package can be passed through `edgeQvalue`. See `edgeQvalue` for additional arguments.

11 Advanced topic: Using the ExpressionSet object

`edge` was designed for complementing `ExpressionSet` objects in significance analysis. The `edgeSet` inherits all the slots from an `ExpressionSet` object and adds vital slots for significance analysis. The rest of this section is for advanced users because it requires knowledge of alternative and null model creation. To begin, lets create an `ExpressionSet` object from the `kidney` dataset:

```
library(edge)
anonDf <- as(data.frame(age=age, sex=sex), "AnnotatedDataFrame")
expSet <- ExpressionSet(assayData = kidexpr,
phenoData = anonDf)
```

In the `kidney` experiment they were interested in finding the effect of age on gene expression. In this case, we handle the time variable, `age`, by fitting a natural spline curve as done in [Storey et al. \(2005\)](#). The relevant models for the experiment can be written as

```
library(splines)
nullMod <- ~1 + sex
altMod <- ~1 + sex + ns(age, intercept = FALSE, df = 4)
```

where `nullMod` is the null model and `altMod` is the alternative model. The `sex` covariate is an adjustment variable while `age` is the biological variable of interest. It is important to note that it is necessary to include the adjustment variables in the formulation of the alternative models as done above.

Having both `expSet` and the hypothesis models, the function `edgeSet` can then be used to create an `edgeSet` object:

```
edgeObj <- edgeSet(expSet, full.model = altMod, null.model = nullMod)
slotNames(edgeObj)

## [1] "null.model"      "full.model"
## [3] "null.matrix"     "full.matrix"
## [5] "individual"      "qvalueObj"
## [7] "experimentData"  "assayData"
## [9] "phenoData"       "featureData"
## [11] "annotation"      "protocolData"
## [13] ".__classVersion__"
```

From the slot names, it is evident that the `edgeSet` object inherits the `ExpressionSet` slots in addition to other slots relating to the significance analysis. See section 6.3 for more details on the `edgeSet` slots. We can now simply run `odp` or `lrt` for significance results:

```
edgeODP <- odp(edgeObj, verbose = FALSE)
edgeLRT <- lrt(edgeObj)
summary(edgeODP)

##
## ExpressionSet Summary
##
## ExpressionSet (storageMode: lockedEnvironment)
## assayData: 500 features, 72 samples
## element names: exprs
## protocolData: none
```

```
## phenoData
##   sampleNames: 1 2 ... 72 (72 total)
##   varLabels: age sex
##   varMetadata: labelDescription
## featureData: none
## experimentData: use 'experimentData(object)'
## Annotation:
##
## edge Analysis Summary
##
## Total number of arrays: 72
## Total number of probes: 500
##
## Biological variables:
##   Null Model:~1 + sex
##
##   Full Model:~1 + sex + ns(age, intercept = FALSE, df = 4)
##
## .....
##
## Statistical significance summary:
## pi0: 0.284494
##
## Cumulative number of significant calls:
##
##           <1e-04 <0.001 <0.01 <0.025 <0.05 <0.1
## p-value         2       8    29     43    84   126
## q-value         0       0     4    10    29   120
## local fdr       0       0     2     6    16    48
##
##           <1
## p-value    500
## q-value    500
## local fdr  500
```

And use the function `qvalueObj` to extract the π_0 estimate, q-values, local false discovery rates and p-values:

```
qvalObj <- qvalueObj(edgeODP)
qvals <- qvalObj$qvalues
lfdr <- qvalObj$lfdr
pvals <- qvalObj$pvalues
pi0 <- qvalObj$pi0
```

Acknowledgements

This software development has been supported in part by funding from the National Institutes of Health and the Office of Naval Research.

References

- SE Calvano, W Xiao, DR Richards, RM Felciano, HV Baker, RJ Cho, RO Chen, BH Brownstein, JP Cobb, SK Tschoeke, C Miller-Graziano, LL Moldawer, MN Mindrinos, RW Davis, RG Tompkins, and SF Lowry. A network-based analysis of systemic inflammation in humans. *Nature*, 437:1032–1037, 2005. doi: 10.1038/nature03985. URL <http://www.nature.com/nature/journal/v437/n7061/full/nature03985.html>.
- Ingrid Hedenfalk, David Duggan, Yidong Chen, Michael Radmacher, Michael Bittner, Richard Simon, Paul Meltzer, Barry Gusterson, Manel Esteller, Mark Raffeld, Zohar Yakhini, Amir Ben-Dor, Edward Dougherty, Juha Kononen, Lukas Bubendorf, Wilfrid Fehrle, Stefania Pittaluga, Sofia Gruvberger, Niklas Loman, Oskar Johannsson, Håkan Olsson, Benjamin Wilfond, Guido Sauter, Olli-P. Kallioniemi, Åke Borg, and Jeffrey Trent. Gene-expression profiles in hereditary breast cancer. *New England Journal of Medicine*, 344(8):539–548, 2001. doi: 10.1056/NEJM200102223440801. URL <http://dx.doi.org/10.1056/NEJM200102223440801>. PMID: 11207349.
- Y Idaghdour, JD Storey, SJ Jadallah, and G Gibson. A genome-wide gene expression signature of environmental geography in leukocytes of moroccan amazighs. *PLoS Genetics*, 4. doi: 10.1371/journal.pgen.1000052.
- Jeffrey T Leek and John D Storey. Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genet*, 3(9):e161, 09 2007. doi: 10.1371/journal.pgen.0030161.
- Jeffrey T. Leek, Eva Monsen, Alan R. Dabney, and John D. Storey. Edge: extraction and analysis of differential gene expression. *Bioinformatics*, 22(4):507–508, 2006. doi: 10.1093/bioinformatics/btk005. URL <http://bioinformatics.oxfordjournals.org/content/22/4/507.abstract>.
- Brigham H. Meacham, Peter S. Nelson, and John D. Storey. Supervised normalization of microarrays. *Bioinformatics*, 26(10):1308–1315, 2010. doi: 10.1093/bioinformatics/btq118. URL <http://bioinformatics.oxfordjournals.org/content/26/10/1308.abstract>.
- Graham E. J Rodwell, Rebecca Sonu, Jacob M Zahn, James Lund, Julie Wilhelmy, Lingli Wang, Wenzhong Xiao, Michael Mindrinos, Emily Crane, Eran Segal, Bryan D Myers, James D Brooks, Ronald W Davis, John Higgins, Art B Owen, and Stuart K Kim. A transcriptional profile of aging in the human kidney. *PLoS Biol*, 2(12):e427, 11 2004. doi: 10.1371/journal.pbio.0020427.
- John D. Storey. The optimal discovery procedure: a new approach to simultaneous significance testing. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(3):347–368, 2007. ISSN 1467-9868. doi: 10.1111/j.1467-9868.2007.005592.x. URL <http://dx.doi.org/10.1111/j.1467-9868.2007.005592.x>.
- John D. Storey, Wenzhong Xiao, Jeffrey T. Leek, Ronald G. Tompkins, and Ronald W. Davis. Significance analysis of time course microarray experiments. *Proceedings of the National Academy of Sciences of the United States of America*, 102(36):12837–12842, 2005. doi: 10.1073/pnas.0504609102. URL <http://www.pnas.org/content/102/36/12837.abstract>.
- John D. Storey, James Y. Dai, and Jeffrey T. Leek. The optimal discovery procedure for large-scale significance testing, with applications to comparative microarray experiments. *Biostatistics*, 8(2):414–432, 2007. doi: 10.1093/biostatistics/kxl019. URL <http://biostatistics.oxfordjournals.org/content/8/2/414.abstract>.
- Sangsoo Woo, Jeffrey T. Leek, and John D. Storey. A computationally efficient modular optimal discovery procedure. *Bioinformatics*, 27(4):509–515, 2011. doi: 10.1093/bioinformatics/btq701. URL <http://bioinformatics.oxfordjournals.org/content/27/4/509.abstract>.