

The locfdr Package

February 16, 2008

Title Computes local false discovery rates

Version 1.1-6

Author Bradley Efron, Brit B. Turnbull and Balasubramanian Narasimhan

Description Computation of local false discovery rates

Maintainer Bradley Efron <brad@stat.Stanford.EDU>

License GPL 2.0

Depends splines

R topics documented:

hivdata	1
lfdrsim	2
locfdr	2
Index	5

hivdata	<i>HIV data set</i>
---------	---------------------

Description

The data comprises 7680 z -values, each relating to a two-sample t -test. The test compares gene expression values for 4 HIV patients with values for 4 normal subjects; the t -score $T[i]$ for gene i has been transformed to a normal scale, $z[i] = \text{qnorm}(\text{pt}(T[i], \text{df}=6))$, so that the $z[i]$'s theoretically would have a standard $N(0, 1)$ distribution under the null hypothesis. The original experiment is described in van't Wout et. al. (2003).

Usage

```
data(hivdata)
```

Format

A vector containing 7680 z -values

References

van't Wout, et. al., Cellular gene expression upon human immuno-deficiency virus type 1 infection of CD4+-T-Cell lines, Journal of Virology 77, 1392-1402.

 lfdrsim

Simulated data set for locfdr

Description

A simulated dataset that involves 2000 "genes", each of which has yielded a test statistic "zex", with $zex[i] \sim N(\mu[i], 1)$ (independently for $i = 1, 2, \dots, 2000$.) The data comprises 2000 μ_i values and 2000 z -values.

Usage

```
data(lfdrsim)
```

Format

A matrix of 2000 rows and 2 columns containing mu and the z-score values (zex)

 locfdr

Local False Discovery Rate Calculation

Description

Compute local false discovery rates, following the definitions and description in references listed below.

Usage

```
locfdr(zz, bre = 120, df = 7, pct = 0, pct0 = 1/4, nulltype = 1, type = 0, plot = 1, mult, mlests, main = " ", sw = 0)
```

Arguments

<code>zz</code>	A vector of summary statistics, one for each case under simultaneous consideration. The calculations assume a large number of cases, say <code>length(zz)</code> exceeding 200. Results may be improved by transforming <code>zz</code> so that its elements are theoretically distributed as $N(0, 1)$ under the null hypothesis. See the <code>locfdr</code> vignette for tips on creating <code>zz</code> .
<code>bre</code>	Number of breaks in the discretization of the z -score axis, or a vector of break-points fully describing the discretization. If <code>length(zz)</code> is small, such as when the number of cases is less than about 1000, set <code>bre</code> to a number lower than the default of 120.
<code>df</code>	Degrees of freedom for fitting the estimated density $f(z)$.
<code>pct</code>	Excluded tail proportions of zz 's when fitting $f(z)$. <code>pct=0</code> includes full range of zz 's. <code>pct</code> can also be a 2-vector, describing the fitting range.
<code>pct0</code>	Proportion of the zz distribution used in fitting the null density $f_0(z)$ by central matching. If a 2-vector, e.g. <code>pct0=c(0.25, 0.60)</code> , the range [<code>pct0[1]</code> , <code>pct0[2]</code>] is used. If a scalar, [<code>pct0</code> , <code>1-pct0</code>] is used.
<code>nulltype</code>	Type of null hypothesis assumed in estimating $f_0(z)$, for use in the <code>fdr</code> calculations. 0 is the theoretical null $N(0, 1)$, 1 is maximum likelihood estimation, 2 is central matching estimation, 3 is a split normal version of 2.
<code>type</code>	Type of fitting used for f ; 0 is a natural spline, 1 is a polynomial, in either case with degrees of freedom <code>df</code> [so total degrees of freedom including the intercept is <code>df+1</code> .]
<code>plot</code>	Plots desired. 0 gives no plots. 1 gives single plot showing the histogram of zz and fitted densities f and $p_0 * f_0$. 2 also gives plot of <code>fdr</code> , and the right and left tail area <code>Fdr</code> curves. 3 gives instead the <code>f1</code> cdf of the estimated <code>fdr</code> curve; <code>plot=4</code> gives all three plots.
<code>mult</code>	Optional scalar multiple (or vector of multiples) of the sample size for calculation of the corresponding hypothetical <code>Efdr</code> value(s).
<code>mlests</code>	Optional vector of initial values for (<code>delta0</code> , <code>sigma0</code>) in the maximum likelihood iteration.
<code>main</code>	Main heading for the histogram plot when <code>plot>0</code> .
<code>sw</code>	Determines the type of output desired. 2 gives a list consisting of the last 5 values listed under <code>Value</code> below. 3 gives the square matrix of dimension <code>bre-1</code> representing the influence function of <code>log(fdr)</code> . Any other value of <code>sw</code> returns a list consisting of the first 5 (6 if <code>mult</code> is supplied) values listed below.

Details

See the `locfdr` vignette for details and tips.

Value

<code>fdr</code>	the estimated local false discovery rate for each case, using the selected <code>type</code> and <code>nulltype</code> .
------------------	--

<code>f_{p0}</code>	the estimated parameters delta (mean of f_0), sigma (standard deviation of f_0), and p_0 , along with their standard errors.
<code>E_{fdr}</code>	the expected false discovery rate for the non-null cases, a measure of the experiment's power as described in Section 3 of the second reference. Overall <code>E_{fdr}</code> and right and left values are given, both for the specified <code>nulltype</code> and for <code>nulltype 0</code> . If <code>nulltype==0</code> , values are given for <code>nulltypes 1</code> and <code>0</code> .
<code>cdf1</code>	a 99x2 matrix giving the estimated cdf of <code>fdr</code> under the non-null distribution <code>f1</code> . Large values of the cdf for small <code>fdr</code> values indicate good power; see Section 3 of the second reference. Set <code>plot</code> to 3 or 4 to see the <code>cdf1</code> plot.
<code>mat</code>	A matrix of estimates of $f(x)$, $f_0(x)$, $f_{dr}(x)$, etc. at the <code>bre - 1</code> midpoints "x" of the break discretization, convenient for comparisons and plotting. Details are in the <code>locfdr</code> vignette.
<code>z.2</code>	the interval along the <code>zz</code> -axis outside of which $f_{dr}(z) < 0.2$, the locations of the yellow triangles in the histogram plot. If no elements of <code>zz</code> on the left or right satisfy the criterion, the corresponding element of <code>z.2</code> is NA.
<code>call</code>	the function call.
<code>mult</code>	If the argument <code>mult</code> was supplied, vector of the ratios of hypothetical <code>E_{fdr}</code> for the supplied multiples of the sample size to <code>E_{fdr}</code> for the actual sample size.
<code>pds</code>	The estimates of p_0 , delta, and sigma.
<code>x</code>	The bin midpoints.
<code>f</code>	The values of $f(z)$ at the bin midpoints.
<code>pds.</code>	The derivative of the estimates of p_0 , delta, and sigma with respect to the bin counts.
<code>stdev</code>	The delta-method estimates of the standard deviations of the p_0 , delta, and sigma estimates.

Author(s)

Bradley Efron, Brit B. Turnbull, and Balasubramanian Narasimhan

References

Efron, B. (2004) "Large-scale simultaneous hypothesis testing: the choice of a null hypothesis", *Jour Amer Stat Assoc*, **99**, pp. 96–104

Efron, B. (2006) "Size, Power, and False Discovery Rates"

Efron, B. (2007) "Correlation and Large-Scale Simultaneous Significance Testing", *Jour Amer Stat Assoc*, **102**, pp. 93–103

<http://www-stat.stanford.edu/~brad/papers/>

Examples

```
## HIV data example
data(hivdata)
w <- locfdr(hivdata)
```

Index

*Topic **datasets**

hivdata, 1

lfdrsim, 2

*Topic **htest**

locfdr, 2

*Topic **models**

locfdr, 2

hivdata, 1

lfdrsim, 2

locfdr, 2