

# Introduction to the les package: Identifying Loci of Enhanced Significance in Tiling Array Data

Julian Gehring

July 20, 2010

## Abstract

In this vignette we describe using the *les* package for finding Loci of Enhanced Significance (LES) in tiling microarray data. With an example of a general framework we illustrate how to apply the package for exploring regions of regulation with differential design experiments.

## 1 Introduction

Tiling microarrays have become an important platform for the investigation of regulation in expression and DNA-protein interaction. Due to their lower bias towards annotation compared to other microarray platforms they provide an powerful tool for biological research.

Beside the analysis of single microarrays the investigation of differential effects between experimental conditions is important for current research. A common approach consists in applying some suitable statistical tests on the level of single probes and thereby computing p-values  $p_i$  for each probe  $i$  independently. Since the targets of such experiments cover regions with several probes a reasonable next step involves combining information from neighboring probes. Many approaches use a smoothing window to obtain this. While these methods may work well most of them produce a feature that lacks any statistical interpretation.

In regions with differential effects the test statistics change their distribution and are referred to as *Loci of Enhanced Significance (LES)*. The changes in the test statistics depend on the underlying test applied.

The *les* package provides the ability to detect such LES independent of the underlying statistical test and can therefore be used for a wide range of applications. This vignette illustrates how LES can be found in tiling microarray data sets.

## 2 Data and statistics on probe level

For this analysis we will use a simulated data set describing differential expression between two conditions.

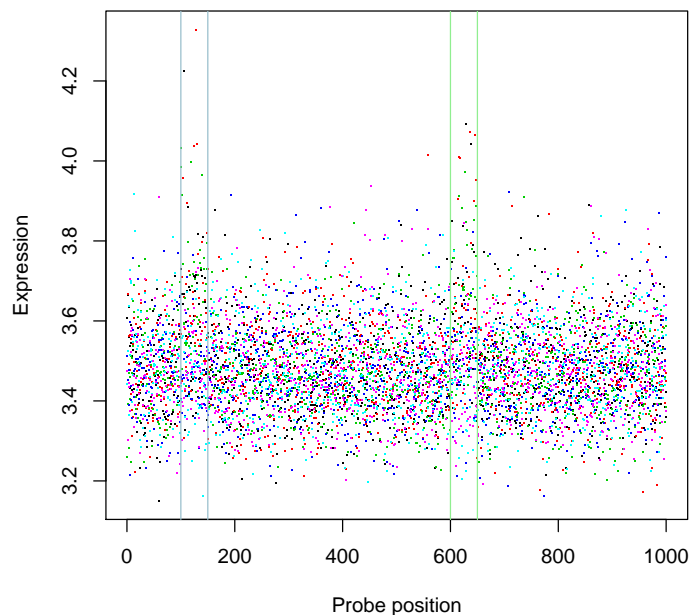
Please note that this data does not reflect real biological data. It is rather meant to illustrate the workflow than to provide real life data.

The data set contains 1000 probes with 3 chips each for the conditions treatment and control. The expression values are stored in an expression set. We will extract the position of the probes, the conditions of the samples and the expression values. There are two regions of regulation present in the data each 50 bp long.

```
> library(les)
> library(Biobase)
> data(simTile)
> treatment <- as.factor(phenoData(simTile)$condition ==
+   "treatment")
> pos <- featureData(simTile)$position
> exprs <- exprs(simTile)
> regions <- c(100, 150, 600, 650)
> cols <- rep(c("lightblue3", "lightgreen"), each = 2)
```

Next we have a look at the expression values.

```
> matplot(exprs, pch = ".", xlab = "Probe position",
+   ylab = "Expression")
> abline(v = regions, col = cols)
```

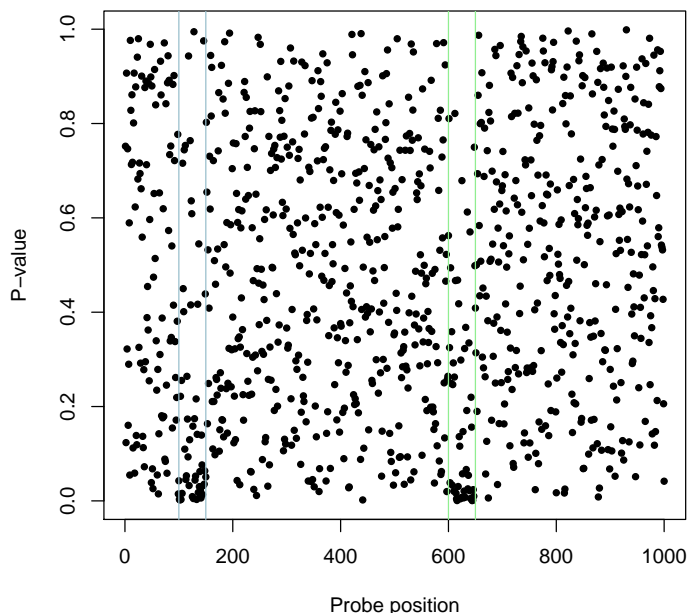


In the next step we will compute the statistics of changes between the two conditions for each probe. Since the sample size as for most tiling microarray experiments is small we will use a modified t-test provided by the *les* package.[2]

We will also plot the p-values  $p_i$  against the probe positions  $i$ . By looking at the raw p-values it may be hard to detect the *LES*.

```
> library(limma)
> design <- cbind(mean = 1, diff = treatment)
> fit <- lmFit(exprs, design)
> fit <- eBayes(fit)
> pval <- as.numeric(fit$p.value[, "diff"])

> plot(pos, pval, pch = 20, xlab = "Probe position",
+       ylab = "P-value")
> abline(v = regions, col = cols)
```



### 3 Incorporation information from neighboring p-values

In any well designed tiling array experiments the potential targets will cover several neighboring probes. Thereby neighboring p-values should contain mutual information and incorporation of such will be beneficial.

In the *les* package this is done in the following manner: For each probe  $i$  the surrounding p-values  $p_i$  get weights assigned by a windowing function. A weighted cumulative density is then computed and the fraction of significant  $p_i$  is estimated by iterative linear fitting. The method is based on the fact that p-values under the null hypothesis  $H_0$  come from a uniform distribution whereas p-values violating  $H_0$  are shifted towards smaller values.[1] This results in the index  $\Lambda_i$  which measures the fraction of p-values violating  $H_0$  in the window and therefore the degree of regulation in the local surrounding. It should be noted that this approach is closely related to the estimation of a false discovery rate and  $\Lambda_i$  can be interpreted in a similar way.

For the analysis we will first store our data in an object of class *Les* by calling the `create` function. The only data required for the analysis are the position of the probes  $i$ , the corresponding p-values  $p_i$  from the statistical test and optionally their chromosomal location.

```
> res <- create(pos, pval)
```

Then we can compute our first estimate of  $\Lambda_i$  for which we have to specify a window size. The power of the detection will be high if the size of the window matches the size of the target. In many experiments a rough prior knowledge on the target size is available which will be sufficient of the first step. Here we will chose a small window size to start with.

By default a triangular weighting function will be chosen. We can also take a rectangular window or write our own function and pass it via the weighting argument. We can further on specify whether we want to include the Grenander correction for the cumulative density or use multicore processing on some platforms.

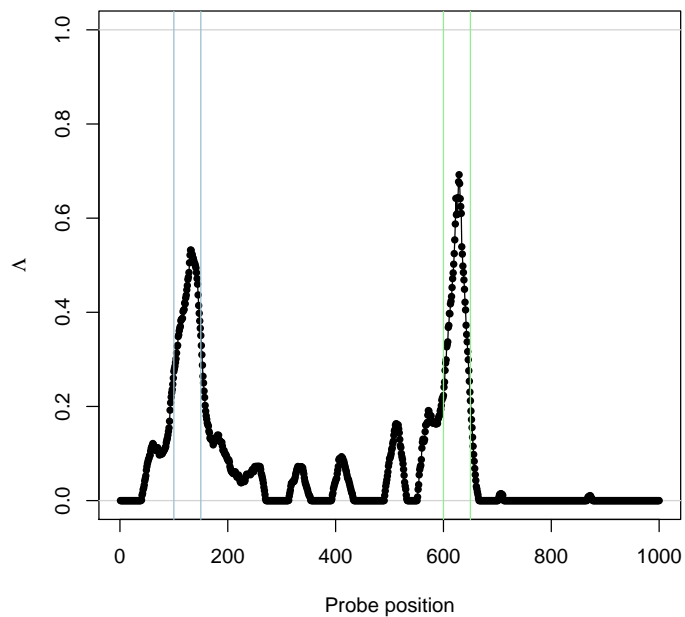
```
> win <- 30
> res <- estimate(res, win)
```

All data, results and parameters are now stored in the object named `res`. We can get a short summary on the results by calling `print`, `summary` or by plotting it.

```
> res

** Object of class 'Les' **
* 1000 probes on 1 chromosomes
* Lambda in range [0, 0.692297] with kernel size 30

> plot(res)
> abline(v = regions, col = cols)
```



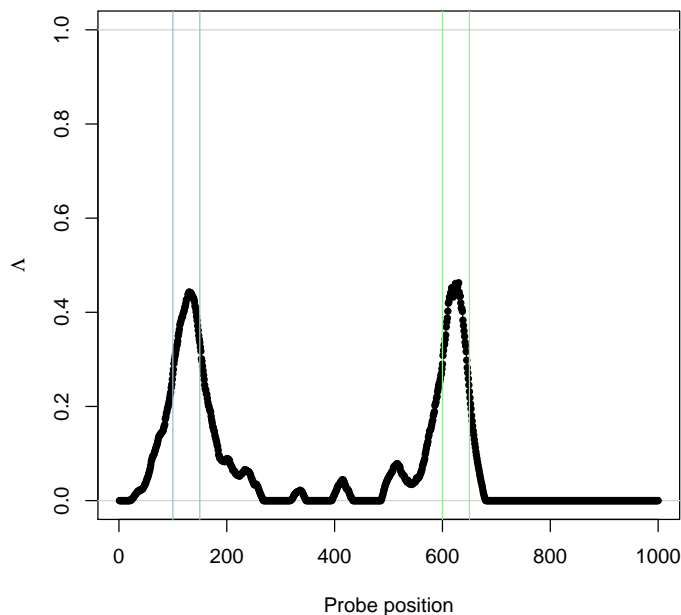
For comparison we will analyze and plot the same data with a different window size. This allows us to explore our data set.

```
> win2 <- 50
> res <- estimate(res, win2)

> res

** Object of class 'Les' **
* 1000 probes on 1 chromosomes
* Lambda in range [0, 0.462787] with kernel size 50

> plot(res)
> abline(v = regions, col = cols)
```



We can already see two distinct peaks that correspond well to the simulated regions of regulation.

The `plot` method provides additional arguments that help customizing the figure.

## 4 Parameter estimation from the data

To turn the continuous  $\Lambda_i$  into distinct regions of interest we have to define a threshold  $\Theta$ . It can be derived from the data by estimating the number of probes with a significant effect  $R$  on the whole array. Then  $\Theta$  can be chosen such that  $|\Lambda_i \geq \Theta| = R$ . The content of any slot can be accessed by using the `/`-method.

```
> res <- threshold(res, grenander = TRUE, verbose = TRUE)
```

```
[1] "51 significant probes estimated with limit Lambda>=0.401031"
```

Based on  $\Theta$  we can search for regions that have a continuous  $\Lambda_i \geq \Theta$ . The `regions` method takes by default the estimated  $\hat{\Theta}$  as shown before. We can also pass our own estimate for  $\Theta$  with the `limit` argument. Further restrictions can be imposed on the regions search such as the minimal length of a region and the maximum gap allowed between probes of one region.

A data frame with the estimated regions can be accessed with the `/`-method. This can also be used to access any other data slot of a LES object.

```
> theta <- 0.3
> res <- regions(res, limit = theta, verbose = TRUE)

[1] "2 regions found for lambda>=0.3"

> res

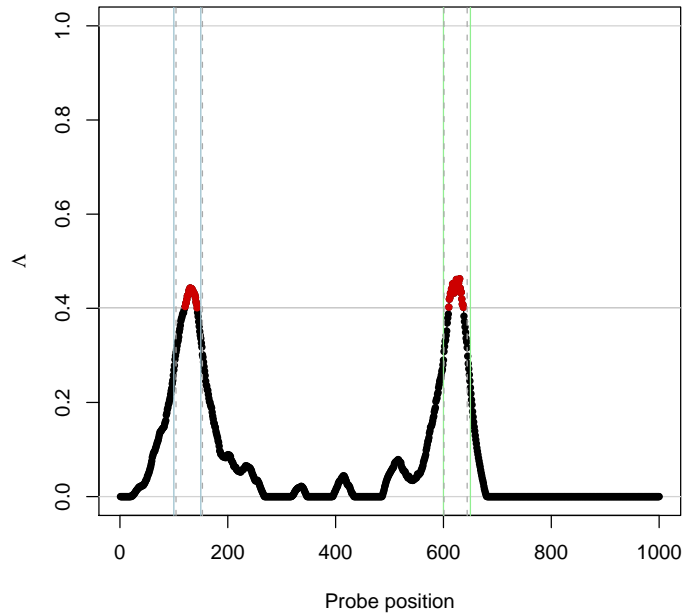
** Object of class 'Les' **
* 1000 probes on 1 chromosomes
* Lambda in range [0, 0.462787] with kernel size 50
* 52 regulated probes estimated for lambda >= 0.401031
* 2 regions detected

> res["regions"]

  chr start end size nProbes    ri    se    rs
1   0   104 153   50      50 0.4568 0.0103 44.4261
2   0   601 644   44      44 0.5616 0.0255 22.0044

> region <- res["regions"]
> borders <- c(region$start, region$end)
> plot(res)
> abline(v = regions, col = cols)
> abline(v = borders, col = "darkgray", lty = 2)
```



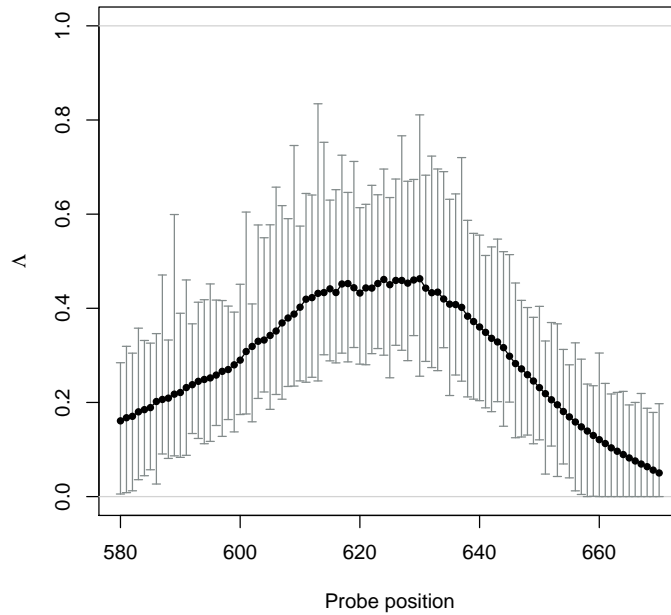


## 5 Calculation of confidence intervals

In some cases it is also useful to provide confidence (CI) intervals for  $\Lambda_i$ . These are computed by bootstrapping the probes in the window. Since bootstrapping is by its nature computationally demanding and CI are primarily interesting in regions of interest it is possible to compute CIs for a subset of probes and to specify the number of bootstraps.

```
> subset <- pos >= 580 & pos <= 670
> res <- ci(res, subset, nBoot = 50)

> plot(res, error = "ci", limit = theta, xlim = c(580,
+       670))
```



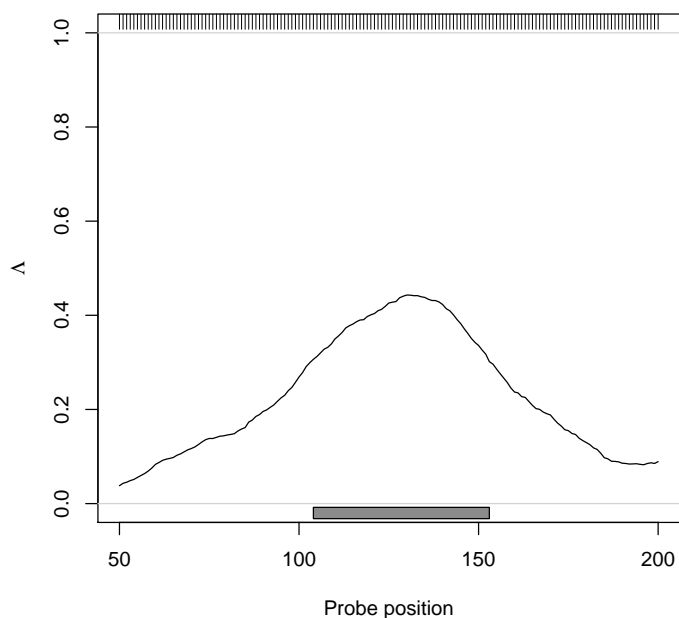
## 6 Plotting capabilities

The `plot` method provides many options for customizing result graphics.

The following command plots a smaller region of the chromosome with confidence intervals. Additionally the estimated regions, the position of probes are shown.

```
> plot(res, error = "ci", region = TRUE, rug = TRUE,
+      rugSide = 3, main = "LES for simulated data",
+      probePch = NA, sigCol = "lightblue", xlim = c(50,
+      200), limit = FALSE)
```

### LES for simulated data



## 7 Exporting result to external software

In some cases it can be useful to analyze results of this analysis in other software than R alone. For this purpose both the estimated regions as well as  $\Lambda$  can be saved to a file with the *export* method. Available formats for the regions are the *bed* and *gff*, for  $\Lambda$  *wig*. Since these formats are widely used they can be directly loaded into many genome software packages and browsers.

## 8 Specification of own window functions

With the `triangWeight`, `rectangWeight` and `gaussWeight` three window functions are already include in the *les* package. We can also specify own window functions and pass it via the `weighting` argument in the *estimate* method. They have to be specified in the following format, here described with a triangular weighting.

```
> weightFcn <- function(distance, win) {
+   weight <- 1 - distance/weight
+   return(weight)
+ }
```

## References

- [1] Kilian Bartholomé, Clemens Kreutz, and Jens Timmer. Estimation of gene induction enables a Relevance-Based ranking of gene sets. *Journal of Computational Biology*, 16(7):959–967, 2009. ISSN 1066-5277. doi: 10.1089/cmb.2008.0226. URL <http://www.liebertonline.com/doi/abs/10.1089/cmb.2008.0226>.
- [2] G. Smyth. limma: Linear models for microarray data. In *Bioinformatics and Computational Biology Solutions Using R and Bioconductor*, pages 397–420. 2005. URL [http://dx.doi.org/10.1007/0-387-29362-0\\_23](http://dx.doi.org/10.1007/0-387-29362-0_23).

## Session information

- R version 2.11.1 (2010-05-31), x86\_64-pc-linux-gnu
- Base packages: base, datasets, graphics, grDevices, methods, stats, tools, utils
- Other packages: Biobase 2.9.0, fdrtool 1.2.6, les 0.9.7, limma 3.5.12
- Loaded via a namespace (and not attached): annotate 1.27.1, AnnotationDbi 1.11.4, boot 1.2-42, DBI 0.2-5, gdata 2.8.0, genefilter 1.31.0, gplots 2.8.0, GSRI 1.1.4, gtools 2.6.2, RColorBrewer 1.0-2, RSQLite 0.9-1, splines 2.11.1, survival 2.35-8, xtable 1.5-6