

LPE test for microarray data with a small number of replicates

Nitin Jain, Michael O Connell, and Jae K. Lee

July 22, 2005

Contents

1	Introduction	1
2	Mouse Immune Response Study dataset	1
3	Discussion	6

1 Introduction

The *LPE* package describes local-pooled-error (LPE) test for identifying significant differentially expressed genes in microarray experiments. Local pooled error test is especially useful when the number of replicates is low (2-3) ?. LPE estimation is based on pooling errors within genes and between replicate arrays for genes in which expression values are similar. This is motivated by the observation that errors between duplicates vary as a function of the average gene expression intensity and by the fact that many gene expression studies are implemented with a limited number of replicated arrays ?.

Step by step analysis is presented in Section 2 using data from a 6-chip oligonucleotide microarray study of a mouse immune response study.

Details of methodology and application of Local Pooled Error (LPE) test can be obtained from the LPE paper, published in Bioinformatics ?.

2 Mouse Immune Response Study dataset

Load the library

```
> set.seed(0) # To have reproducible results
> library(LPE)
```

```

> data(Ley)

> dim(Ley)

[1] 12488 7

> Ley[1:3,]
      ID    c1    c2    c3    t1    t2    t3
1 AFX-MurIL2_at 16.0 14.1 19.3 2782.7 2861.3 2540.2
2 AFX-MurIL10_at 22.7  6.9 28.2   18.6   12.7    7.5
3 AFX-MurIL4_at 33.9 17.1 23.9   24.9   25.2   24.9

> Ley[,2:7] <- preprocess(Ley[,2:7], data.type = "MAS5")

> Ley[1:3,]
      ID    c1    c2    c3    t1    t2    t3
1 AFX-MurIL2_at 4.058556 4.077996 4.419651 11.681194 11.668376 11.528191
2 AFX-MurIL10_at 4.563176 3.058342 4.960576  5.125427  4.303960  3.707980
3 AFX-MurIL4_at 5.141769 4.327081 4.718751  5.484593  5.255005  5.330675

```

Mouse immune response study was conducted by Dr. Klaus Ley, Univeristy of Virginia. Three replicates of Affymetrix oligonucleotide chips per condition were used. Based on M vs A scater plot matrix, IQR normalization was performed, so that interquartile ranges on all chips are set to their widest range. It is performed by multiplying by a scaling factor. Note that this is a simple constant-scale & location normalization step. Finally log based 2 transformation was done. Replcates of Naive condition are named as c1, c2, c3 and those of Actiavted condition are named as t1, t2 and t3 respectively.

Remove the control spots

```

> Ley <- Ley[substring(Ley$ID,1,4) != "AFX",]

> dim(Ley)

[1] 12422 7

```

```
> Ley[1:3,]
```

	ID	c1	c2	c3	t1	t2	t3
67	92539_at	11.999273	12.253685	12.398052	11.924385	12.042756	11.824377
68	92540_f_at	8.948516	9.034942	8.674348	11.284850	11.323132	11.289058
69	92541_at	6.242440	6.223671	6.185748	5.866883	5.775228	6.430501

Calculate the baseline error distribution of Naive condition, which returns a dataframe of A vs M for selected number of bins ($= 1/q$), where $q = \text{quantile}$.

```
> var.Naive <- baseOlig.error(Ley[,2:4],q=0.01)
```

```
> dim(var.Naive)
```

```
[1] 12422      2
```

```
> var.Naive[1:3,]
```

	A	var.M
67	12.253685	0.04541592
68	8.948516	0.05166602
69	6.223671	0.24738118

Similarly calculate the base-line distribution of Activated condition:

```
> var.Activated <- baseOlig.error(Ley[,5:7], q=0.01)
```

```
> dim(var.Activated)
```

```
[1] 12422      2
```

```
> var.Activated[1:3,]
```

	A	var.M
67	11.924385	0.01296899
68	11.289058	0.01884620
69	5.866883	0.27232983

Calculate the lpe variance estimates as described above. The function *lpe* takes the first two arguments as the replicated data, next two arguments as the baseline distribution of the replicates calculated from the *baseOlig.error* function, Gene IDs as probe.set.name. Adjustment for multiple comparison is applied using Bioconductor's multtest package (Dudoit et. al.)

```
> lpe.val <- data.frame(lpe(Ley[,5:7], Ley[,2:4], var.Activated, var.Naive,
  probe.set.name=Ley$ID))
```

```
> lpe.val <- round(lpe.val, digits=2)
```

```
> dim (lpe.val)
```

```
[1] 12422 13
```

```
> lpe.val[1:3,]
```

	x.t1	x.t2	x.t3	median.1	std.dev.1	y.c1	y.c2	y.c3	median.2
92539_at	11.92	12.04	11.82	11.92	0.11	12.00	12.25	12.40	12.25
92540_f_at	11.28	11.32	11.29	11.29	0.14	8.95	9.03	8.67	8.95
92541_at	5.87	5.78	6.43	5.87	0.52	6.24	6.22	6.19	6.22

	std.dev.2	median.diff	pooled.std.dev	z.stats
92539_at	0.21	-0.33	0.17	-1.88
92540_f_at	0.23	2.34	0.19	12.18
92541_at	0.50	-0.36	0.52	-0.68

Doing FDR correction

```
> fdr.BH <- fdr.adjust(lpe.val, adjp="BH")
```

```
> dim(fdr.BH)
```

```
[1] 12422 2
```

```
> round(fdr.BH[1:4, ],2)
```

	FDR	z.real
7915	0	26.66
5557	0	24.68

```
344    0 24.22
5985   0 24.22
```

Resampling based FDR adjustment takes a while to run, and returns the critical z-values and corresponding FDR.

```
> fdr.2 <- fdr.adjust(lpe.val, adjp="resamp", iterations=2)
```

```
iteration number 1 is in progress
iteration number 1 finished
iteration number 2 is in progress
iteration number 2 finished
Computing FDR...
```

```
> fdr.2
```

	target.fdr	z.critical
[1,]	0.001	3.54
[2,]	0.010	2.69
[3,]	0.020	2.37
[4,]	0.030	2.17
[5,]	0.040	2.05
[6,]	0.050	1.94
[7,]	0.060	1.85
[8,]	0.070	1.76
[9,]	0.080	1.69
[10,]	0.090	1.63
[11,]	0.100	1.56
[12,]	0.150	1.31
[13,]	0.200	1.12
[14,]	0.500	0.25

Note that above table may differ slightly due to generation of 'NULL distribution' by resampling. For each target.fdr, we can note critical z-value, above which all genes are considered significant.

3 Discussion

Using our LPE approach, the sensitivity of detecting subtle expression changes can be dramatically increased and differential gene expression patterns can be identified with both small false-positive and small false-negative error rates. This is because, in contrast to the individual gene's error variance, the local pooled error variance can be estimated very accurately.

Acknowledgments. We wish to acknowledge the following colleagues: P. Aboyoun, J. Betcher, D. Clarkson, J. Gibson, A. Hoering, S. Kaluzny, L. Kannapel, D. Kinsey, P. McKinnis, D. Stanford, S. Vega and H. Yan.