

Using lumi, a package processing Illumina Microarray

Pan Du^{†*}, Gang Feng^{††}, Warren A. Kibbe^{‡‡}, Simon Lin^{§§}

March 3, 2010

[‡]Robert H. Lurie Comprehensive Cancer Center
Northwestern University, Chicago, IL, 60611, USA

Contents

1	Overview of lumi	2
2	Citation	2
3	Installation of lumi package	3
4	Object models of major classes	3
5	Data preprocessing	3
5.1	Intelligently read the BeadStudio output file	3
5.2	Quality control of the raw data	8
5.3	Background correction	18
5.4	Variance stabilizing transform	18
5.5	Data normalization	19
5.6	Quality control after normalization	22
5.7	Encapsulate the processing steps	22
5.8	Inverse VST transform to the raw scale	29
6	Handling large data sets	30
7	Performance comparison	32
8	Gene annotation	32
9	A use case: from raw data to functional analysis	32
9.1	Preprocess the Illumina data	33
9.2	Identify differentially expressed genes	34
9.3	Gene Ontology and other functional analysis	35
9.4	GEO submission of the data	36

*dupan (at) northwestern.edu

†g-feng (at) northwestern.edu

‡wakibbe (at) northwestern.edu

§s-lin2 (at) northwestern.edu

10 Session Info	37
11 Acknowledgments	38
12 References	38

1 Overview of lumi

Illumina microarray is becoming a popular microarray platform. The BeadArray technology from Illumina makes its preprocessing and quality control different from other microarray technologies. Unfortunately, until now, most analyses have not taken advantage of the unique properties of the BeadArray system. The *lumi* Bioconductor package especially designed to process the Illumina microarray data. The *lumi* package provides an integrated solution for the bead-level Illumina microarray data analysis. The package covers data input, quality control, variance stabilization, normalization and gene annotation.

The *lumi* package provides unique functions for microarray processing. It includes a variance-stabilizing transformation (VST) algorithm [2] that takes advantage of the technical replicates available on every Illumina microarray. A robust spline normalization (RSN), which combines the features of the quantile and loess normalization, and simple scaling normalization (SSN) algorithms are also implemented in this package. Options available in other popular normalization methods are also provided. Multiple quality control plots for expression and control probe data are provided in the package. To better annotate the Illumina data, a new, vendor independent nucleotide universal identifier (nuID) [3] was devised to identify the probes of Illumina microarray. The nuID indexed Illumina annotation packages is compatible with other Bioconductor annotation packages. Mappings from Illumina Target Id or Probe Id to nuID are also included in the annotation packages. The output of lumi processed results can be easily integrated with other microarray data analysis, like differentially expressed gene identification, gene ontology analysis or clustering analysis.

2 Citation

For the people using lumi package, please cite the following papers in your publications.

* For the package:

Du, P., Kibbe, W.A. and Lin, S.M., (2008) 'lumi: a pipeline for processing Illumina microarray', *Bioinformatics* 24(13):1547-1548

* For the VST (variance stabilization transformation) algorithm, please cite:

Lin, S.M., Du, P., Kibbe, W.A., (2008) 'Model-based Variance-stabilizing Transformation for Illumina Microarray Data', *Nucleic Acids Res.* 36, e11

* For nuID annotation packages, please cite:

Du, P., Kibbe, W.A. and Lin, S.M., (2007) 'nuID: A universal naming schema of oligonucleotides for Illumina, Affymetrix, and other microarrays', *Biology Direct*, 2, 16.

Thanks for your help!

3 Installation of lumi package

In order to install the lumi package, the user needs to first install R, some related Bioconductor packages. You can easily install them by the following codes.

```
source("http://bioconductor.org/biocLite.R")
biocLite("lumi")
```

For the users want to install the latest developing version of lumi, which can be downloaded from the developing section of Bioconductor website. Some additional packages may be required to be installed because of the update the Bioconductor. These packages can also be found from the developing section of Bioconductor website. You can also directly install the source packages from the Bioconductor website by specify the developing version number, which can be found at the Bioconductor website. Suppose the developing version is 2.3, to install the latest lumi package in the Bioconductor developing version, you can use the following command:

```
install.packages("lumi", repos="http://www.bioconductor.org/packages/2.3/bioc", type="source")
```

An Illumina benchmark data package *lumiBarnes* can be downloaded from Bioconductor Experiment data website.

4 Object models of major classes

The *lumi* package has one major class: **LumiBatch**. **LumiBatch** is inherited from **ExpressionSet** class in Bioconductor for better compatibility. Their relations are shown in Figure 1. **LumiBatch** class includes *se.exprs*, *beadNum* and *detection* in **assayData** slot for additional informations unique to Illumina microarrays. A *controlData* slot is used to keep the control probe information, and a *QC* slot is added for keeping the quality control information. The S4 function **plot** supports different kinds of plots by specifying the specific plot type of **LumiBatch** object. See help of **plot-methods** function for details. The *history* slot records all the operations made on the **LumiBatch** object. This provides data provenance. Function **getHistory** is to retrieve the *history* slot. Please see the help files of **LumiBatch** class for more details. A series of functions: **lumiR**, **lumiR.batch**, **lumiB**, **lumiT**, **lumiN** and **lumiQ** were designed for data input, preprocessing and quality control. Function **lumiExpresso** encapsulates the preprocessing methods for easier usability.

5 Data preprocessing

The first thing is to load the *lumi* package.

```
> library(lumi)
```

5.1 Intelligently read the BeadStudio output file

The **lumiR** function supports directly reading the Illumina raw data output of the Illumina Bead Studio toolkit from version 1 to version 3. It can automatically detect the BeadStudio output version and format and create a new



Figure 1: Object models in lumi package

LumiBatch object for it. An example of the input data format is shown in Figure 2. For simplicity, only part of the data of first sample is shown. The data in the highlighted columns are kept in the corresponding slots of **LumiBatch** object, as shown in Figure 2. The `lumiR` function will automatically determine the starting line of the data. The columns with header including **AVG_Signal** and **BEAD_STD** are required for the **LumiBatch** object. By default, the sample IDs and sample labels are extracted from the column names of the data file. For example, based on the column name: **AVG_Signal-1304401001_A**, we will extract "1304401001" as the sample ID and "A" as the sample label (The function assumes the separation of the sample ID and the sample label is "_" if it exists in the column name.). The function will check the uniqueness of sample IDs. If the sample ID is not unique, the entire portion after removing "AVG_Signal" will be used as a sample ID. The user can suppress this parsing by setting the parameter "parseColumnName" as FALSE.

The `lumiR` will automatically initialize the QC slot of the **LumiBatch** object by calling `lumiQ`. If BeadStudio outputted the control probe data, their information will be kept in the `controlData` slot of the **LumiBatch** object. If BeadStudio outputted the sample summary information, which is called [Samples Table] in the output text file, the information will be kept in `BeadStudio-Summary` within the QC slot of the **LumiBatch** object.

The BeadStudio can output the gene profile or the probe profile. As the probe profile provides unique mapping from the probe Id to the expression profile, outputting probe profile is preferred. When the probe profile is outputted, as show in Figure 2(B), the `ProbeId` column will be used as the identifier of **LumiBatch** object.

We strongly suggest outputting the header information when using BeadStudio, as shown in Figure 2. Please refer to the separate document ("Resolve the Inconsistency of Illumina Identifiers through nuID Annotation") in the `lumi` package for more details of the changing of BeadStudio output formats.

The recent version of BeadStudio can also output the annotation information together with the expression data. In the users also want to input the annotation information, they can set the parameter "inputAnnotation" as TRUE. At the same time, they can also specify which columns to be inputted by setting parameter "annotationColumn". The BeadStudio annotation columns include: SPECIES, TRANSCRIPT, ILMN_GENE, UNIGENE_ID, GI, ACCESSION, SYMBOL, PROBE_ID, ARRAY_ADDRESS_ID, PROBE_TYPE, PROBE_START, PROBE_SEQUENCE, CHROMOSOME, PROBE_CHR_ORIENTATION, PROBE_COORDINATES, DEFINITION, ONTOLOGY_COMPONENT, ONTOLOGY_PROCESS, ONTOLOGY_FUNCTION, SYNONYMS, OBSOLETE_PROBE_ID. As the annotation data is huge, by default, we only input: ACCESSION, SYMBOL, PROBE_START, CHROMOSOME, PROBE_CHR_ORIENTATION, PROBE_COORDINATES, DEFINITION. As some annotation information may be outdated. We recommend using Bioconductor annotation packages to retrieve the annotation information.

For convenience, another function `lumiR.batch` is designed to input files in batch. Basically it combines the output of each file. See the help of `lumiR.batch` for details. `lumiR.batch` function also allows users to add sample information in the `phenoData` slot of the **LumiBatch** object. This will be useful in the data analysis. The `sampleInfoFile` parameter is optional. The file is a Tab-separated text file. The first ID column is required. It represents sample ID, which is de-

Illumina Inc. BeadStudio version 1.4.0.1						
Normalization = none						
Array Content = 11188230_100CP_MAGE-ML.XML						
Error Model = none						
D	<i>Id</i>	e = 2/3	<i>expr</i> slot	1	<i>se.expr</i> slot	<i>beadNum</i> slot
Local Settings = en-U						
TargetID	AVG_Signal-1	BEAD_STDEV-	Avg_NBEADS-	Detection-10	MIN_Signal-9	
GI_10047089	179.5	9.7	19	0.97076323	182.5	
GI_10047091	144.5	12.3	19	0.55569952	141.8	
GI_10047093	699.7	31.9	18	1	811.9	
GI_10047097	2069.9	78.1	14	1	2405.6	
GI_10047099	163.7	6	34	0.86485123	595.1	
GI_10047103	3487.6	112.6	15	1	4427.8	
GI_10047105	212.4	34	13	0.99980148	227.4	

(A) BeadStudio version 1

[Header]						
BSGX Version 3.0.14						
Report Date 3/8/07 6:56						
Project DianePalmer3_7_07						
Group Set NonNormalized						
Analysis	No	<i>Id</i>	<i>expr</i> slot	<i>se.expr</i> slot	<i>beadNum</i> slot	<i>detection</i> slot
Normalization none						
[Sample Probe Profile]						
TargetID	ProbeID	AVG_Signal	BEAD_STDEV	Avg_NBEADS	Detection Pvi	
ILMN_10000	6960451	46.68013	1.546319	53	0.4011299	
ILMN_10001	2600731	44.7272	1.645874	49	0.569209	
ILMN_10002	2120309	38.04584	1.262413	43	0.9533898	
ILMN_10004	7510608	51.82488	2.436115	36	0.1115819	
ILMN_995	1980743	38.54818	1.346273	30	0.9449152	

(B) BeadStudio version 3

Figure 2: An example of the input data format

defined based on the column names of BeadStudio output file. For example, sample ID of column "1881436070_A_STA.AVG_Signal" is "1881436070_A_STA". The sample ID column can also be found in the "Samples Table.txt" file output by BeadStudio. Another "Label" column (if provided) will be used as the sample-Names of LumiBatch object. All information of sampleInfoFile will be directly added in the phenoData slot in LumiBatch object, which can be retrieved by the command like: pData(phenoData(x.lumi)).

```
> ## specify the file name
> # fileName <- 'Barnes_gene_profile.txt' # Not Run
> ## load the data
> # x.lumi <- lumiR.batch(fileName, sampleInfoFile='sampleInfo.txt') # Not
```

Here, we just load the pre-saved example data, example.lumi, which is a subset of the experiment data package *lumiBarnes* in the Bioconductor. The example data includes four samples "A01", "A02", "B01" and "B02". "A" and "B" represent different Illumina slides (8 microarrays on each slide), and "01" and "02" represent different samples. That means "A01" and "B01" are technique replicates at different slides, the same for "A02" and "B02".

```
> ## load example data (a LumiBatch object)
> data(example.lumi)
> ## summary of the example data
> example.lumi
```

Summary of data information:

```
Illumina Inc. BeadStudio version 1.4.0.1
      Normalization = none
      Array Content = 11188230_100CP_MAGE-ML.XML
      Error Model = none
      DateTime = 2/3/2005 3:21 PM
      Local Settings = en-US
```

Major Operation History:

	submitted	finished		command	lumiVersion
1	2007-04-22 00:08:36	2007-04-22 00:10:36		lumiR("../data/Barnes_gene_profile.txt")	1.1.6
2	2007-04-22 00:10:36	2007-04-22 00:10:38		lumiQ(x.lumi = x.lumi)	1.1.6
3	2007-04-22 00:13:06	2007-04-22 00:13:10		addNuId2lumi(x.lumi = x.lumi, lib = "lumiHumanV1")	1.1.6
4	2007-04-22 00:59:20	2007-04-22 00:59:36		Subsetting 8000 features and 4 samples.	1.1.6

Object Information:

```
LumiBatch (storageMode: lockedEnvironment)
assayData: 8000 features, 4 samples
  element names: beadNum, detection, exprs, se.exprs
phenoData
```

```

sampleNames: A01, A02, B01, B02
varLabels and varMetadata description:
  sampleID: The unique Illumina microarray Id
  label: The label of the sample
featureData
  featureNames: oZsQEQXp9ccVIlwoQo, 9qedFRd_5Cul.ueZeQ, ..., 33KnLHy.RFaieogAF4 (8000 tot
  fvarLabels and fvarMetadata description:
    TargetID: The Illumina microarray identifier
experimentData: use 'experimentData(object)'
```

Annotation: lumiHumanV1
Control Data: Available
QC information: Please run summary(x, 'QC') for details!

5.2 Quality control of the raw data

The quality control of a **LumiBatch** object includes a data summary (the mean and standard deviation, sample correlation, detectable probe ratio of each sample (microarray)), different quality control plots, and the control probe information.

BeadStudio will usually separately output (or attached after the expressed data in the same file) the control probe (gene) information, usually named as "Control Probe Profile.txt". The controlData slot in LumiBatch class is designed to keep the control probe (gene) information. The control probe file can be inputted by using function `getControlData` or directly add it to a LumiBatch object by using function `addControlData2lumi`. Several functions `plotControlData`, `plotHousekeepingGene` and `plotStringencyGene` are designed to plot control probe data. Please see their help files for more details.

`LumiQ` function will produce the data summary of a **LumiBatch** object and organize the results in a QC slot of **LumiBatch** object. When creating the **LumiBatch** object, the `LumiQ` function will be called to initialize the QC slot of the **LumiBatch** object.

Summary of the quality control information of example.lumi data. If the QC slot of the **LumiBatch** object is empty, function `lumiQ` will be automatically called to estimate the quality control information.

```
> ## summary of the quality control
> summary(example.lumi, 'QC')
```

Data dimension: 8000 genes x 4 samples

Summary of Samples:

	A01	A02	B01	B02
mean	8.3240	8.568	8.2580	8.3470
standard deviation	1.5580	1.686	1.7230	1.6690
detection rate(0.01)	0.5432	0.564	0.5774	0.5758
distance to sample mean	76.9500	65.280	88.3200	49.1100

Major Operation History:

	submitted	finished
1	2007-04-22 00:08:36	2007-04-22 00:10:36
2	2007-04-22 00:10:36	2007-04-22 00:10:38

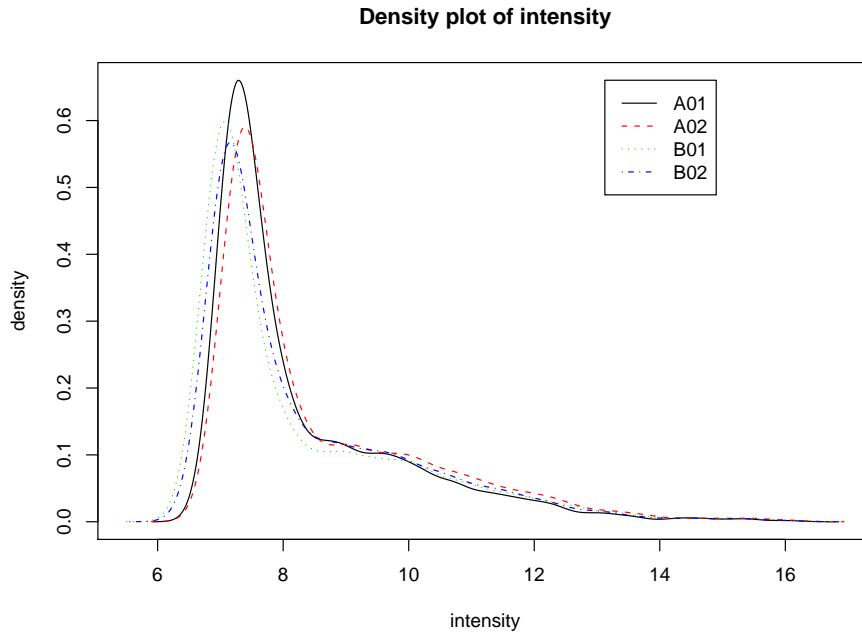


Figure 3: Density plot of Illumina microarrays before normalization

```

                                command lumiVersion
1 lumiR("../data/Barnes_gene_profile.txt")      1.1.6
2                               lumiQ(x.lumi = x.lumi) 1.1.6

```

The S4 method `plot` can produce the quality control plots of **LumiBatch** object. The quality control plots includes: the density plot (Figure 3), box plot (Figure 4), pairwise scatter plot between microarrays (Figure 5) or pair scatter plot with smoothing (Figure 6), pairwise MApot between microarrays (Figure 7) or MApot with smoothing (Figure 8), density plot of coefficient of variance, (Figure 9), and the sample relations (Figure 10). More details are in the help of `plot,LumiBatch-method` function. Most of these plots can also be plotted by the extended general functions: `density` (for density plot), `boxplot`, `MAplot`, `pairs` and `plotSampleRelation`.

Figure 3 shows the density plot of the **LumiBatch** object by using `plot` or `density` functions.

```

> ## plot the density
> plot(example.lumi, what='density')
> ## or
> density(example.lumi)

```

Figure 4 shows the box plot of the **LumiBatch** object by using `plot` or `boxplot` functions.

```

> ## plot the box plot
> plot(example.lumi, what='boxplot')

```

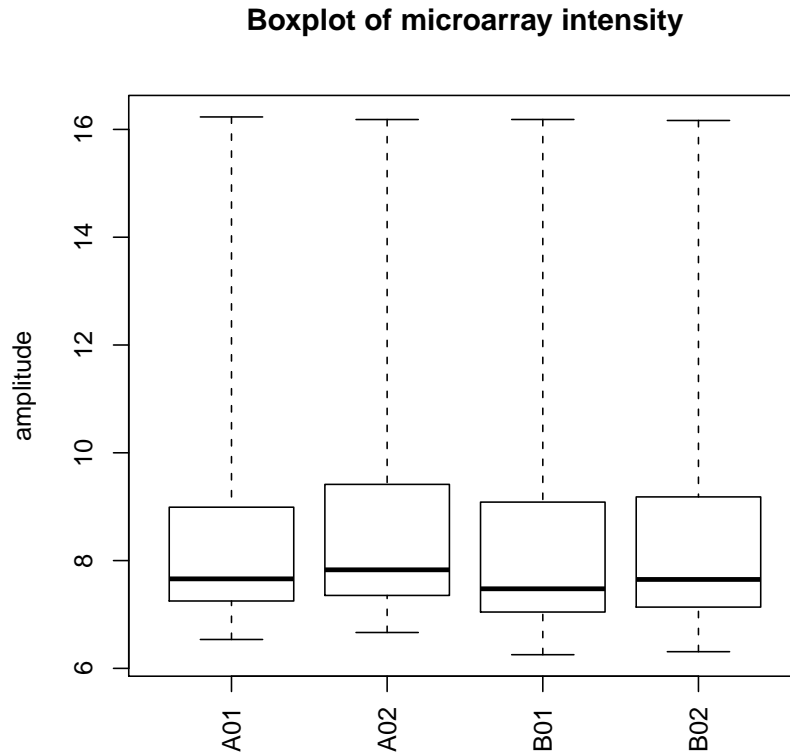


Figure 4: Density plot of Illumina microarrays before normalization

```
> ## or
> boxplot(example.lumi)
```

Figure 5 shows the pairwise sample correlation of the **LumiBatch** object by using `plot` or `pairs` functions.

```
> ## plot the pair plot
> plot(example.lumi, what='pair')
> ## or
> pairs(example.lumi)
> ## pairwise scatter plot with smoothing
> pairs(example.lumi, smoothScatter=T)
```

Figure 7 shows the MA plot of the **LumiBatch** object by using `plot` or `MAplot` functions.

```
> ## plot the MAplot
> plot(example.lumi, what='MAplot')
> ## or
```

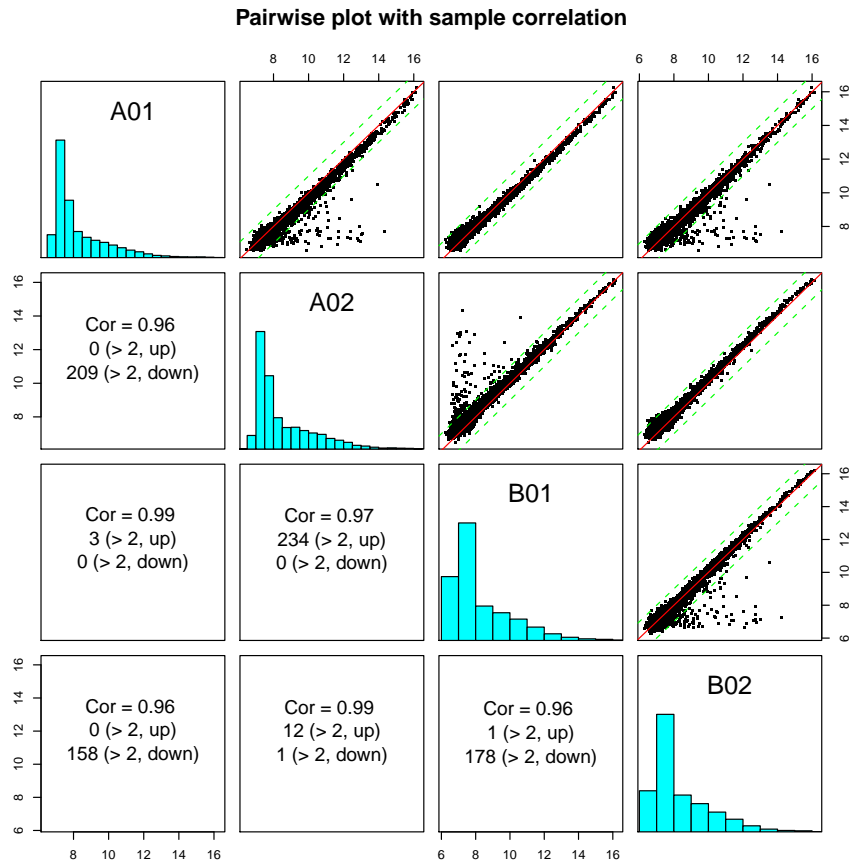


Figure 5: Pairwise plot with microarray correlation before normalization

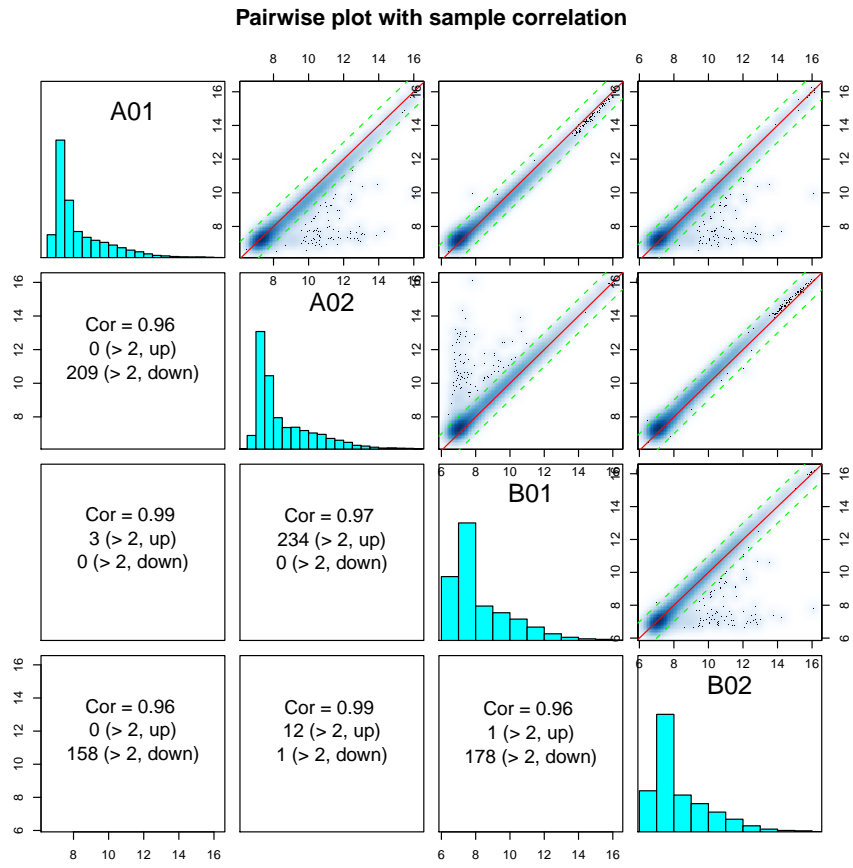


Figure 6: Pairwise plot with microarray correlation before normalization

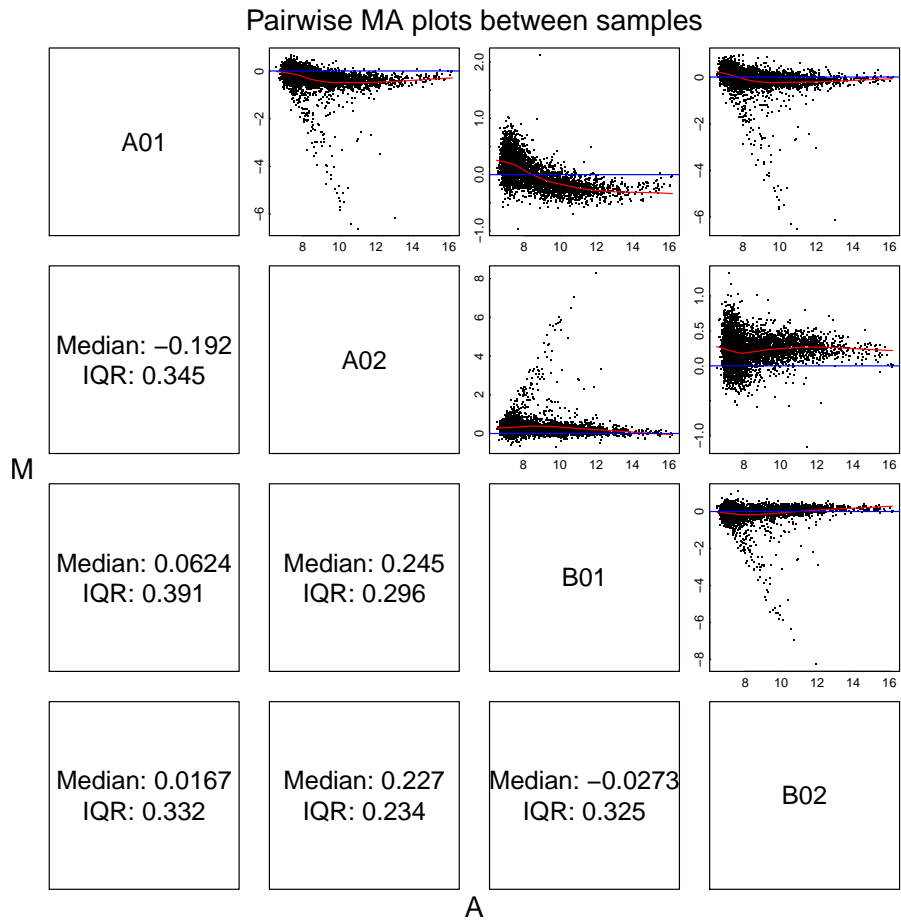


Figure 7: Pairwise MAplot before normalization

```
> MAplot(example.lumi)
> ## with smoothing
> MAplot(example.lumi, smoothScatter=TRUE)
```

The density plot of the coefficient of variance of the **LumiBatch** object. See Figure 9. Figure 9 shows the density plot of the coefficient of variance of the **LumiBatch** object by using `plot` function.

Figure 10 shows the sample relations using hierarchical clustering.

Figure 11 shows the sampleRelation using MDS. The color of the sample is based on the sample type, which is "01", "02", "01", "02" for the sample data. Please see the help of `plotSampleRelation` and `plot-methods` for more details.

```
> ## plot the sample relations
> plot(example.lumi, what='sampleRelation', method='mds', color=c("01", "02", "01", "02"))
> ## or
> plotSampleRelation(example.lumi, method='mds', color=c("01", "02", "01", "02"))
```

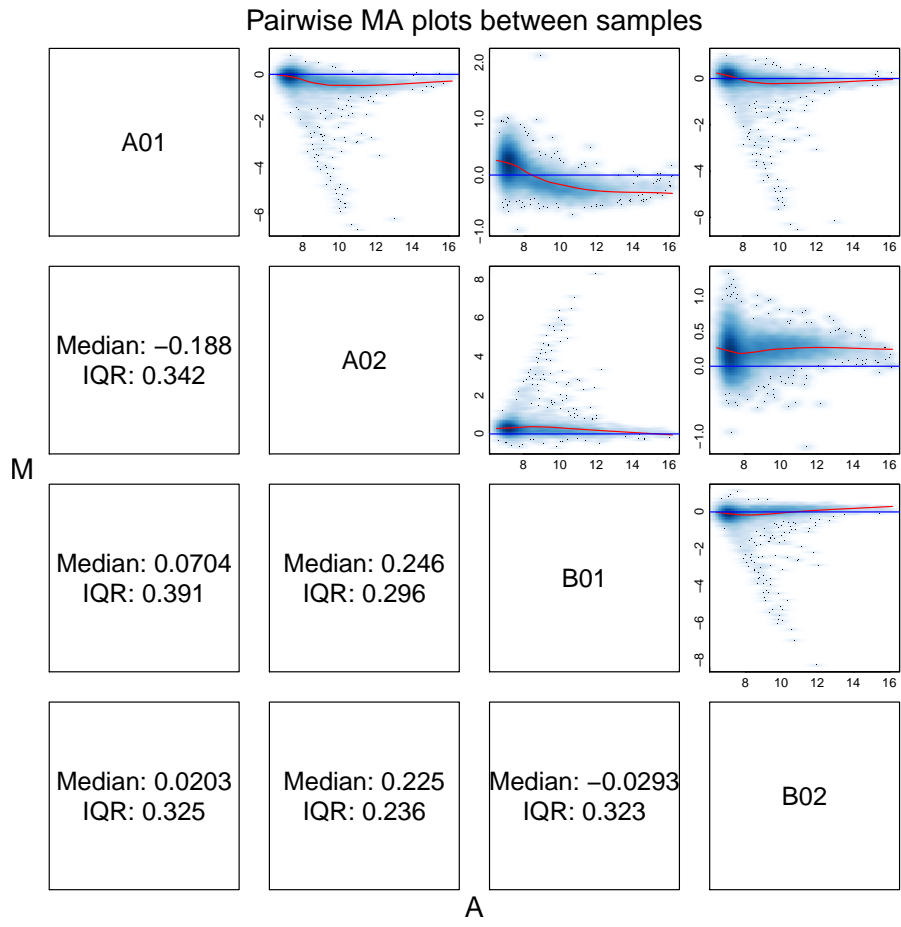


Figure 8: Pairwise MAplot with smoothing before normalization

```
> ## density plot of coefficient of variance  
> plot(example.lumi, what='cv')
```

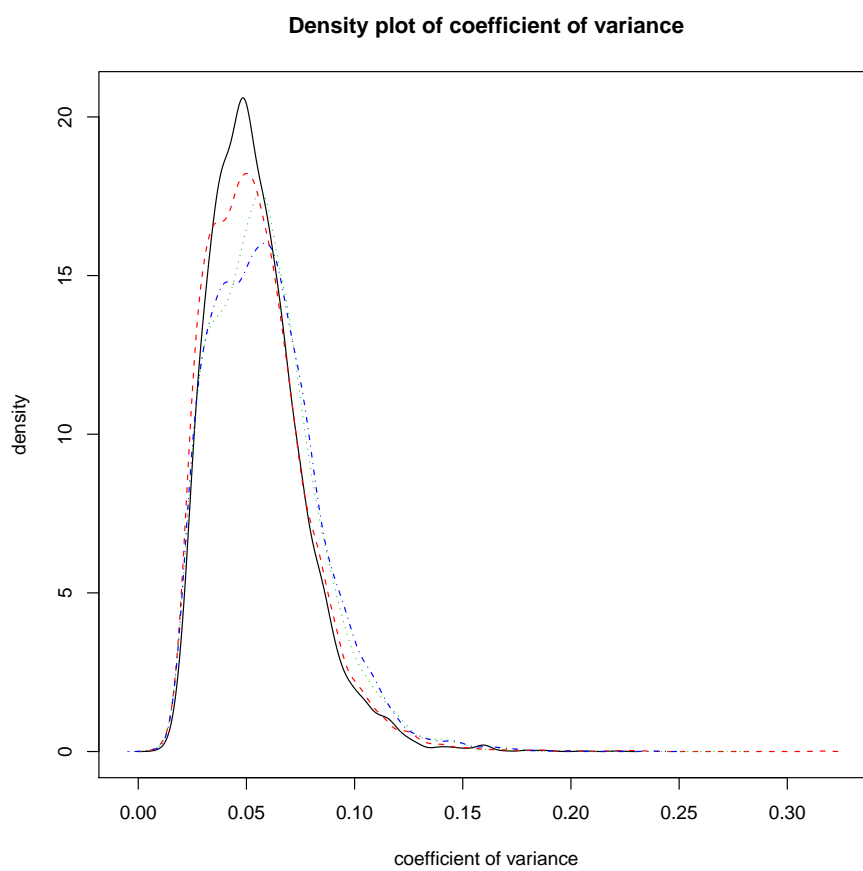


Figure 9: Density Plot of Coefficient of Variance

```
> plot(example.lumi, what='sampleRelation')
```

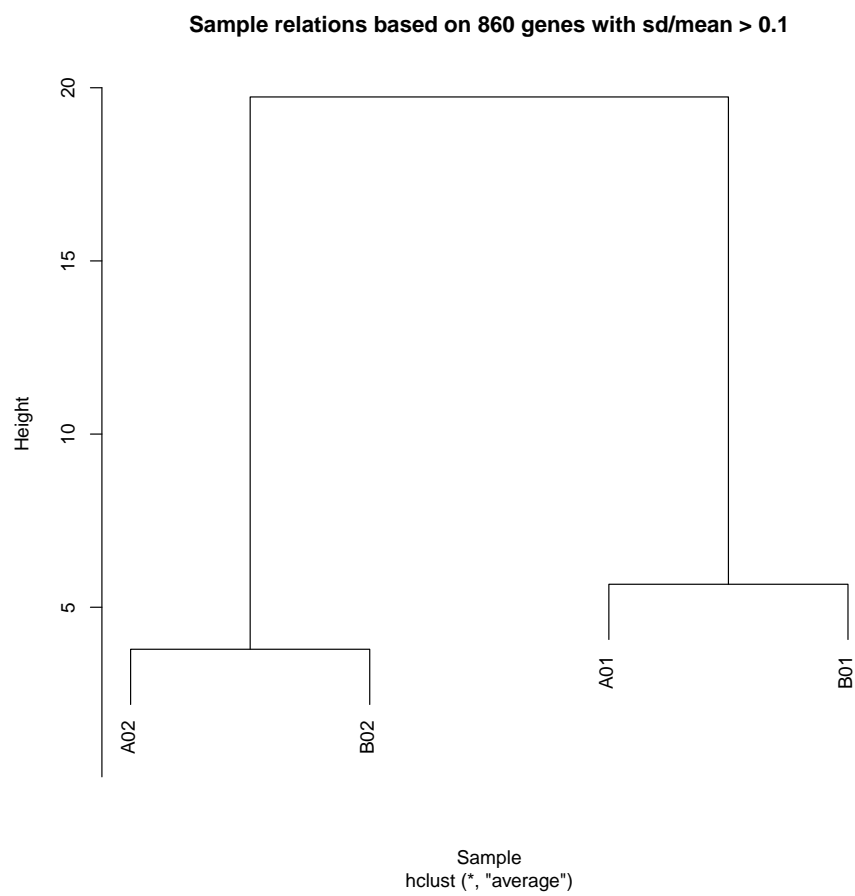


Figure 10: Sample relations before normalization

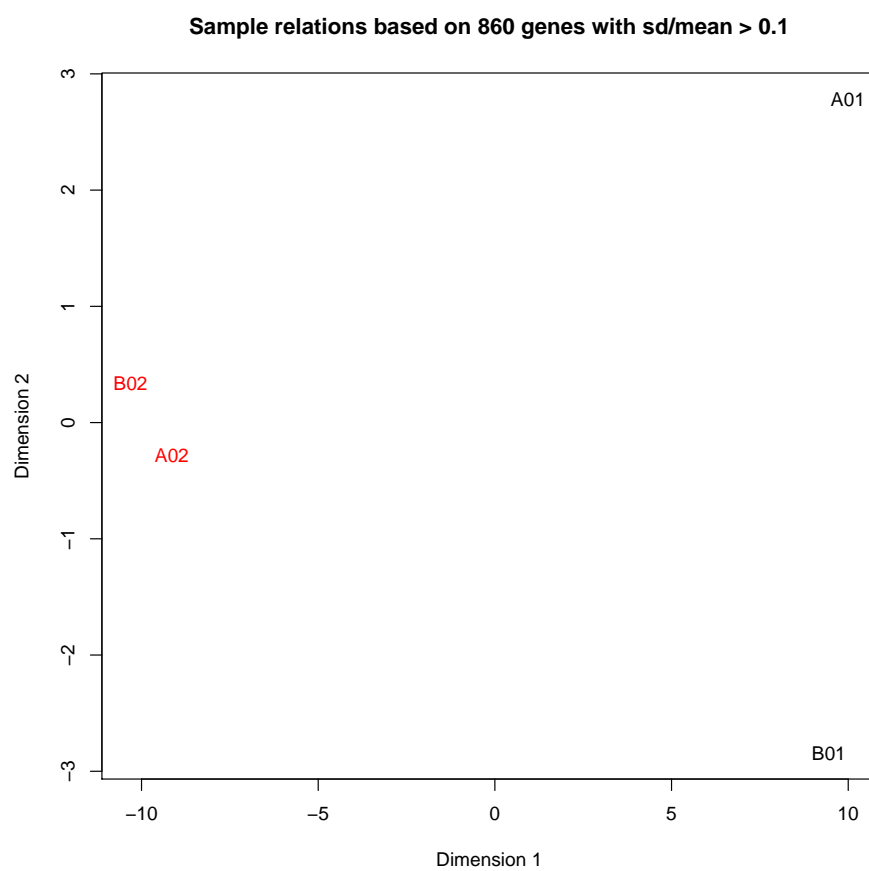


Figure 11: Sample relations before normalization

5.3 Background correction

The *lumi* package provides `lumiB` function for background correction. We suppose the BeadStudio output data has been background corrected. Therefore, no background correction used by default. A method 'bgAdjust' is designed to approximate what BeadStudio does for background adjustment. In the case when 'log2' transform is used in the `lumiT` step, the background correction method ('forcePositive') will be automatically used, which basically adds an offset (minus minimum value plus one) if there is any negative values to force all expression values to be positive. If users are more interested in the low level background correction, please refer to the package *beadarray* for more details. Users can also provide their own background correction function with a `LumiBatch` Object as the first argument and return a `LumiBatch` Object with background corrected. See `lumiB` help document for more details.

5.4 Variance stabilizing transform

Variance stabilization is critical for subsequent statistical inference to identify differential genes from microarray data. We devised a variance-stabilizing transformation (VST) by taking advantages of larger number of technical replicates available on the Illumina microarray. Please see [2] for details of the algorithm.

Because the `STDEV` (or `STDERR`) columns of the BeadStudio output file is the standard error of the mean of the bead intensities corresponding to the same probe. (Thanks Gordon Smyth kindly provided this information!). As the variance stabilization (see help of `vst` function) requires the information of the standard deviation instead of the standard error of the mean, the value correction is required. The corrected value will be $x * \sqrt{N}$, where x is the old value (standard error of the mean), N is the number of beads corresponding to the probe. The parameter 'stdCorrection' of `lumiT` determines whether to do this conversion and is effective only when the 'vst' method is selected. By default, the parameter 'stdCorrection' is `TRUE`.

Function `lumiT` performs variance stabilizing transform with both input and output being `LumiBatch` object.

Do default VST variance stabilizing transform

```
> ## Do default VST variance stabilizing transform
> lumi.T <- lumiT(example.lumi)
```

```
Perform vst transformation ...
2010-03-03 02:30:11 , processing array 1
2010-03-03 02:30:11 , processing array 2
2010-03-03 02:30:11 , processing array 3
2010-03-03 02:30:11 , processing array 4
```

The `plotVST` can plot the transformation function of VST, see Figure 12, which is close to log2 at high expression values, see Figure 13. Function `lumiT` also provides options to do "log2" or "cubicRoot" transform. See help of `lumiT` for details.

```
> ## plot VST transformation
> trans <- plotVST(lumi.T)
```

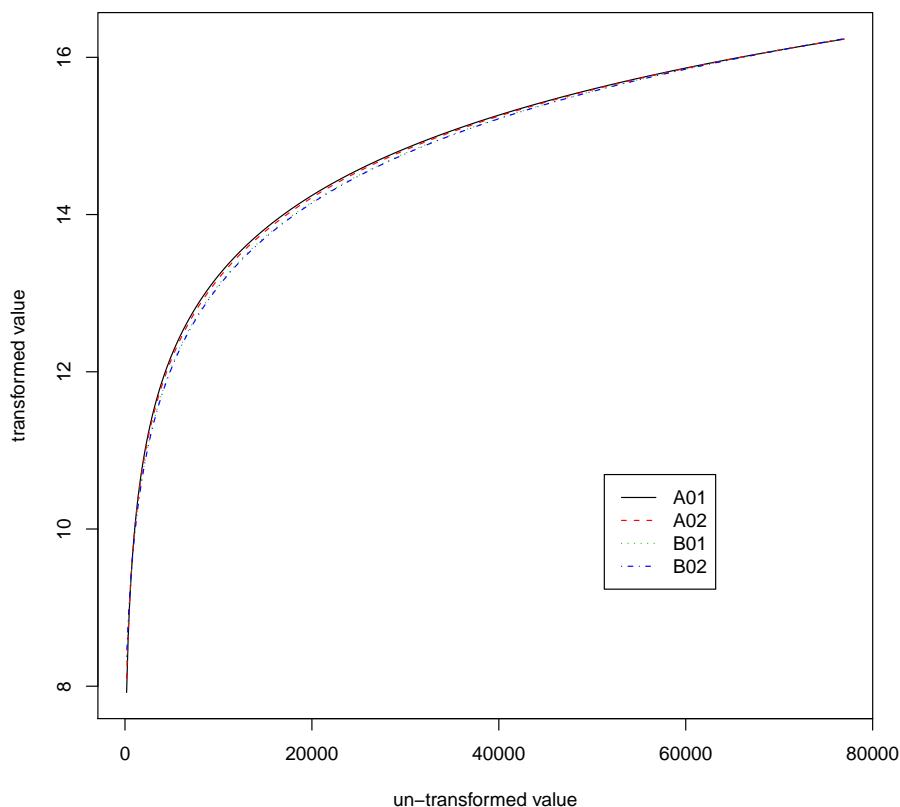


Figure 12: VST transformation

```
> ## compare the log2 and VST transform
> matplot(log2(trans$untransformed), trans$transformed, main='compare VST and log2 transfo
```

5.5 Data normalization

lumi package provides several normalization method options, which include quantile, SSN (Simple Scaling Normalization), RSN (Robust Spline Normalization), loess normalization and Rank Invariant Normalization.

Comparing with other normalization methods, like quantile and curve-fitting methods, SSN is a more conservative method. The only assumption is that each sample has the same background level and the same scale (if do scaling). It basically make all the samples have the same background level and the same scale comparing to the background (if do scaling). There are three methods ('density', 'mean' and 'median') for background estimation. If `bgMethod` is 'none', then no background adjustment. For the 'density' `bgMethod`, it estimates the background based on the mode of probe intensities based on the assumption that the background level intensity is the most frequent value across all the probes

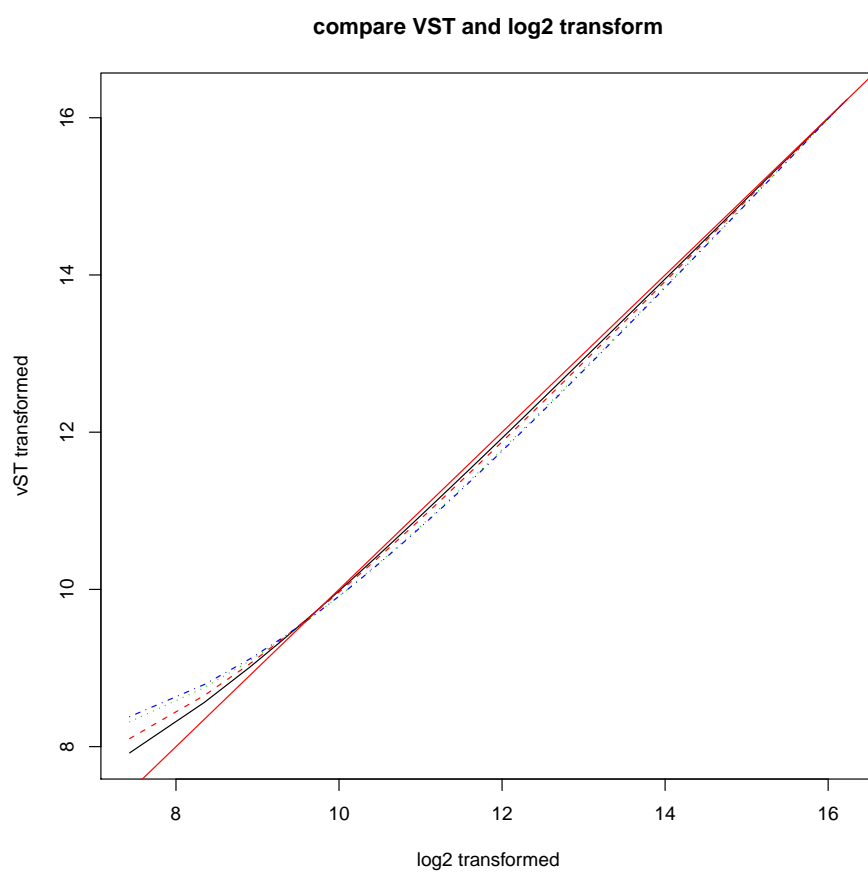


Figure 13: Compare VST and log2 transform

in the chip. For the foreground level estimation, it also provides three methods ('mean', 'density', 'median'). For the 'density' fgMethod, it assumes the background probe levels are symmetrically distributed. The foreground levels were estimated by taking the intensity mean of all other probes except from the background probes. For the 'mean' and 'median' methods (for both bgMethod and fgMethod), it basically estimates the level based on the mean or median of all probes of the sample. If the fgMethod is the same as bgMethod (except 'density' method), no scaling will be performed.

Another normalization method which is unique in the *lumi* package is the Robust Spline Normalization (RSN) algorithm. RSN combines the features of quantile and loess normalization. The advantages of quantile normalization include computational efficiency and preserving the rank order of genes. However, the intensity transformation of a quantile normalization is discontinuous because the normalization forces the intensity values for different samples (microarrays) having exactly the same distribution. This can cause small differences among intensity values to be lost. In contrast, the loess or spline normalization provides a continuous transformation. However, these methods cannot ensure that the rank of the probes remain unchanged across samples. Moreover, the loess normalization assumes the majority of the genes measured by the probes are non-differentially expressed and their distribution is approximately symmetric, which may not be a good assumption. To address some of these concerns, we developed a Robust Spline Normalization (RSN) method, which combines features from loess and quantile normalization methods. We use a monotonic spline to calibrate one microarray to the reference microarray. To increase the robustness of the spline method, we down-weight the contributions of probes of putatively differentially expressed genes. The probe intensities that are from potentially differentially expressed genes are heuristically determined as follows: First, we run a quantile normalization. Next, we estimate the fold-change of a gene measured by a probe based on the quantile-normalized data. The weighting factor for a probe is calculated based on a Gaussian window function. More details will be shown in a separate manuscript.

The default normalization method used in the Illumina BeadStudio software is Rank Invariant Normalization. In order to support similar functionalities, the *lumi* package also provides a similar normalization implementation call "rankinvariant" (We thanks Arno Velds implemented this function.). Please check the help of `rankinvariant` for more details.

By default, function `lumiN` performs popular quantile normalization. `lumiN` also provides other options to do "rsn", "ssn", "loess", "vsn", "rankinvariant" normalization. See help of `lumiN` for details.

Do default quantile between microarray normalization

```
> ## Do quantile between microarray normalization
> lumi.N <- lumiN(lumi.T)
```

Perform quantile normalization ...

Users can also easily select other normalization method. For example, the following command will run RSN normalization.

5.6 Quality control after normalization

To make sure the data quality meets our requirement, we do a second round of quality control of normalized data with different QC plots. Compare the plots before and after normalization, we can clearly see the improvements.

```
> ## Do quality control estimation after normalization
> lumi.N.Q <- lumiQ(lumi.N)
```

Perform Quality Control assessment of the LumiBatch object ...

```
> ## summary of the quality control
> summary(lumi.N.Q, 'QC')          ## summary of QC
```

Data dimension: 8000 genes x 4 samples

Summary of Samples:

	A01	A02	B01	B02
mean	8.8940	8.894	8.8940	8.8940
standard deviation	1.3000	1.300	1.3000	1.3000
detection rate(0.01)	0.5432	0.564	0.5774	0.5758
distance to sample mean	14.8400	14.600	14.8300	14.9600

Major Operation History:

	submitted	finished	command	lumiVersion
1	2007-04-22 00:08:36	2007-04-22 00:10:36	lumiR("../data/Barnes_gene_profile.txt")	1.1.6
2	2007-04-22 00:10:36	2007-04-22 00:10:38	lumiQ(x.lumi = x.lumi)	1.1.6
3	2007-04-22 00:13:06	2007-04-22 00:13:10	addNuId2lumi(x.lumi = x.lumi, lib = "lumiHumanV1")	1.1.6
4	2007-04-22 00:59:20	2007-04-22 00:59:36	Subsetting 8000 features and 4 samples.	1.1.6
5	2010-03-03 02:30:11	2010-03-03 02:30:11	lumiT(x.lumi = example.lumi)	1.12.4
6	2010-03-03 02:30:11	2010-03-03 02:30:12	lumiN(x.lumi = lumi.T)	1.12.4
7	2010-03-03 02:30:12	2010-03-03 02:30:12	lumiQ(x.lumi = lumi.N)	1.12.4

5.7 Encapsulate the processing steps

The `lumiExpresso` function is to encapsulate the major functions of Illumina preprocessing. It is organized in a similar way as the `expresso` function in *affy* package. The following code basically did the same processing as the previous multi-steps and produced the same results `lumi.N.Q`.

```
> ## Do all the default preprocessing in one step
> lumi.N.Q <- lumiExpresso(example.lumi)
```

```
> plot(lumi.N.Q, what='density') ## plot the density
```

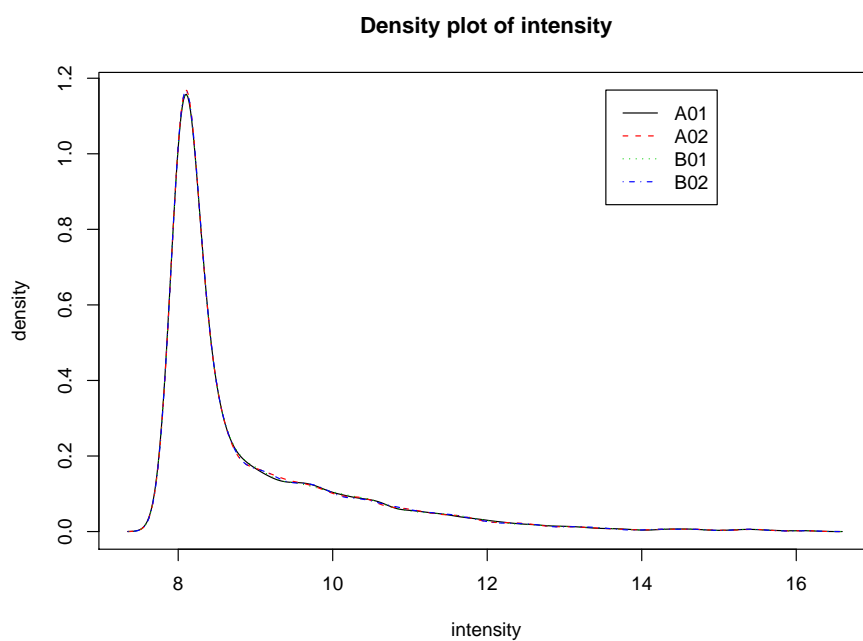


Figure 14: Density plot of Illumina microarrays after normalization

```
> plot(lumi.N.Q, what='boxplot') ## box plot
```

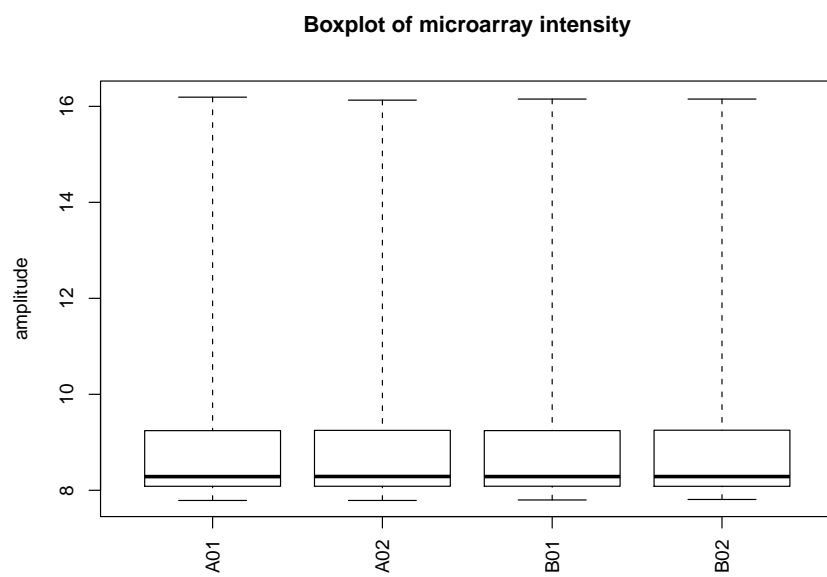


Figure 15: Density plot of Illumina microarrays after normalization

```
> plot(lumi.N.Q, what='pair')
```

```
## pairwise plots
```

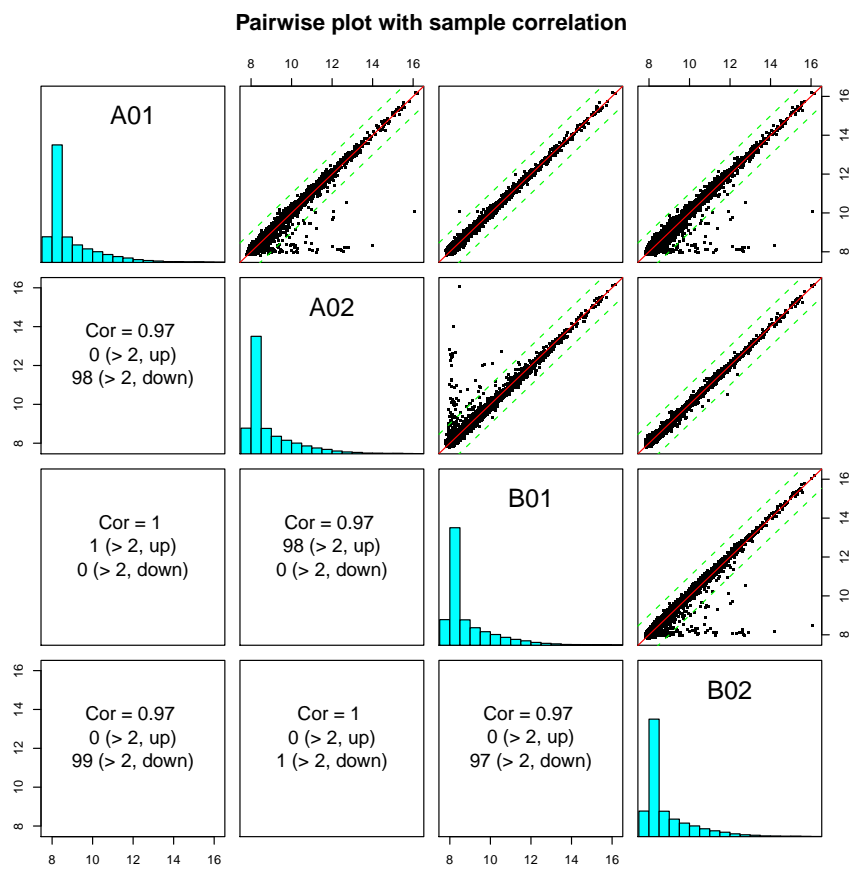


Figure 16: Pairwise plot with microarray correlation after normalization

```
> plot(lumi.N.Q, what='MAplot') ## plot the pairwise MAplot
```

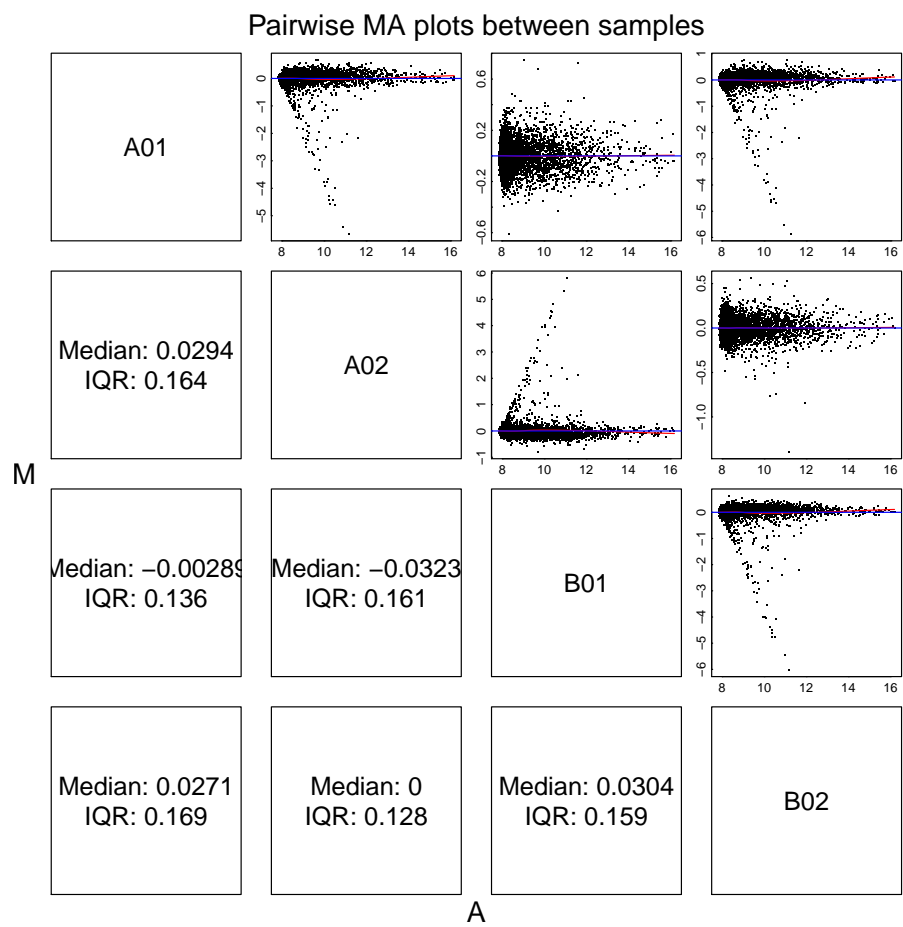


Figure 17: Pairwise MAplot after normalization

```
> ## plot the sampleRelation using hierarchical clustering  
> plot(lumi.N.Q, what='sampleRelation')
```

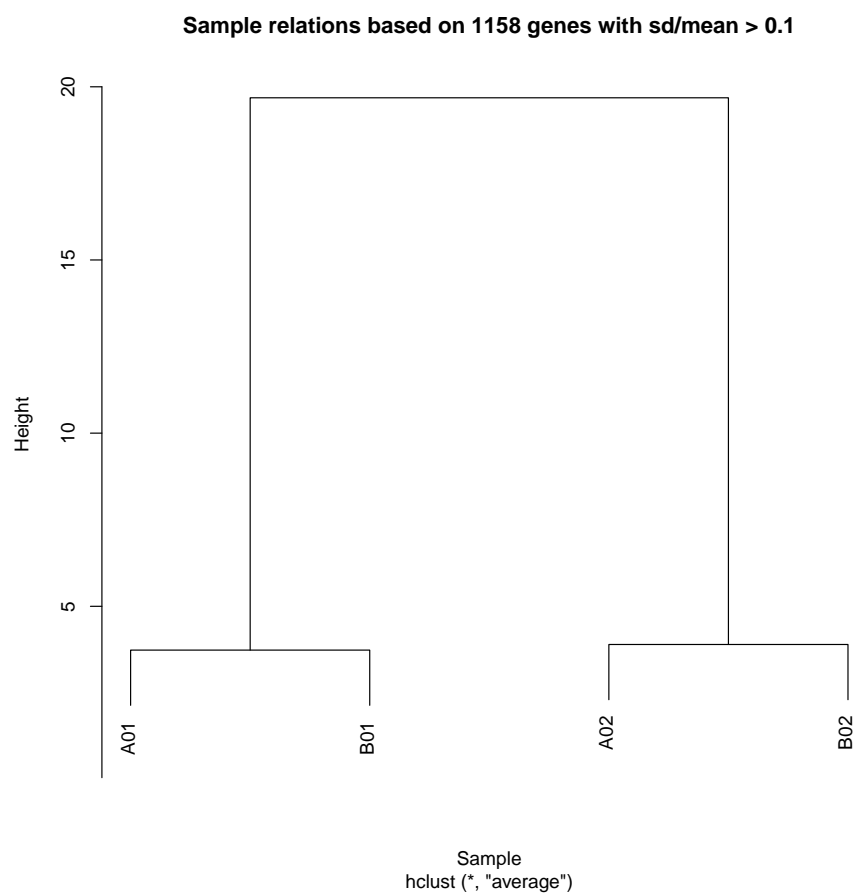


Figure 18: Sample relations after normalization

```
> ## plot the sampleRelation using MDS  
> plot(lumi.N.Q, what='sampleRelation', method='mds', color=c("01", "02", "01", "02"))
```

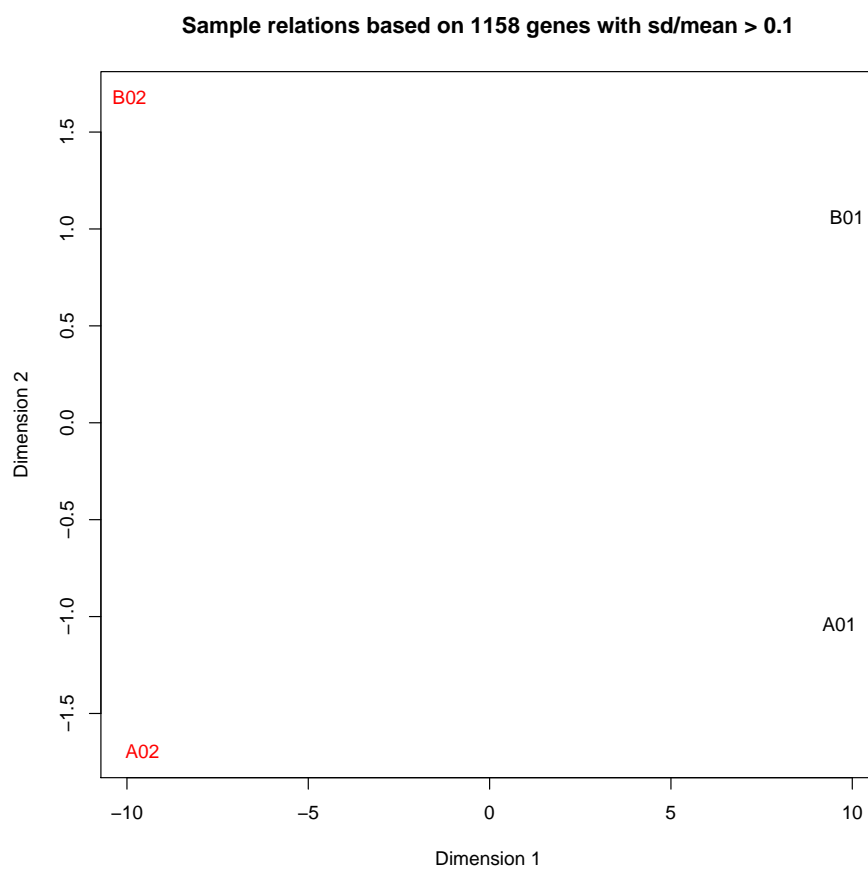


Figure 19: Sample relations after normalization


```
Background Correction: bgAdjust
Variance Stabilizing Transform method: vst
Normalization method: quantile
```

```
Background correction ...
Perform bgAdjust background correction ...
done.
```

```
Variance stabilizing ...
Perform vst transformation ...
2010-03-03 02:30:22 , processing array 1
2010-03-03 02:30:22 , processing array 2
2010-03-03 02:30:22 , processing array 3
2010-03-03 02:30:22 , processing array 4
done.
```

```
Normalizing ...
Perform quantile normalization ...
done.
```

```
Quality control after preprocessing ...
Perform Quality Control assessment of the LumiBatch object ...
done.
```

Users can easily customize the processing parameters. For example, if the user wants to do "rsn" normalization, the user can run the following code. For more details, please read the help document of `lumiExpresso` function.

5.8 Inverse VST transform to the raw scale

Figure 13 shows VST is very close to \log_2 in the high expression range. In convenience, users usually can directly use 2^x to approximate the data in raw scale and estimate the fold-change. For the users concern more in the low expression range, we also provide the function `inverseVST` to resume the data in the raw scale. Need to mention, the inverse transform should be performed after statistical analysis, or else it makes no sense to transform back and forth. The `inverseVST` function can directly applied to the **LumiBatch** object after `lumiT` with VST transform, or VST transform plus RSN normalization (default method of `lumiN`). For the RSN normalized data, the inverse transform is based on the parameters of the Target Array because the Target Array is the benchmark data and is not changed after normalization. Other normalization methods, like quantile or loess, will change the values of all the arrays. As a result, no inverse VST transform available for them. Users may use some kind of approximation for the quantile normalized data by themselves. Here we just provide some examples of VST parameters retrieving and inverse VST transform.

```
> ## Parameters of VST transformed LumiBatch object
> names(attributes(lumi.T))
```

```

[1] "history"          "controlData"      "QC"
[4] "assayData"        "phenoData"        "featureData"
[7] "experimentData"   "annotation"       "protocolData"
[10] ".__classVersion__" "class"            "vstParameter"
[13] "transformFun"

> ## VST parameters: "vstParameter" and "transformFun"
> attr(lumi.T, 'vstParameter')

          a          b          g Intercept
A01 2.396259 0.02244804 1.480568 4.166654
A02 3.574381 0.02063079 1.505113 4.094366
B01 6.513429 0.02172944 1.554504 3.614302
B02 6.878816 0.02030299 1.566239 3.626431

> attr(lumi.T, 'transformFun')

      A01      A02      B01      B02
"asinh" "asinh" "asinh" "asinh"

> ## Parameters of VST transformed and RSN normalized LumiBatch object
> names(attributes(lumi.N.Q))

[1] "history"          "controlData"      "QC"
[4] "assayData"        "phenoData"        "featureData"
[7] "experimentData"   "annotation"       "protocolData"
[10] ".__classVersion__" "class"            "vstParameter"
[13] "transformFun"

> ## VSN "targetArray" , VST parameters: "vstParameter" and "transformFun"
> attr(lumi.N.Q, 'vstParameter')

          a          b          g Intercept
7.20873386 0.02136944 1.53480676 3.80340471

> attr(lumi.N.Q, 'transformFun')

[1] "asinh"

> ## After doing statistical analysis of the data, users can recover to the raw scale for
> ## Inverse VST to the raw scale
> lumi.N.raw <- inverseVST(lumi.N.Q)

```

6 Handling large data sets

Several users asked about processing large data set, e.g., over 100 samples. Directly handling such big data set usually will cause "out of memory" error in most computers. In this case, when read the BeadStudio output file, we can ignore the "beadNum" (related columns. The function `lumiR` provides a parameter called "columnNameGrepPattern". we can set the string grep pattern of "detection" and "beadNum" as NA. You can also ignore "detection" columns. However, the "detection" information is useful for the estimation of present count of each

probe and used in the VST parameter estimation. To further save memory, you can suppress the input of annotation data by setting "inputAnnotation" as FALSE.

Here is some example code:

```
## load the data with empty detection and beadNum slots, and without annotation informatio
> x.lumi <- lumiR("fileName.txt", columnNameGrepPattern=list(beadNum=NA), inputAnnotation
```

Usually, the large data set is composed of many small data files. In this case, the transformations, like log2 and vst, can be performed right after the input of each data file and some information can be removed in the object after transformation. *lumi* provides the `lumiR.batch` function for this purpose.

Here is some example code:

```
## load the list of data files (a vector of file names)
##           and do VST transformation for each file and combine the results.
> x.lumi <- lumiR.batch(fileList, transform='vst')
```

Another good news is that the normalization, like `rsn` and `ssn` in the *lumi* package, can sequentially process the data and handle such large data set.

The solution can be like this:

1. Read the data file by smaller batches (e.g. 10 or just one by one), and then do the variance stabilization for each data batch using `lumiR.batch` or `lumiR` function.
2. Pick one sample as the target array for normalization and then using "RSN" or "SSN" normalization method to normalize all batches of data using the same target array.
3. Combine the normalized data. (In order to save memory, the user can first remove those probes not expressed in all samples.)

In the `rsn` and `ssn` functions, there is a parameter called "targetArray", which is the model for other chips to normalize. It can be a column index, a vector or a `LumiBatch` object with one sample. In our case, we need to use one `LumiBatch` object with one sample as the "targetArray". The selection of the target array is flexible. We suggest to choose the one most similar to the mean of all samples. For convenience, we can also just select the first sample as "targetArray" (suppose it has no quality problem). The selected target array will also be used for all other data batches. Since different data batches use the same target array as model, the results are comparable and can be combined!

Here is the example code:

```
## Read in the Batch ith data file, suppose named as "fileName.i.txt"
> x.lumi.i <- lumiR("fileName.i.txt")
## variance stabilization (using vst or log2 transform)
> x.lumiT.i <- lumiT(x.lumi.i)
## select the "targetArray"
## This target array will also be used for other batches of data.
## For convenience, here we just select the first sample as targetArray.
> targetArray <- x.lumiT.i[,1]
## Do RSN normalization
> x.lumiN.i <- lumiN(x.lumiT.i, method='rsn', targetArray=targetArray)
```

The normalized data batches can be combined by using function `combine(x, y)`.

7 Performance comparison

We have selected the Barnes data set [4], which is a series dilution of two tissues at five different dilutions, to compare different preprocessing methods. In order to better compare the algorithms, we selected the samples with the smallest dilution difference (the most challenging comparison), i.e., the samples with the dilution ratios of 100:0 and 95:5 (each condition has two technical replicates) for comparison. For the Barnes data set, because we do not know which of the signals are coming from 'true' differentially expressed genes, we cannot use an ROC curve to compare the performance of different algorithms. Instead, we evaluated the methods based on the concordance of normalized intensity profile and real dilution profile of the selected probes. More detailed evaluations with other criteria and based on other data sets can be found in our paper [2].

Following Barnes et al. (2005)[4], we defined a concordant gene (really a concordant probe) as a signal from a probe with a correlation coefficient larger than 0.8 between the normalized intensity profile and the real dilution profile (five dilution ratios with two replicates at each dilution ratio). If a selected differentially expressed probe is also a concordant one, it is more likely to be truly differentially expressed. Figure 20 shows the percentage of concordant probes among the selected probes, which were selected by ranking the probes' p-value (calculated based on *limma* package) from low to high. We can see the VST transformed data outperforms the Log2-transformed and VSN processed data. For the normalization methods, RSN and quantile normalization have similar performance for the VST transformed data, and RSN outperforms quantile for the Log transformed data.

Please see another vignette in the lumi package: "[lumi_vST_evaluation.pdf](#)" for more details of the evaluation of VST (Variance Stabilizing Transformation).

8 Gene annotation

One challenge of Illumina microarray is the inconsistency and changes of Illumina identifiers across versions, even across different releases. This makes the integration of the Illumina data difficult. In order to resolve these problems, we invented a nuID (nucleotide universal IDentifier) annotation system, released related annotation packages and a website to provide identifier mapping and the latest annotation. Please refer to the separate document ("Resolve the Inconsistency of Illumina Identifiers through nuID Annotation") in the lumi package for more details.

9 A use case: from raw data to functional analysis

Figure 21 shows the data processing flow chart of the use case. Since the classes in *lumi* package are inherited from class **ExpressionSet**, packages and functions compatible with class **ExpressionSet** or accepting matrix as input all can be used for *lumi* results. Here we just give two examples: using *limma* to identify differentiated genes and using *GStats* to annotate the significant genes.

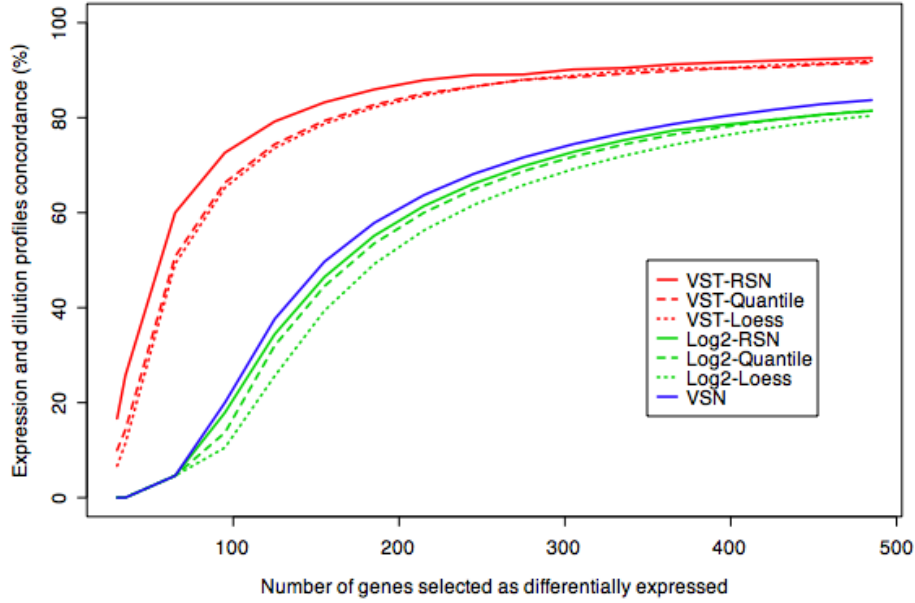


Figure 20: Comparison of the concordance between the expression and dilution profiles of the selected differentially expressed genes

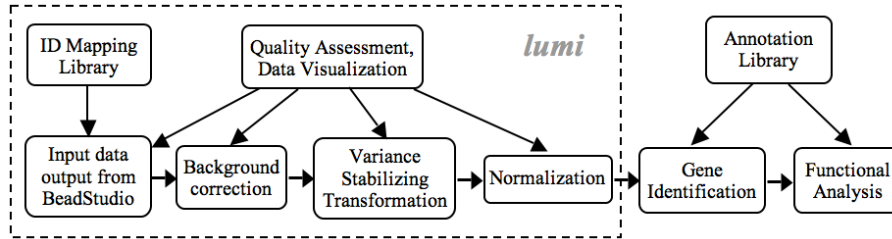


Figure 21: Flow chart of the use case

We use the Barnes data set [4] as an example, which has been created as a Bioconductor experiment data package *lumiBarnes*. The Barnes data set measured a dilution series of two human tissues, blood and placenta. It includes six samples with the titration ratio of blood and placenta as 100:0, 95:5, 75:25, 50:50, 25:75 and 0:100. The samples were hybridized on HumanRef-8 BeadChip (Illumina, Inc) in duplicate. We select samples with titration ratio, 100:0 and 95:5 (each has two technique replicates) in this data set to evaluate the detection of differential expressions.

9.1 Preprocess the Illumina data

```

> library(lumi)
> ## specify the file name
> # fileName <- 'Barnes_gene_profile.txt' # Not run

```

```

> ## load the data if you want to provide the sample Information and map Illumina IDs to n
> # example.lumi <- lumiR.batch(fileName, lib.mapping='lumiHumanIDMapping', sampleInfoFile

> ## load saved data
> data(example.lumi)
> ## sumary of the daa
> example.lumi
> ## summary of quality control information
> summary(example.lumi, 'QC')

> ## preprocessing and quality control after normalization
> lumi.N.Q <- lumiExpresso(example.lumi, QC.evaluation=TRUE)
> ## summary of quality control information after preprocessing
> summary(lumi.N.Q, 'QC')

> ## Output the data as Tab separated text file
> write.exprs(lumi.N.Q, file='processedExampleData.txt')

```

9.2 Identify differentially expressed genes

Identify the differentiated genes based on moderated t-test using *limma*.

Retrieve the normalized data

```

> dataMatrix <- exprs(lumi.N.Q)

```

To speed up the processing and reduce false positives, remove the unexpressed genes

```

> presentCount <- detectionCall(example.lumi)
> selDataMatrix <- dataMatrix[presentCount > 0,]
> probeList <- rownames(selDataMatrix)

> ## Specify the sample type
> sampleType <- c('100:0', '95:5', '100:0', '95:5')
> if (require(limma)) {
+     ## compare '95:5' and '100:0'
+     design <- model.matrix(~ factor(sampleType))
+     colnames(design) <- c('100:0', '95:5-100:0')
+     fit <- lmFit(selDataMatrix, design)
+     fit <- eBayes(fit)
+     ## Add gene symbols to gene properties
+     if (require(lumiHumanAll.db) & require(annotate)) {
+         geneSymbol <- getSYMBOL(probeList, 'lumiHumanAll.db')
+         geneName <- sapply(lookup(probeList, 'lumiHumanAll.db', 'GENENAME'), func
+         fit$genes <- data.frame(ID= probeList, geneSymbol=geneSymbol, geneName=ge
+     }
+     ## print the top 10 genes
+     print(topTable(fit, coef='95:5-100:0', adjust='fdr', number=10))
+
+     ## get significant gene list with FDR adjusted p.values less than 0.01

```

```

+       p.adj <- p.adjust(fit$p.value[,2])
+       sigGene.adj <- probeList[ p.adj < 0.01]
+       ## without FDR adjustment
+       sigGene <- probeList[ fit$p.value[,2] < 0.01]
+ }

```

	ID	geneSymbol
1116	ol_iQkR.siio.kvH6k	PLAC4
3080	EY761AIG0XSLUfnuyc	CGA
3772	WlCoF7taz2MeYf3l6I	SDC1
47	NSjRKdq2eSGf0ur4aQ	PRG2
1401	6QNTThLQLd61eU6IXhI	PSG9
2520	QaYYojcJJvVELV3I98	DLK1
1027	uioiKiIlzFXx8k5EC4	CRH
3831	TueuSaiCheWBxB6B18	KISS1
3257	rSU1F9I7txuZ3lPQdo	DCN
3236	Q.oCSr13l5wQlRuhS0	FSTL1

1116		place
3080		glycoprotein hormones, alp
3772		
47	proteoglycan 2, bone marrow (natural killer cell activator, eosinophil granule major	
1401		pregnancy specific beta-1-
2520		delta-like 1 homolo
1027		corticotropin rel
3831		KiSS-1 metasta
3257		
3236		fol

	logFC	AveExpr	t	P.Value	adj.P.Val	B
1116	5.632889	10.797487	69.39455	5.602391e-16	2.939015e-12	26.87647
3080	6.046529	11.055906	63.94646	1.383393e-15	2.975232e-12	26.10723
3772	4.850967	10.171877	62.76042	1.701429e-15	2.975232e-12	25.92735
47	4.583911	10.326147	59.63395	2.992662e-15	3.144939e-12	25.42974
1401	4.563612	10.124415	59.62530	2.997463e-15	3.144939e-12	25.42832
2520	4.303575	10.193919	57.10027	4.834715e-15	3.860305e-12	24.99966
1027	4.461162	10.035413	56.60483	5.323258e-15	3.860305e-12	24.91256
3831	4.705771	10.168084	56.09144	5.886855e-15	3.860305e-12	24.82124
3257	4.210354	9.893562	53.10355	1.077490e-14	6.280568e-12	24.26698
3236	4.190180	10.446315	52.39880	1.248742e-14	6.550900e-12	24.13031

Based on the significant genes identified using *limma* or t-test, we can do further analysis, like GO analysis (*GOstats* package) and machine learning (*MLInterface* package). Next, we will use GO analysis as an example.

9.3 Gene Ontology and other functional analysis

Based on the interested genes identified using *limma* or other tests, we can further do functional analysis. We can use package *GOstats*, *Category* and other packages to do this type of analysis.

There is one important thing need to mention during this type of analysis. As we described in previous section, the lumi annotation packages are nuID

indexed and are built by pooling all types of chips of the same species. This makes it different from the traditional Affymetrix annotation packages, which is one package for one type of chip. Because lots of methods in Category/GOstats were originally designed based on the Affymetrix annotation packages, the default setting of these function may not work well for lumi annotation packages. However, this can be solved by first transferring the probe ids as Entrez Ids, and then do analysis at the Entrez Id level instead of the probe Id level. Please see our example for how to transfer Probe Ids as Entrez Ids.

Following is an example of performing Hypergeometric test of Gene Ontology based on the significant gene list (for e. Table 1 shows the significant GO terms of Molecular Function with p-value less than 0.01. Here only show the significant GO terms of BP (Biological Process). For other GO categories MF(Molecular Function) and CC (Cellular Component), it just follows the same procedure.

```
> if (require(GOstats) & require(lumiHumanAll.db)) {
+
+     ## Convert the probe Ids as Entrez Ids and make them unique
+     sigLL <- unique(unlist(lookup(sigGene, 'lumiHumanAll.db', 'ENTREZID')))
+     sigLL <- as.character(sigLL[!is.na(sigLL)])
+     params <- new("GOHyperGParams",
+         geneIds= sigLL,
+         annotation="lumiHumanAll.db",
+         ontology="BP",
+         pvalueCutoff= 0.01,
+         conditional=FALSE,
+         testDirection="over")
+
+     hgOver <- hyperGTest(params)
+
+     ## Get the p-values of the test
+     gGhyp.pv <- pvalues(hgOver)
+
+     ## Adjust p-values for multiple test (FDR)
+     gGhyp.fdr <- p.adjust(gGhyp.pv, 'fdr')
+
+     ## select the Go terms with adjusted p-value less than 0.01
+     sigGO.ID <- names(gGhyp.fdr[gGhyp.fdr < 0.01])
+
+     ## Here only show the significant GO terms of BP (Molecular Function)
+     ##           For other categories, just follow the same procedure.
+     sigGO.Term <- getGOTerm(sigGO.ID)[["BP"]]
+ }
```

9.4 GEO submission of the data

As more and more journals require the microarray data should be submitted to GEO website, we created GEO submission functions for users convenience. The submission file will be in the SOFT format. So users can submit all the data in a batch. There are two major functions users need to use. `produceGEOSampleInfoTemplate` produces a template of sample information with some default

	GO ID	Term	p-value	Significant Genes No.	Total Genes No.
1	GO:0009611	response to wound...	8.4244e-06	42	443
2	GO:0006955	immune response	8.8296e-06	68	859
3	GO:0006952	defense response	1.7525e-05	72	945
4	GO:0006950	response to stres...	1.9132e-05	81	1103
5	GO:0009607	response to bioti...	5.0811e-05	72	976
6	GO:0009613	response to pest,...	7.2813e-05	45	533
7	GO:0006954	inflammatory resp...	0.00025402	25	250
8	GO:0009605	response to exter...	0.00026005	46	580
9	GO:0051707	response to other...	0.00040553	45	575
10	GO:0051674	localization of c...	0.00082563	30	348
11	GO:0006928	cell motility	0.00082563	30	348
12	GO:0040011	locomotion	0.00099205	30	352

Table 1: GO terms, p-values and counts.

fillings. Some fields have been filled in with default settings. Users should fill in or modify the detailed sample descriptions by referring some previous submissions. No blank fields are allowed. Users are also suggested to fill in the "Sample_platform_id" by checking information of the GEO Illumina platform. `produceGEOSubmissionFile` is the main function of produce GEO submission file including both normalized and raw data information in the SOFT format. By default, the R objects, `lumiNormalized`, `lumiRaw` and `sampleInfo`, will be saved in a file named 'supplementaryData.Rdata'. (See help information of `produceGEOSubmissionFile`) Users can include this R data file as a GEO supplementary data file.

```
## Produce the sample information template
```

```
> produceGEOSampleInfoTemplate(lumi.example, lib.mapping = 'lumiHumanIDMapping', fileName
```

```
## After editing the 'GEOSampleInfo.txt' by filling in sample information
```

```
> produceGEOSubmissionFile(lumi.N.Q, lumi.example, lib='lumiHumanIDMapping', sampleInfo='G
```

10 Session Info

```
> toLatex(sessionInfo())
```

- R version 2.10.1 (2009-12-14), i386-pc-mingw32
- Locale: LC_COLLATE=English_United States.1252,
LC_CTYPE=English_United States.1252,
LC_MONETARY=English_United States.1252, LC_NUMERIC=C,
LC_TIME=English_United States.1252
- Base packages: base, datasets, graphics, grDevices, methods, stats, tools,
utils
- Other packages: affy 1.24.2, annotate 1.24.1, AnnotationDbi 1.8.1,
Biobase 2.6.1, DBI 0.2-5, limma 3.2.2, lumi 1.12.4,
lumiHumanAll.db 1.8.1, lumiHumanIDMapping 1.4.0, MASS 7.3-5,
mgcv 1.6-1, org.Hs.eg.db 2.3.6, preprocessCore 1.8.0, RSQLite 0.8-3

- Loaded via a namespace (and not attached): affyio 1.14.0, grid 2.10.1, KernSmooth 2.23-3, lattice 0.18-3, Matrix 0.999375-37, nlme 3.1-96, xtable 1.5-6

11 Acknowledgments

We would like to thanks the users and researchers around the world contribute to the lumi package, provide great comments and suggestions and report bugs. Especially, we would like to thanks Michal Blazejczyk, Peter Bram, Ligia Bras, Vincent Carey, Kevin Coombes, Sean Davis, Jean-Eudes DAZARD, Ryan Gordon, Wolfgang Huber, DeokHoon Kim, Matthias Kohl, Danilo Licastro, Ezhou Lori Long, Renee McElhaney, Martin Morgan, Ingrid H. G. ístense, Denise Scholtens, Wei Shi, Gordon Smyth, Michael Stevens, Jiexin Zhang (sorted by last name) and many other people not mentioned here.

12 References

1. Du, P., Kibbe, W.A. and Lin, S.M., (2008) 'lumi: a pipeline for processing Illumina microarray', *Bioinformatics* 24(13):1547-1548
2. Lin, S.M., Du, P., Kibbe, W.A., (2008) 'Model-based Variance-stabilizing Transformation for Illumina Microarray Data', *Nucleic Acids Res.* 36, e11
3. Du, P., Kibbe, W.A. and Lin, S.M., (2007) 'nuID: A universal naming schema of oligonucleotides for Illumina, Affymetrix, and other microarrays', *Biology Direct*, 2, 16.
4. Barnes, M., Freudenberg, J., Thompson, S., Aronow, B. and Pavlidis, P. (2005) "Experimental comparison and cross-validation of the Affymetrix and Illumina gene expression analysis platforms", *Nucleic Acids Res*, 33, 5914-5923.