

HowTo BGX

Ernest Turro
Imperial College London

April 25, 2007

1 Introduction

This vignette describes how to use *bgx*, a C++ implementation of a Bayesian hierarchical integrated approach to the modelling and analysis of Affymetrix GeneChip arrays. The model and methodology is described in Hein et al, 2005.

There are two ways to run *bgx*: (1) through R and (2) as a standalone binary. Both ways make use of probe level GeneChip data, which you must obtain as GeneChip CEL files.

2 Reading in the CEL files

When you load *bgx*, several required packages from the Bioconductor¹ project are automatically loaded.

```
> library(bgx)
```

The *affy* package allows you to read CEL files into an `AffyBatch` object. This can be achieved by changing your working directory to wherever the CEL files are stored and executing:

```
> aData <- ReadAffy()
```

This will read in the CEL files in alphabetical order and save the data in the `aData` object. Alternatively, you can specify the specific files you would like to read in by adding their paths to the argument list, for example:

```
> aData <- ReadAffy("CEL/choe/chipC-rep1.CEL", "CEL/choe/chipS-rep2.CEL")
```

¹<http://bioconductor.org>

3 Running BGX through R

A basic execution of the program can be performed by simply passing an `AffyBatch` object as a single parameter to the `bgx` function and saving the result in an `ExpressionSet` object. The result will hold array-specific gene expression values and their corresponding standard errors in `assayData(eset)$exprs` and `assayData(eset)$se.exprs` respectively.

```
> eset <- bgx(aData)
```

A more elaborate scenario would involve splitting the arrays into a number of conditions using the *samplesets* argument²; specifying which genes to analyse with the *genes* argument; specifying whether to take into account probe affinity with *probeAff*; setting the number of burn-in and post burn-in runs with the *burnin* and *iter* arguments respectively; setting the set of parameters to save with the *output* argument³; and specifying where to save the runs with *rundir*. Execute `help(bgx)` in R for a full explanation of all the parameters.

As an example, let us analyse the `Dilution` data set and save the results in the current working directory ("."):

```
> library(affydata)
> library(hgu95av2cdf)
> data(Dilution)
> eset <- bgx(Dilution, samplesets=c(2,2), probeAff=FALSE, burnin=2048, iter=8192,ge
```

The `eset` object will contain gene expression information for each gene under each condition (not necessarily each array). You may obtain the gene expression measure using the `exprs` function. For instance:

```
> exprs(eset)[10:40,] # Shorthand for assayData(eset)$exprs[10:40,]
```

| | condition 1 | condition 2 |
|----------|-------------|-------------|
| 947_at | 6.54444 | 6.25029 |
| 948_s_at | 4.84583 | 4.45123 |
| 949_s_at | 4.84665 | 4.55750 |
| 950_at | 4.48827 | 4.27474 |
| 951_at | 2.70152 | 2.49357 |
| 952_at | 1.58622 | 1.94747 |
| 953_g_at | 5.28392 | 4.89309 |

²Note that if your `AffyBatch` object contains information on the experimental design in the `phenoData` slot, you do not need to use the *samplesets* argument.

³*output* can be set to either "minimal", "trace" or "all". See the documentation for an explanation of what these levels mean

| | | |
|----------|---------|---------|
| 954_s_at | 6.36809 | 6.09223 |
| 955_at | 6.60929 | 6.33914 |
| 956_at | 7.00133 | 6.70313 |
| 957_at | 4.62867 | 4.25959 |
| 958_s_at | 5.53495 | 5.18097 |
| 959_at | 1.90068 | 1.65003 |
| 960_g_at | 5.22854 | 4.93502 |
| 961_at | 1.71897 | 1.60414 |
| 962_at | 2.46971 | 2.01055 |
| 963_at | 4.55256 | 4.25666 |
| 964_at | 4.27345 | 3.96223 |
| 965_at | 2.17875 | 1.15486 |
| 966_at | 4.44520 | 3.96493 |
| 967_g_at | 4.85540 | 4.59308 |
| 968_i_at | 3.33662 | 3.63703 |
| 969_s_at | 4.80061 | 4.41228 |
| 970_r_at | 6.29191 | 6.17022 |
| 971_s_at | 2.03316 | 2.89112 |
| 973_at | 4.35829 | 4.09485 |
| 974_at | 2.27413 | 2.23694 |
| 975_at | 4.26401 | 3.96236 |
| 976_s_at | 3.52073 | 3.17471 |
| 977_s_at | 4.83872 | 4.59583 |
| 978_at | 3.19150 | 2.26953 |

Run `help(ExpressionSet)` in R for more information.

Note that *samplesets* should be set to an array specifying the number of replicates in each condition. If set to (3,2), **bgx** will treat the first three arrays read into R as replicates under condition 1 and the next two as replicates under condition 2. You should make sure that all condition 1 files are read in first and all condition 2 files are read in second by `ReadAffy()`. You may check the order of the samples in your `AffyBatch` object by using the `sampleNames` function:

```
> sampleNames(Dilution)
[1] "20A" "20B" "10A" "10B"
```

4 Running BGX as a standalone binary

Occasionally it may be useful to run **bgx** as a standalone binary from the command line⁴. In this case, you should use the `standalone.bgx` function instead of the `bgx` function.

⁴You can compile it by tweaking 'src/Makefile.standalone' to your specifications and running 'make -f Makefile.standalone' from the 'src' directory.

It takes the same arguments as `bgx`, with the addition of *dirname*, which should specify where you would like to save the input files required by the standalone binary.

```
aData <- ReadAffy() # Read in 6 arrays across two conditions
                  # in alphabetical order
standalone.bgx(aData, samplesets=c(3,3), genes=c(1:650,1000:1200),
  burnin=16384, iter=65536, output="minimal",
  dirname="input-choe3replicates")
```

Once you have saved the input files, you should locate the binary, make sure it is executable⁵, and pass the path to the newly created `infile.txt` file as a single argument. For example:

```
./bgx ../input-choe3replicates/infile.txt
```

5 Detailed analysis of the output

If you wish to analyse the output in detail, you should first read the output into a list as follows:

```
> bgxOutput <- readOutput.bgx("run.1")
```

You may then pass the `bgxOutput` object to any of several analysis functions. For instance, to view the gene expression distributions under the various conditions for gene 10, you could do:

```
> plotExpressionDensity(bgxOutput, gene=10)
```

⁵Under Unix-like environments, you can type `chmod +x bgx` at the command prompt to do this.



In order to get a list of ranked differential expression values, you could do:

```
> rankedGeneList <- rankByDE(bgxOutput)
> print(rankedGeneList[1:25,]) # print top 25 DEG
```

| | Position | DiffExpression |
|---------------------------|----------|----------------|
| 941_at | 4 | 36.347898 |
| 956_at | 19 | 31.947193 |
| 955_at | 18 | 30.719990 |
| AFFX-HUMGAPDH/M33197_5_at | 89 | 29.730292 |
| AFFX-HSAC07/X00351_5_at | 83 | 28.877642 |
| AFFX-HUMGAPDH/M33197_M_at | 91 | 25.443945 |
| 947_at | 10 | 25.033984 |
| AFFX-HSAC07/X00351_M_at | 85 | 24.970413 |
| 954_s_at | 17 | 22.723817 |
| 946_at | 9 | 22.693250 |
| 958_s_at | 21 | 21.205970 |
| AFFX-HUMGAPDH/M33197_3_at | 87 | 18.793519 |

| | | |
|-----------------------------|----|-----------|
| AFFX-HUMISGF3A/M97935_MB_at | 96 | 17.405426 |
| 953_g_at | 16 | 16.603250 |
| AFFX-BioDn-3_at | 70 | 15.625753 |
| 982_at | 44 | 14.335130 |
| AFFX-HUMISGF3A/M97935_3_at | 93 | 13.916216 |
| AFFX-HUMISGF3A/M97935_MA_at | 95 | 12.615134 |
| 948_s_at | 11 | 12.404184 |
| AFFX-HSAC07/X00351_3_at | 81 | 12.382537 |
| 969_s_at | 32 | 12.107360 |
| 993_at | 54 | 12.089784 |
| 957_at | 20 | 11.909911 |
| 960_g_at | 23 | 9.688060 |
| 949_s_at | 12 | 9.650797 |

Run `help(analysis.bgx)` for more detailed usage instructions on the analysis functions.