

<b>Recommendations for normalization of microarray data</b>			
Author(s): Tim Beissbarth, Markus Ruschhaupt, David Jackson, Chris Lawerenz, Ulrich Mansmann			
Created on: 11.11.2005	Version: 1.1	No.: 1.1	Page 1 of 4

Normalization is an essential procedure in the analysis of DNA microarrays to compare data from different arrays or colour channels. Measurements on microarrays may be systematically biased by diverse effects such as efficiency of RNA extraction, reverse transcription, label incorporation, exposure, scanning, spot detection, etc. Furthermore, there are systematic effects due to characteristics of the array, such as effects of different probes (i.e. cDNAs or oligos), spotting effects, region effects, pin effects, etc.. Normalization attempts to compensate for such effects through use of internal controls.

There are three types of internal controls that can be used for normalization:

1. Most commonly normalization is based on all genes on the array. The assumption is used that between two conditions the majority of genes do not change in terms of their expression level.
2. Previously known Housekeeping-genes.
3. Spiked-in control genes.

Usually one of these methods is used for normalization and the other two can be used to validate the results. In most cases, normalization is most stable using the majority of genes, at least for microarrays where “all” genes in a genome are represented. For specialized arrays, e.g. with only disease specific genes, it might be necessary to spot a sufficient number of housekeeping genes or spiked-in controls and use these for normalization. Special care should be taken, when using housekeeping genes for normalization, as they often vary depending on conditions. To avoid errors use big enough numbers (50+) and not just 1 or 2.

It is, however, necessary to make sure that the data being compared is actually comparable. Normalization will always make data look increasingly similar. However, this doesn't make sense if there are fundamental reasons why data is not comparable.

Normalization consists of several steps:

- Background Correction
  - Local background (image analysis) – if the image contains heterogeneous background, scratches, etc.
  - Global background (e.g. 5% quantile)
  - No background correction – some people avoid background correction. This usually leads to some underestimation of the ratio of differential expression, but avoids some of the problems connected to background correction.

- Transformation: Microarray intensities should always be looked at using log<sub>2</sub> scale. This scaling should roughly adjust the variance to be the same for all intensities. Differences of log<sub>2</sub> intensities reflect the log<sub>2</sub> ratios (M values) for a comparison. As an alternative to log<sub>2</sub> scaling, the variance stabilizing transformation VSN (see below) may be used (this is roughly equivalent to using the natural logarithm instead of log<sub>2</sub>!).
- Robust estimation of a “rescaling” factor (e.g. median of differences). One of the method for internal control (all genes, housekeeping genes, or spiked-in controls) is used for this purpose.

There are many normalization methods! Which of the methods is most stable and gives best results is dependent on the type of data, the image analysis program, etc. To determine the best method, it is a good idea to try several methods initially on a few datasets and inspect the results visually using controls.

- Scale normalization: the simplest way to normalize data is simply to adjust the scale of the data, e.g. set the median of differences to 0. This does not consider any region or intensity dependent effects, however.
- Lowess (aka loess): Local regression does take into account intensity dependent effects and might partially correct for background (additive) effects. There are also variants that take into account local effects, e.g. print-tip lowess. This type of normalization is most commonly used for two-colour arrays.
- Quantile: Similar idea to scale normalization but more drastic, as all of the various quantiles are adjusted and not only the 50% quantile (median). This type of normalization is most commonly used for affymetrix arrays.
- VSN: Estimates an additive and multiplicative offset and transforms the data to have equal variance for all intensities. This transformation is similar to using the natural log transformation, but tries to adjust for effects that are often observed after background subtraction (i.e. high-variance for lowly expressed genes). VSN can be used with cDNA or affymetrix data, and is advisable if you observe unstable results with lowly expressed genes.

Quality control on your experiment has to be performed.

- Negative controls, heterologous DNA, genes from different organisms can be spotted to check background hybridization levels.
- One approach is to use housekeeping genes whose expression is assumed to be constant under all conditions.
- Spiked-in controls (positive dynamic range controls, negative controls, ratio controls) provide a good indication of the quality of an experiment.
- Check whether the controls behave as expected after normalization.

Visual inspection of the data before and after normalization, for example, by means of a scatter plot or MA plot is essential and can help to avoid serious errors during the normalization procedure. Scatter Plots (log<sub>2</sub> red channel vs. log<sub>2</sub> green channel, or log<sub>2</sub> expr. level 1 vs. log<sub>2</sub> expr. level 2) or MA-plots (average log<sub>2</sub> expression A vs. log<sub>2</sub> difference M) are common ways to quickly inspect the data of a comparison. Further, to look at the spatial distribution of expression values or quality values on a chip or across different chips can help to detect problems with printing or hybridization.

## Conclusions

Normalization is essential to compare the varying conditions of microarray experiments. While there are several methods for normalization, the choice of which method gives the best results really depends on your local settings. Use visual inspection of the data and comparisons before and after normalization to control that the procedure worked correctly. It is standard to display data in log<sub>2</sub> scale. Use different types of controls, i.e. negative controls, spiked in controls, and housekeeping genes when spotting arrays, to be able to observe possible problems with the hybridizations or normalization procedure. Whenever possible normalize on the majority of genes or use a sufficiently large number of non-differentially expressed controls to normalize the data.

## References

Beissbarth T, Fellenberg K, Brors B, Arribas-Prat R, Boer JM, Hauser NC, Scheideler M, Hoheisel JD, Schuetz G, Poustka A, Vingron M: Processing and quality control of DNA array hybridization data. *Bioinformatics*; 11.2000; 16(11): 1014-1022.

Huber W., von Heydebreck A., Sültmann H., Poustka A. and Vingron M. (2002): Variance stabilization applied to microarray data calibration and to the quantification of differential expression. *Bioinformatics* 18: suppl. 1 (2002), S96-S104 (ISMB 2002).

