

Variation-preserving normalization unveils blind spots in gene expression profiling

Carlos P. Roca^{1,4,5}, Susana I. L. Gomes², Mónica J. B. Amorim², and
Janeck J. Scott-Fordsmand^{3,5}

¹Department of Chemical Engineering, Universitat Rovira i Virgili, 43007
Tarragona, Spain

²Department of Biology & CESAM, University of Aveiro, 3810-193
Aveiro, Portugal

³Department of Bioscience, Aarhus University, 8600 Silkeborg, Denmark

Keywords: differential gene expression; gene expression microarrays; RNA-Seq; data
normalization

⁴Current address: Department of Bioscience, Aarhus University, 8600 Silkeborg, Denmark.

⁵Correspondence should be addressed to C.P.R. (carlosp.roca@urv.cat) or J.J.S.-F. (jsf@bios.au.dk)

Abstract

RNA-Seq and gene expression microarrays provide comprehensive profiles of gene activity, but lack of reproducibility has hindered their application. A key challenge in the data analysis is the normalization of gene expression levels, which is currently performed following an implicit assumption that most genes are not differentially expressed. Here, we present a mathematical approach to normalization that makes no assumption of this sort. We have found that variation in gene expression is much greater than currently believed, and that it can be measured with available technologies. Our results also explain, at least partially, the problems encountered in transcriptomics studies. We expect this improvement in detection to help efforts to realize the full potential of gene expression profiling, especially in analyses of cellular processes involving complex modulations of gene expression.

Introduction

Since the discovery of DNA structure by Watson and Crick, molecular biology has progressed increasingly quickly, with rapid advances in sequencing and related genomic technologies. Among these, microarrays and RNA-Seq have been widely adopted to obtain gene expression profiles, by measuring the concentration of tens of thousands of mRNA molecules in single assays (Schena et al., 1995; Lockhart et al., 1996; Duggan et al., 1999; Mortazavi et al., 2008; Wang et al., 2009). Despite their enormous potential (Golub et al., 1999; van 't Veer et al., 2002; Ivanova et al., 2002; Chi et al., 2003), problems of reproducibility and reliability (Tan et al., 2003; Frantz, 2005; Couzin, 2006) have discouraged their use in some areas, e.g. biomedicine (Michiels et al., 2005; Weigelt and Reis-Filho, 2010; Brettingham-Moore et al., 2011). In more mature microarray technologies, issues such as probe design, cross-hybridization, non-linearities and batch effects (Draghici et al., 2006) have been identified as possible culprits, but the problems persist (Shi et al., 2006; Su et al., 2014).

The normalization of gene expression, which is required to set a common reference level

among samples (Smyth and Speed, 2003; Irizarry et al., 2003; Bullard et al., 2010; Garber et al., 2011; Dillies et al., 2013), is also reportedly problematic, affecting the reproducibility of results with both microarray (Shi et al., 2006; Shippy et al., 2006) and RNA-Seq (Su et al., 2014; Bullard et al., 2010; Dillies et al., 2013). Batch effects and their influence on normalization have recently received a great deal of attention (Leek et al., 2010; Reese et al., 2013; Li et al., 2014), resulting in approaches aiming to remove unwanted technical variation caused by differences between batches of samples or by other sources of expression heterogeneity (Listgarten et al., 2010; Gagnon-Bartsch and Speed, 2012; Risso et al., 2014). A different issue, however, is the underlying assumption made by the most widely used normalization methods to date, such as median and quantile normalization (Bolstad et al., 2003) for microarrays, or RPKM (Mortazavi et al., 2008) and TMM (Robinson and Oshlack, 2010) for RNA-Seq, which posit that most genes are not differentially expressed (Dillies et al., 2013; Hicks and Irizarry, 2015). This *lack-of-variation assumption* may seem reasonable for many applications, but it has not been confirmed. Furthermore, results obtained with other technologies, particularly qRT-PCR, suggest that it may not be valid (Shi et al., 2006; Bullard et al., 2010).

Some methods have been proposed to address the issue of the lack-of-variation assumption, based on the use of spike-ins (Lovén et al., 2012), negative control probes (Wu and Aryee, 2010) or negative control genes (Gagnon-Bartsch and Speed, 2012), that is, on external or internal controls that are *known a priori* not to be differentially expressed (Lippa et al., 2010). The applicability of these methods, however, has been limited by this requirement of a priori knowledge, which is rarely available for a sufficiently large number of controls. Thus, in attempts to clarify and overcome limitations imposed by the lack-of-variation assumption, we have developed an approach to normalization that does not assume lack-of-variation and that does not require the use of spike-ins or a priori knowledge of control genes. The analysis of a large gene expression dataset using this approach shows that the assumption can severely undermine the detection of variation in gene expression. We have found that large numbers of differentially expressed genes with substantial expression changes are missed when data are normalized with methods that assume lack-of-variation.

Results

Datasets and Normalization Methods

The dataset was obtained from biological triplicates of *Enchytraeus crypticus* (a globally distributed soil organism used in standard ecotoxicity tests), sampled under 51 experimental conditions (42 treatments and 9 controls), involving exposure to several substances, at several concentrations and durations according to a factorial design (Supp. Table 1). Gene expression was measured using a customized high-density oligonucleotide microarray, and the resulting dataset was normalized with four methods. Two of these methods are the most widely used procedures for microarrays, median (or scale) normalization and quantile normalization (Bolstad et al., 2003), whereas the other two, designated *median condition-decomposition normalization* and *standard-vector condition-decomposition normalization*, have been developed for this study.

With the exception of quantile normalization, all used methods apply a multiplicative factor to the expression levels in each sample, equivalent to the addition of a number in the usual \log_2 -scale for gene expression levels. Solving the *normalization problem* consists of finding these correction factors. The problem can be exactly and linearly decomposed into several sub-problems: one within-condition normalization for each experimental condition and one final between-condition normalization for the condition averages. In the within-condition normalizations, the samples (replicates) subjected to each experimental condition are normalized separately, whereas in the final between-condition normalization average levels for all conditions are normalized together. Because there are no genes with differential expression in any of the within-condition normalizations, the lack-of-variation assumption only affects the final between-condition normalization. The assumption is avoided by using, in this normalization, expression levels only from *no-variation genes*, i.e. genes that show no evidence of differential expression under a statistical test. Both methods of normalization proposed here follow this condition-decomposition approach.

With median condition-decomposition normalization, all normalizations are performed

with median values, as in conventional median normalization, but only no-variation genes are included in the between-condition step. Otherwise, if all genes were used in this final step, the resulting total normalization factors would be exactly the same as those obtained with conventional median normalization.

For standard-vector condition-decomposition normalization, a vectorial procedure was developed to carry out each normalization step. The samples of any experimental condition, in a properly normalized dataset, must be *exchangeable*. In mathematical terms, the expression levels of each gene can be considered as an s -dimensional vector, where s is the number of samples for the experimental condition. After standardization (mean subtraction and variance scaling), these standard vectors are located in a $(s - 2)$ -dimensional hypersphere. The exchangeability mentioned above implies that, when properly normalized, the distribution of standard vectors must be invariant with respect to permutations of the sample labels and must have zero expected value. These properties allow to obtain, under fairly general assumptions, a robust estimator of the normalization factors.

To further explore and compare outcomes of the normalization methods, they were also applied to a synthetic random dataset. This dataset was generated with identical means and variances gene-by-gene to the real dataset, and with the assumption that all genes were no-variation genes. In addition, normalization factors were applied, equal to those obtained from the real dataset. Thus, the synthetic dataset was very similar to the real one, while complying by construction with the lack-of-variation assumption.

Normalization Results

Figure 1 displays the results of applying the four normalization methods to the real and synthetic datasets. Each panel shows the interquartile ranges of expression levels for the 153 samples, grouped in triplicates exposed to each experimental condition. Both median (second row) and quantile normalization (third row) yielded similar outputs, for both datasets. In contrast, the condition-decomposition normalizations (fourth and fifth rows) identified marked differences, detecting much greater variation between conditions in the

real dataset. Conventional median normalization makes, by design, the median of each sample the same, while quantile normalization makes the full distribution of each sample the same. Hence, if there were differences in medians or distributions of gene expression between experimental conditions, both methods would have removed them. Figures 1G,I show that such variation between conditions was present in the real dataset.

Influence of no-variation genes on normalization

To clarify how the condition-decomposition normalizations preserved the variation between conditions, we studied the influence of the choice of no-variation genes in the final between-condition normalization. To this end, we obtained the between-condition variation with both methods in two families of cases. In one family, no-variation genes were chosen in decreasing order of p -values from an ANOVA test. In the other family, genes were chosen at random. The first option was similar to the approach implemented to obtain the results presented in Figures 1G–J, with the difference that there the number of genes was chosen automatically by a statistical test. As shown in Figure 2A, for the real dataset the random choice of genes resulted in $n^{-1/2}$ decays (n being the number of chosen genes), followed by a plateau. The $n^{-1/2}$ decays reflect the standard errors of the estimators of the normalization factors. Selecting the genes by decreasing p -values, however, yielded a completely different result. Up to a certain number of genes, the variance remained similar, but for larger numbers of genes the variance dropped rapidly. Figure 2A shows, therefore, that between-condition variation was removed as soon as the between-condition normalizations used genes that varied in expression level across experimental conditions. The big circles in Figure 2A indicate the working points of the normalizations used to generate the results displayed in Figures 1G,I. In fact, these points slightly underestimated the variation between conditions. Although the statistical test for identifying no-variation genes ensured that there was no evidence of variation, inevitably the expression of some selected genes varied across conditions.

Figure 2B displays the results obtained with the synthetic dataset. There were no plateaus

when no-variation genes were chosen randomly, only $n^{-1/2}$ decays, and small differences when no-variation genes were selected by decreasing p -values. Big circles show that working points were selected with much larger numbers of genes in the synthetic dataset (Figs. 1H,J) than in the real dataset (Figs. 1G,I). The residual variation, produced by errors in the estimation of the normalization factors, was much smaller than the variation detected in the real dataset, especially for standard-vector condition-decomposition normalization. Overall, Figure 2 shows that the between-condition variation pictured in Figures 1G,I is not an artifact caused by using an exceedingly small or extremely particular set of genes in the final between-condition normalization, but that this variation originated from the real dataset.

Differential Gene Expression

Finally, Figure 3A shows the numbers of differentially expressed gene probes (DEGP), identified after normalizing with the four methods, for each of the 42 experimental treatments versus the corresponding control (Supp. Table 2). Compared to conventional methods, the number of DEGP detected with the condition-decomposition normalizations was much larger under most treatments, including some whose number of DEGP was larger by more than one order of magnitude. These are statistically significant changes of gene expression, i.e. changes that cannot be explained by chance. More important is the scale of the detected variation, as illustrated by the boxplots in Figure 3C showing absolute fold changes of DEGP detected after standard-vector condition-decomposition normalization. For all treatments, the entire interquartile range of absolute fold change is above 1.5-fold, and for more than two thirds of the treatments the median absolute fold change is greater than 2. This amount of gene expression variation cannot be neglected, and warrants further research to explore its biological significance.

Discussion

The variation between medians displayed in Figures 1G,I may seem surprising, given routine expectations based on current methods (Figs. 1C,E). Nevertheless, this variation inevitably results from the imbalance between over- and under-expressed genes. As an illustration, let us consider a case with two experimental conditions, in which the average expression of a given gene is less than the distribution median under one condition, but greater than the median under the other. The variation of this gene alone will change the value of the median to the expression level of the next ranked gene. Therefore, if the number of over-expressed genes is different from the number of under-expressed genes, and enough changes cross the median boundary, then the median will substantially differ between conditions. Only when the differential expression is balanced or small enough, will the median stay the same. This argument applies equally to any other quantile in the distribution of gene expression. Transcriptional amplification is an extreme example of change in the distribution of expression levels (Lovén et al., 2012), which can nevertheless be properly normalized with condition-decomposition methods, and without resorting to spike-ins as long as some genes are not differentially expressed.

An important feature of the approaches to normalization proposed here (linear decomposition into normalization sub-problems per condition, and standard-vector normalization for each sub-problem) is that they do not depend on any particular aspect of the technology of gene expression microarrays or RNA-Seq. The numbers in the input data are interpreted as measured concentrations of mRNA molecules, in order to identify the normalization factors and irrespectively of whether the concentrations were obtained from fluorescence intensities of hybridized cDNA (microarrays) or from counts of fragments read of mRNA sequences (RNA-Seq). Nevertheless, we consider that specific within-sample corrections for each technology are still necessary and must be applied *before* the between-sample normalizations proposed here. Examples include background correction for microarrays or gene-length normalization (RPKM) for RNA-Seq. Equally, methods that address the influence of biological or technical confounding factors on downstream

analyses, such as SVA (Leek and Storey, 2007) or PEER (Stegle et al., 2010), should be applied when necessary, *after* normalizing.

The lack-of-variation assumption underlying the current methods of normalization was self-fulfilling, removing variation in gene expression that was present in the real dataset. Moreover, it had negative consequences for downstream analyses, as it both removed potentially important biological information and introduced errors in the detection of gene expression. A removal of variation can be understood as errors in the estimation of normalization factors. Considering data and errors vectorially, the length of each vector equals, after centering and up to a constant factor, the standard deviation of the data or error. The addition of an error of small magnitude, compared to the data variance, would have only a minor effect. However, errors of similar or greater magnitude than the data variance may, depending on the lengths and relative angles of the vectors, severely distort the observed data variance. This will in turn cause spurious results in the statistical analyses. Furthermore, the angles between the data and the correct normalization factors (considered as vectors) are random. Data reflect biological variation, while normalization factors respond to technical variation. If the experiment is repeated, even with exactly the same experimental settings, the errors in the normalization factors will vary randomly, causing random spurious results in the downstream analyses. This explains, at least partially, the lack of reproducibility found in transcriptomics studies, especially for the detection of small changes of gene expression, because small variations are most likely to be distorted by errors in the estimates of normalization factors. Accordingly, the largest differences in numbers of DEGP detected by conventional compared to condition-decomposition methods (Fig. 3A) occurred consistently in the treatments with the smallest magnitudes of gene expression changes, e.g. treatments 28, 29 and 33 (Figs. 3B,C).

In summary, this study proves that large numbers of genes change in expression level (often strongly) across experimental conditions, and too extensively to ignore in the normalization of gene expression data. Further, our approach, which avoids the prevailing lack-of-variation assumption, demonstrates that current normalization methods likely remove and distort important variation in gene expression. It also offers a means to inves-

tigate broad changes in gene expression that have remained hidden to date. We expect this to provide revealing insights about diverse biomolecular processes, particularly those involving substantial numbers of genes, such as cell differentiation, toxic responses, diseases with non-Mendelian inheritance patterns and cancer. After years of lagging behind the advances in genome sequencing, we believe that the procedures presented here will assist efforts to realize the full potential of gene expression profiling.

References

- Bolstad, B. M., Irizarry, R. A., Astrand, M. and Speed, T. P. (2003). A comparison of normalization methods for high density oligonucleotide array data based on variance and bias. *Bioinformatics* *19*, 185–193.
- Brettingham-Moore, K. H., Duong, C. P., Heriot, A. G., Thomas, R. J. S. and Phillips, W. A. (2011). Using gene expression profiling to predict response and prognosis in gastrointestinal cancers-the promise and the perils. *Ann of Surg Oncol* *18*, 1484–1491.
- Bullard, J. H., Purdom, E., Hansen, K. D. and Dudoit, S. (2010). Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* *11*, 94.
- Castro-Ferreira, M. P., de Boer, T. E., Colbourne, J. K., Vooijs, R., van Gestel, C. A. M., van Straalen, N. M., Soares, A. M. V. M., Amorim, M. J. B. and Roelofs, D. (2014). Transcriptome assembly and microarray construction for *Enchytraeus crypticus*, a model oligochaete to assess stress response mechanisms derived from soil conditions. *BMC Genomics* *15*, 302.
- Chang, Y., Lye, M. L. and Zeng, H. C. (2005). Large-scale synthesis of high-quality ultralong copper nanowires. *Langmuir* *21*, 3746–3748.
- Chi, J.-T., Chang, H. Y., Haraldsen, G., Jahnsen, F. L., Troyanskaya, O. G., Chang, D. S., Wang, Z., Rockson, S. G., van de Rijn, M., Botstein, D. and et al. (2003). Endothelial

- cell diversity revealed by global expression profiling. *Proc Natl Acad Sci USA* *100*, 10623–10628.
- Couzin, J. (2006). Genomics. Microarray data reproduced, but some concerns remain. *Science* *313*, 1559.
- Dillies, M.-A., Rau, A., Aubert, J., Hennequet-Antier, C., Jeanmougin, M., Servant, N., Keime, C., Marot, G., Castel, D., Estelle, J. and et al. (2013). A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. *Brief Bioinform* *14*, 671–683.
- Draghici, S., Khatri, P., Eklund, A. C. and Szallasi, Z. (2006). Reliability and reproducibility issues in DNA microarray measurements. *Trends Genet* *22*, 101–109.
- Duggan, D. J., Bittner, M., Chen, Y., Meltzer, P. and Trent, J. M. (1999). Expression profiling using cDNA microarrays. *Nat Genet* *21*, 10–14.
- Durbin, J. (1973). *Distribution Theory for Tests Based on the Sample Distribution Function*. Society for Industrial and Applied Mathematics, Philadelphia.
- Eaton, M. L. (2007). *Multivariate Statistics: A Vector Space Approach*. Institute of Mathematical Statistics, Beachwood, Ohio.
- Fang, K., Kotz, S. and Ng, K. W. (1990). *Symmetric Multivariate and Related Distributions*. Chapman and Hall, New York.
- Feller, W. (1971). *An Introduction to Probability Theory and Its Applications*, vol. 2,. 2 edition, Wiley, New York.
- Frantz, S. (2005). An array of problems. *Nat Rev Drug Discov* *4*, 362–363.
- Gagnon-Bartsch, J. A. and Speed, T. P. (2012). Using control genes to correct for unwanted variation in microarray data. *Biostatistics* *13*, 539–52.
- Garber, M., Grabherr, M. G., Guttman, M. and Trapnell, C. (2011). Computational methods for transcriptome annotation and quantification using RNA-seq. *Nat Methods* *8*, 469–477.

- Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J. P., Coller, H., Loh, M. L., Downing, J. R., Caligiuri, M. A. and et al. (1999). Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science* 286, 531–537.
- Gomes, S. I. L., Caputo, G., Pinna, N., Scott-Fordsmand, J. J. and Amorim, M. J. B. (2015a). Effect of 10 different TiO₂ and ZrO₂ (nano)materials on the soil invertebrate *Enchytraeus crypticus*. *Environ Toxicol Chem* doi:10.1002/etc.3080.
- Gomes, S. I. L., Scott-Fordsmand, J. J. and Amorim, M. J. B. (2015b). Cellular energy allocation to assess the impact of nanomaterials on soil invertebrates (Enchytraeids): The effect of Cu and Ag. *Int J Environ Res Public Health* 12, 6858–6878.
- Gupta, A. K., Varga, T. and Bodnar, T. (2013). Elliptically Contoured Models in Statistics and Portfolio Theory. Springer, New York.
- Hicks, S. C. and Irizarry, R. A. (2015). quantro: a data-driven approach to guide the choice of an appropriate normalization method. *Genome Biol* 16, 117.
- Huber, W., Carey, V. J., Gentleman, R., Anders, S., Carlson, M., Carvalho, B. S., Bravo, H. C., Davis, S., Gatto, L., Girke, T. and et al. (2015). Orchestrating high-throughput genomic analysis with Bioconductor. *Nat Methods* 12, 115–121.
- Irizarry, R. A., Hobbs, B., Collin, F., Beazer-Barclay, Y. D., Antonellis, K. J., Scherf, U. and Speed, T. P. (2003). Exploration, normalization, and summaries of high density oligonucleotide array probe level data. *Biostatistics* 4, 249–264.
- Ivanova, N. B., Dimos, J. T., Schaniel, C., Hackney, J. A., Moore, K. A. and Lemischka, I. R. (2002). A stem cell molecular signature. *Science* 298, 601–604.
- Kallenberg, O. (2005). Probabilistic Symmetries and Invariance Principles. Springer, New York.
- Leek, J. T., Scharpf, R. B., Bravo, H. C., Simcha, D., Langmead, B., Johnson, W. E., Geman, D., Baggerly, K. and Irizarry, R. A. (2010). Tackling the widespread and critical impact of batch effects in high-throughput data. *Nat Rev Genet* 11, 733–739.

- Leek, J. T. and Storey, J. D. (2007). Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genet* 3, 1724–35.
- Li, S., Labaj, P. P., Zumbo, P., Sykacek, P., Shi, W., Shi, L., Phan, J., Wu, P.-Y., Wang, M., Wang, C., Thierry-Mieg, D., Thierry-Mieg, J., Kreil, D. P. and Mason, C. E. (2014). Detecting and correcting systematic variation in large-scale RNA sequencing data. *Nat Biotechnol* 32, 888–895.
- Lippa, K. A., Duewer, D. L., Salit, M. L., Game, L. and Causton, H. C. (2010). Exploring the use of internal and external controls for assessing microarray technical performance. *BMC Res Notes* 3, 349.
- Listgarten, J., Kadie, C., Schadt, E. E. and Heckerman, D. (2010). Correction for hidden confounders in the genetic analysis of gene expression. *Proc Natl Acad Sci USA* 107, 16465–70.
- Lockhart, D. J., Dong, H., Byrne, M. C., Follettie, M. T., Gallo, M. V., Chee, M. S., Mittmann, M., Wang, C., Kobayashi, M., Horton, H. and et al. (1996). Expression monitoring by hybridization to high-density oligonucleotide arrays. *Nat Biotechnol* 14, 1675–1680.
- Lovén, J., Orlando, D. A. A., Sigova, A. A. A., Lin, C. Y. Y., Rahl, P. B. B., Burge, C. B. B., Levens, D. L. L., Lee, T. I. I. and Young, R. A. A. (2012). Revisiting global gene expression analysis. *Cell* 151, 476–482.
- Michiels, S., Koscielny, S. and Hill, C. (2005). Prediction of cancer outcome with microarrays: A multiple random validation strategy. *Lancet* 365, 488–492.
- Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L. and Wold, B. (2008). Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods* 5, 621–628.
- OECD (2004a). Guidelines for the Testing of chemicals No 202. *Daphnia* sp. Acute Immobilization Test. Organization for Economic Cooperation and Development, Paris.
- OECD (2004b). Guidelines for the Testing of chemicals No. 220. *Enchytraeid* Reproduction Test. Organization for Economic Cooperation and Development, Paris.

- R Core Team (2015). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna. URL: <http://www.R-project.org/>.
- Reese, S. E., Archer, K. J., Therneau, T. M., Atkinson, E. J., Vachon, C. M., de Andrade, M., Kocher, J.-P. A. and Eckel-Passow, J. E. (2013). A new statistic for identifying batch effects in high-throughput genomic data that uses guided principal component analysis. *Bioinformatics* 29, 2877–83.
- Risso, D., Ngai, J., Speed, T. P. and Dudoit, S. (2014). Normalization of RNA-seq data using factor analysis of control genes or samples. *Nat Biotechnol* 32, 896–902.
- Ritchie, M. E., Phipson, B., Wu, D., Hu, Y., Law, C. W., Shi, W. and Smyth, G. K. (2015). *limma* powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* 43, e47.
- Robinson, M. D. and Oshlack, A. (2010). A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol* 11, R25.
- Schena, M., Shalon, D., Davis, R. W. and Brown, P. O. (1995). Quantitative monitoring of gene expression patterns with a complementary DNA microarray. *Science* 270, 467–470.
- Scott-Fordsmand, J. J., Krogh, P. H. and Weeks, J. M. (2000). Responses of *Folsomia fimetaria* (Collembola: Isotomidae) to copper under different soil copper contamination histories in relation to risk assessment. *Environ Toxicol Chem* 19, 1297–1303.
- Shi, L., Reid, L. H., Jones, W. D., Shippy, R., Warrington, J. A., Baker, S. C., Collins, P. J., de Longueville, F., Kawasaki, E. S., Lee, K. Y. and et al. (2006). The MicroArray Quality Control (MAQC) project shows inter- and intraplatform reproducibility of gene expression measurements. *Nat Biotechnol* 24, 1151–1161.
- Shippy, R., Fulmer-Smentek, S., Jensen, R. V., Jones, W. D., Wolber, P. K., Johnson, C. D., Pine, P. S., Boysen, C., Guo, X., Chudin, E. and et al. (2006). Using RNA sample titrations to assess microarray platform performance and normalization techniques. *Nat Biotechnol* 24, 1123–1131.

- Smyth, G. K. and Speed, T. (2003). Normalization of cDNA microarray data. *Methods* *31*, 265–273.
- Stegle, O., Parts, L., Durbin, R. and Winn, J. (2010). A Bayesian framework to account for complex non-genetic factors in gene expression levels greatly increases power in eQTL studies. *PLoS Comput Biol* *6*, e1000770.
- Storey, J. D. (2003). The positive false discovery rate: a Bayesian interpretation and the q-value. *Ann Stat* *31*, 2013–2035.
- Storey, J. D. and Tibshirani, R. (2003). Statistical significance for genomewide studies. *Proc Natl Acad Sci USA* *100*, 9440–9445.
- Su, Z., Labaj, P. P., Li, S., Thierry-Mieg, J., Thierry-Mieg, D., Shi, W., Wang, C., Schroth, G. P., Setterquist, R. A., Thompson, J. F. and et al. (2014). A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium. *Nat Biotechnol* *32*, 903–914.
- Tan, P. K., Downey, T. J., Spitznagel, E. L., Xu, P., Fu, D., Dimitrov, D. S., Lempicki, R. A., Raaka, B. M. and Cam, M. C. (2003). Evaluation of gene expression measurements from commercial microarray platforms. *Nucleic Acids Res* *31*, 5676–5684.
- van 't Veer, L. J., Dai, H., van de Vijver, M. J., He, Y. D., Hart, A. A. M., Mao, M., Peterse, H. L., van der Kooy, K., Marton, M. J., Witteveen, A. T. and et al. (2002). Gene expression profiling predicts clinical outcome of breast cancer. *Nature* *415*, 530–536.
- Wang, Z., Gerstein, M. and Snyder, M. (2009). RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev Genet* *10*, 57–63.
- Weigelt, B. and Reis-Filho, J. S. (2010). Molecular profiling currently offers no more than tumour morphology and basic immunohistochemistry. *Breast Cancer Res* *12 Suppl. 4*, S5.
- Wu, Z. and Aryee, M. J. (2010). Subset Quantile Normalization Using Negative Control Features. *J Comput Biol* *17*, 1385–1395.

Data Deposition and Code Availability

MIAME-compliant microarray data from the experiment were submitted to the Gene Expression Omnibus (GEO) at the NCBI website (platform: GPL20310; series: GSE69746, GSE69792, GSE69793 and GSE69794). Custom code that reproduces all the reported results starting from the raw microarray data is available at the GitHub repository <https://github.com/carlosproca/gene-expr-norm-paper>.

Acknowledgements

This work was funded by the European Union FP7 projects MODERN (Ref. 309314-2) (C.P.R., J.J.S.-F.) and MARINA (Ref. 263215) (J.J.S.-F.), by FEDER through COMPETE (Programa Operacional Factores de Competitividade) and FCT (Fundação para a Ciência e Tecnologia) through project bio-CHIP (Ref. FCT EXPL/AAG-MAA/0180/2013) (S.I.L.G., M.J.B.A.), and by a post-doctoral grant (Ref. SFRH/BPD/95775/2013) (S.I.L.G.).

Author Contributions

S.I.L.G., M.J.B.A. and J.J.S.-F. designed the toxicity experiment. S.I.L.G. carried out the experimental work and collected the microarray data. C.P.R. designed and implemented the novel normalization methods. C.P.R. performed the statistical analyses. All the authors jointly discussed the results. C.P.R. drafted the paper, with input from all the authors. All the authors edited the final version of the paper.

Figures

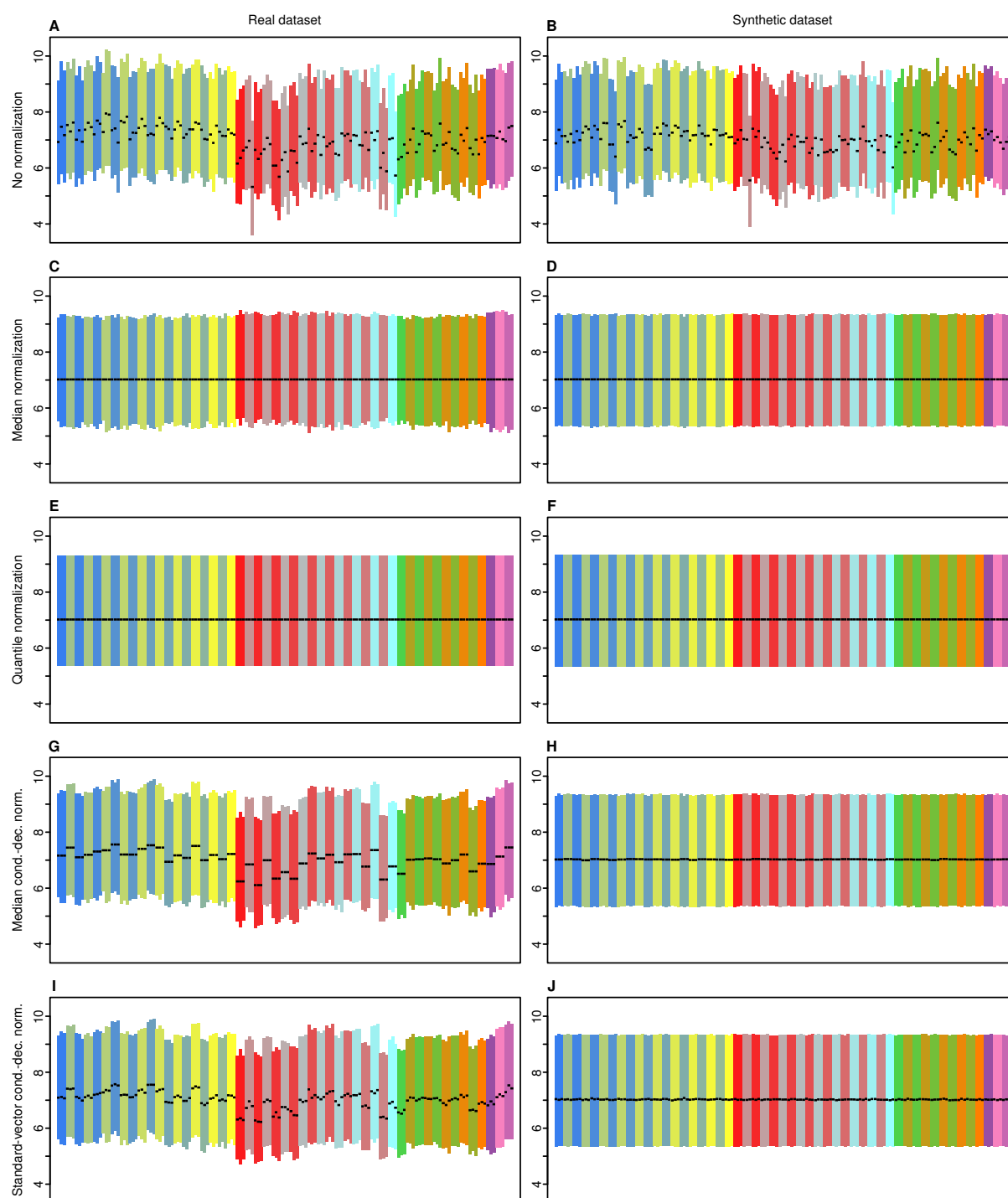


Figure 1

Figure 1: The condition-decomposition normalizations detected a large amount of between-condition variation in the real expression data, in contrast with conventional methods. All 10 panels show interquartile ranges of expression levels of the 153 samples, grouped by the 51 experimental conditions (Ag, blue-yellow; Cu, red-cyan; Ni, green-orange; UV, purple; see Supp. Table 1). Black lines indicate medians. Rows and columns correspond to normalization methods and datasets (as labeled), respectively. In the synthetic dataset no gene was differentially expressed between any two conditions.

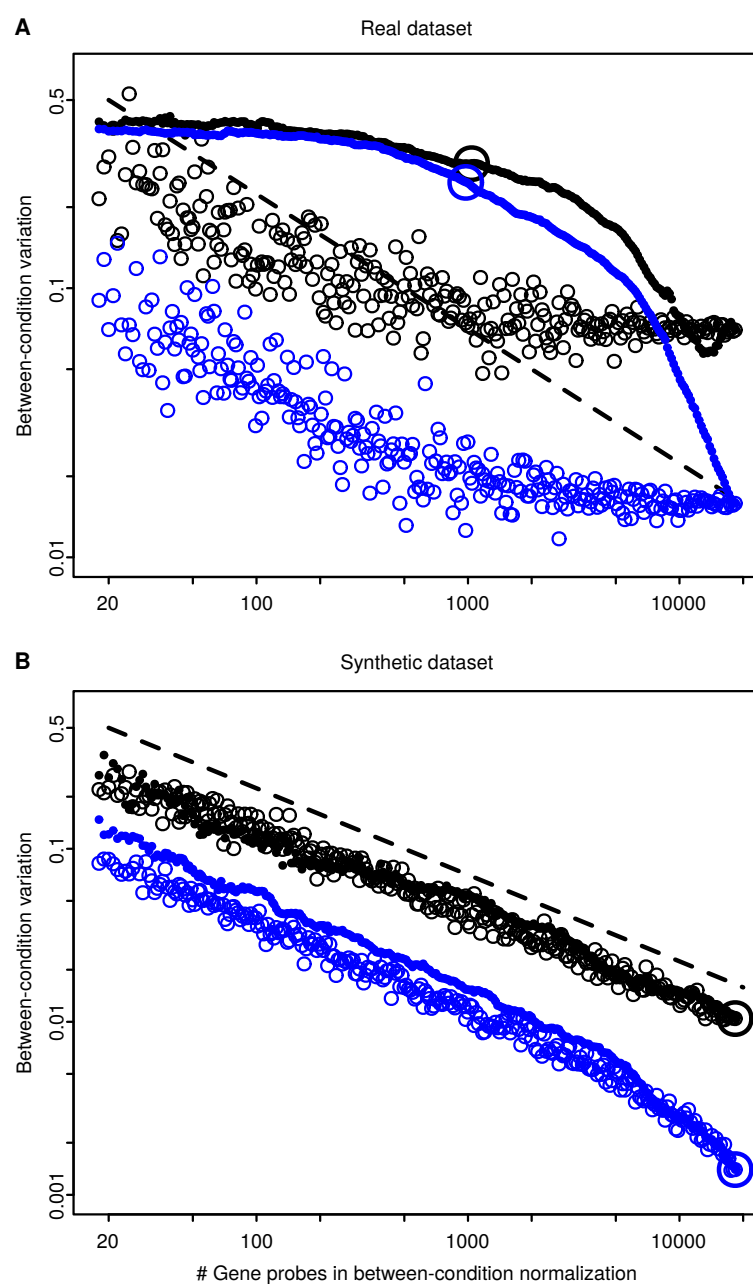


Figure 2

Figure 2: The selection of genes in the final between-condition normalization was crucial to preserve variation between conditions. The panels show the detected variation as a function of the number of gene probes used in the between-condition normalization of the real dataset (A) and synthetic dataset (B). Between-condition variation is represented as the standard deviation of the within-condition mean averages (averages of sample mean expression levels, for all samples under the condition). See Supplementary Figure 1 for within-condition median averages, with similar results. Each point in either of the panels indicates the variation obtained with one complete normalization (black circles, median condition-decomposition normalization; blue circles, standard-vector condition-decomposition normalization). Gene probes were selected in two ways: randomly (empty circles) or in decreasing order of p -values (filled circles). Big circles show the working points of the algorithms whose results are depicted in Figures 1G–J. Black dashed lines show references for $n^{-1/2}$ decays, with the same values in both panels.

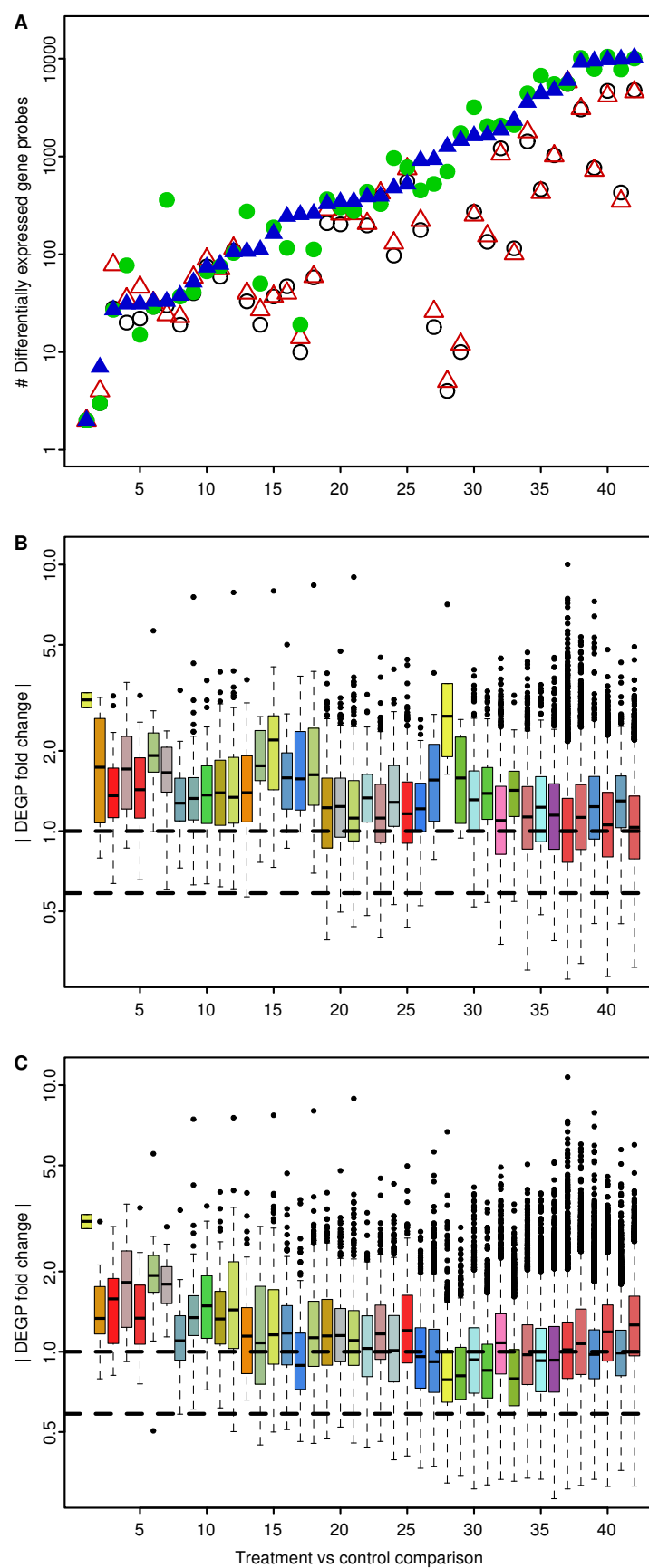


Figure 3

Figure 3: The condition-decomposition normalizations detected much larger numbers of differentially expressed gene probes (DEGP), with substantial fold changes. A: Number of DEGP for each treatment vs control comparison, obtained after applying the four normalization methods (empty black circles, median normalization; empty red triangles, quantile normalization; filled green circles, median condition-decomposition normalization; filled blue triangles, standard-vector condition decomposition normalization). Significant differential expression was identified with R/Bioconductor package limma. (see Supp. Fig. 2 for results with t-tests). Lower panel shows boxplots of absolute values of DEGP fold changes (absolute differences of \log_2 expression levels), also per treatment vs control comparison, obtained with quantile normalization (B) and standard-vector condition-decomposition normalization (C). Boxplots are colored by treatment, with the same color code as in Figure 1. In both panels comparisons are ordered according to the number of DEGP identified with standard-vector condition-decomposition normalization, increasing from left to right (Supp. Table 2). Dashed horizontal lines in the lower panel indicate references of 1.5-fold and 2-fold changes.

Materials and Methods

Test Organism and Exposure Media

The test species was *Enchytraeus crypticus*. Individuals were cultured in Petri dishes containing agar medium, in controlled conditions (Gomes et al., 2015b).

For copper (Cu) exposure, a natural soil collected at Hygum, Jutland, Denmark was used (Gomes et al., 2015b; Scott-Fordsmand et al., 2000). For silver (Ag) and nickel (Ni) exposure, the natural standard soil LUFA 2.2 (LUFA Speyer, Germany) was used (Gomes et al., 2015b). The exposure to ultra-violet (UV) radiation was done in ISO reconstituted water (OECD, 2004a).

Test Chemicals

The tested Cu forms (Gomes et al., 2015b) included copper nitrate ($\text{Cu}(\text{NO}_3)_2 \cdot 3\text{H}_2\text{O}$ > 99%, Sigma Aldrich), Cu nanoparticles (Cu-NPs, 20–30 nm, American Elements) and Cu nanowires (Cu-Nwires, synthesized by reduction of copper (II) nitrate with hydrazine in alkaline medium (Chang et al., 2005)).

The tested Ag forms (Gomes et al., 2015b) included silver nitrate AgNO_3 > 99%, Sigma Aldrich), non-coated Ag nanoparticles (Ag-NPs Non-Coated, 20–30 nm, American Elements),

Polyvinylpyrrolidone (PVP)-coated Ag nanoparticles (Ag-NPs PVP-Coated, 20–30 nm, American Elements), and Ag NM300K nanoparticles (Ag NM300K, 15 nm, JRC Repository). The Ag NM300K was dispersed in 4% Polyoxyethylene Glycerol Triolaete and Polyoxyethylene (20) orbitan mono-Laurat (Tween 20), thus the dispersant was tested alone as control (CTdisp).

The tested Ni forms included nickel nitrate ($\text{Ni}(\text{NO}_3)_2 \cdot 6\text{H}_2\text{O}$ \geq 98.5%, Fluka) and Ni nanoparticles (Ni-NPs, 20 nm, American Elements).

Spiking Procedure

Spiking for the Cu and Ag materials was done as in previous work (Gomes et al., 2015b). For the Ni materials, the Ni-NPs were added to the soil as powder, following the same procedure as for the Cu materials. NiNO₃, being soluble, was added to the pre-moistened soil as aqueous dispersions.

The concentrations tested were selected based on the reproduction effect concentrations EC₂₀ and EC₅₀, for *E. crypticus*, within 95% of confidence intervals, being: CuNO₃ EC_{20/50} = 290/360 mgCu/kg, Cu-NPs EC_{20/50} = 980/1760 mgCu/kg, Cu-Nwires EC_{20/50} = 850/1610 mgCu/kg, Cu-Field EC_{20/50} = 500/1400 mgCu/kg, AgNO₃ EC_{20/50} = 45/60 mgAg/kg, Ag-NP PVP-coated EC_{20/50} = 380/550 mgAg/kg, Ag-NP Non-coated EC_{20/50} = 380/430 mgAg/kg, Ag NM300K EC_{20/50} = 60/170 mgAg/kg, CTdisp = 4% w/w Tween 20, NiNO₃ EC_{20/50} = 40/60 mgNi/kg, Ni-NPs EC_{20/50} = 980/1760 mgNi/kg.

Four biological replicates were performed per test condition, including controls. For Cu exposure, the control condition for all the treatments consisted of soil from a control area at Hygum site, which has a Cu background concentration of 15 mg/kg (Scott-Fordsmand et al., 2000). For Ag exposure, two control sets were performed: CT (un-spiked LUFA soil, to be the control condition for AgNO₃, Ag-NPs PVP-Coated and Ag-NPs Non-Coated treatments) and CTdisp (LUFA soil spiked with the dispersant Tween 20, to be the control condition for the Ag NM300K treatments). For Ni exposure, the control consisted of un-spiked LUFA soil.

Exposure Details

In soil (i.e. for Cu, Ag and Ni) exposure followed the standard ERT (OECD, 2004b) with adaptations as follows: twenty adults with well-developed clitellum were introduced in each test vessel, containing 20 g of moist soil (control or spiked). The organisms were exposed for three and seven days under controlled conditions of photoperiod (16:8 h light:dark) and temperature 20 ± 1 °C without food. After the exposure period, the

organisms were carefully removed from the soil, rinsed in deionized water and frozen in liquid nitrogen. The samples were stored at -80°C , until analysis.

For UV exposure, the test conditions (OECD, 2004a) were adapted for *E. crypticus* (Gomes et al., 2015a). The exposure was performed in 24-well plates, where each well correspond to a replicate and contain 1 ml of ISO water and five adult organisms with clitellum. The test duration was five days, at $20 \pm 1^{\circ}\text{C}$. The organisms were exposed to UV on a daily basis, during 15 minutes per day to two UV intensities (280–400nm) of 1669.25 ± 50.83 and 1804.08 ± 43.10 mW/m², corresponding to total UV doses of 7511.6 and 8118.35 J/m², respectively. The remaining time was spent under standard laboratory illumination (16:8 h photoperiod). UV radiation was provided by an UV lamp (Spectroline XX15F/B, Spectronics Corporation, NY, USA, peak emission at 312 nm) and a cellulose acetate sheet was coupled to the lamp to cut-off UVC-range wavelengths (Gomes et al., 2015a). Thirty two replicates per test condition (including control without UV radiation) were performed to obtain 4 biological replicates with 40 organisms each for RNA extraction. After the exposure period, the organisms were carefully removed from the water and frozen in liquid nitrogen. The samples were stored at -80°C , until analysis.

RNA Extraction, Labeling and Hybridization

RNA was extracted from each replicate, which contained a pool of 20 and 40 organisms, for soil and water exposure, respectively. Three biological replicates per test treatment (including controls) were used. Total RNA was extracted using SV Total RNA Isolation System (Promega). The quantity and purity were measured spectrophotometrically with a nanodrop (NanoDrop ND-1000 Spectrophotometer) and its quality checked by denaturing formaldehyde agarose gel electrophoresis.

500 ng of total RNA were amplified and labeled with Agilent Low Input Quick Amp Labeling Kit (Agilent Technologies, Palo Alto, CA, USA). Positive controls were added

with the Agilent one-color RNA Spike-In Kit. Purification of the amplified and labeled cRNA was performed with RNeasy columns (Qiagen, Valencia, CA, USA).

The cRNA samples were hybridized on custom Gene Expression Agilent Microarrays (4 x 44k format), with a single-color design (Castro-Ferreira et al., 2014). Hybridizations were performed using the Agilent Gene Expression Hybridization Kit and each biological replicate was individually hybridized on one array. The arrays were hybridized at 65 °C with a rotation of 10 rpm, during 17 h. Afterwards, microarrays were washed using Agilent Gene Expression Wash Buffer Kit and scanned with the Agilent DNA microarray scanner G2505B.

Data Acquisition and Analysis

Fluorescence intensity data was obtained with Agilent Feature Extraction Software v. 10.7.3.1, using recommended protocol GE1_107_Sep09. Quality control was done by inspecting the reports on the Agilent Spike-in control probes. Background correction was provided by Agilent Feature Extraction software. To ensure an optimal comparison between the different normalization methods, only gene probes with good signal quality (flag IsPosAndSignif = True) in all samples were employed in the analyses. This implied the selection of 18,339 gene probes from a total of 43,750. Analyses were performed with R (R Core Team, 2015) v. 3.2.2, using R packages plotrix and RColorBrewer, and with Bioconductor (Huber et al., 2015) v. 3.1 packages genefilter and limma (Ritchie et al., 2015).

The synthetic data was generated gene by gene as normal variates with mean and variance equal, respectively, to the sample mean and sample variance of the real data. The applied normalization factors were those detected from the real data with standard-vector condition-decomposition normalization.

Median normalization was performed by subtracting the median of each sample distribution, and then adding the overall median to preserve the global expression level. Quantile

normalization was performed as implemented in the limma package.

The two condition-decomposition normalizations proceeded in the same way: first, 51 independent within-condition normalizations using all genes; then, final between-condition normalization, iteratively detecting no-variation genes and normalizing until convergence.

No-variation genes were identified with one-sided Kolmogorov-Smirnov tests, as goodness-of-fit tests against the uniform distribution, carried out on the greatest p -values obtained from an ANOVA test on the complete dataset (see below). The ANOVA test benefited from the already corrected within-condition variances, provided by the within-condition normalizations. The KS test was rejected at $\alpha = 0.001$.

The criterion for convergence for the median condition-decomposition (CD) normalizations was to require that the relative changes in the standard deviation of the normalization factors were less than 1%, or less than 10% for 10 steps in a row. In the case of standard-vector CD normalizations, convergence required that numerical errors were, compared to the estimated statistical errors (see below), less than 1%, or less than 10% for 10 steps in a row. For Figure 2 and Supplementary Figure 1, due to the very low number of gene probes in some cases, the thresholds for convergence for 10 steps in a row were increased to 80% and 50%, respectively, for median CD and standard-vector CD normalization.

In standard-vector CD normalization, the distribution of standard vectors was trimmed in each step to remove the 1% more extreme values of variance.

Differentially expressed gene probes were identified with limma (Fig. 3) or t-tests (Supp. Fig. 2), using in all cases a FDR threshold of 5%.

The reference distribution with permutation symmetry shown in the polar plots of the probability density function in Supplementary Movies 1–3 was calculated with the 6 permutations of the empirical standard vectors. The Watson U^2 statistic was calculated with the two-sample test (Durbin, 1973). An equal number of samples for comparison was obtained by sampling with replacement the permuted standard vectors.

Mathematical Methods

In a gene expression dataset with g genes, c experimental conditions and n samples per condition, the *observed* expression levels of gene j in condition k , $\mathbf{y}_j^{(k)} = (y_{1j}^{(k)}, \dots, y_{nj}^{(k)})'$, can be expressed in \log_2 -scale as

$$\mathbf{y}_j^{(k)} = \mathbf{x}_j^{(k)} + \mathbf{a}^{(k)}, \quad (1)$$

where $\mathbf{x}_j^{(k)}$ is the vector of *true* gene expression levels and $\mathbf{a}^{(k)}$ is the vector of normalization factors.

Given a sample vector \mathbf{x} , the mean vector is $\bar{\mathbf{x}} = \bar{x}\mathbf{1}$, and the residual vector is $\tilde{\mathbf{x}} = \mathbf{x} - \bar{\mathbf{x}}$. Then, (1) can be linearly decomposed into

$$\bar{\mathbf{y}}_j^{(k)} = \bar{\mathbf{x}}_j^{(k)} + \bar{\mathbf{a}}^{(k)}, \quad (2)$$

$$\tilde{\mathbf{y}}_j^{(k)} = \tilde{\mathbf{x}}_j^{(k)} + \tilde{\mathbf{a}}^{(k)}. \quad (3)$$

Equations (3) define the within-condition normalizations for each condition k . The scalar values in (2) are used to obtain the equations on condition means,

$$\bar{\mathbf{y}}_j^* = \bar{\mathbf{x}}_j^* + \bar{\mathbf{a}}^*, \quad (4)$$

$$\tilde{\mathbf{y}}_j^* = \tilde{\mathbf{x}}_j^* + \tilde{\mathbf{a}}^*. \quad (5)$$

The between-condition normalization is defined by (5). Equations (4) reduce to a single number, which is irrelevant to the normalization. The complete solution for each condition is obtained with $\mathbf{a}^{(k)} = \bar{\mathbf{a}}^{(k)} + \tilde{\mathbf{a}}^{(k)}$.

The n samples of gene j in a given condition can be modeled with the random vectors $\mathbf{X}_j, \mathbf{Y}_j \in \mathbb{R}^n$. Again, $\mathbf{Y}_j = \mathbf{X}_j + \mathbf{a}$, where \mathbf{a} is a fixed vector of normalization factors. It can be proved, under fairly general assumptions, that the true standard vectors have zero expected value

$$\mathbb{E} \left(\sqrt{n-1} \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \right) = \mathbf{0}, \quad (6)$$

whereas the observed standard vectors verify, as long as $\mathbf{a} \neq \mathbf{0}$,

$$0 < \mathbb{E} \left(\sqrt{n-1} \frac{\tilde{\mathbf{Y}}_j}{\|\tilde{\mathbf{Y}}_j\|} \right)' \frac{\tilde{\mathbf{a}}}{\|\tilde{\mathbf{a}}\|} < \mathbb{E} \left(\sqrt{n-1} \frac{1}{\|\tilde{\mathbf{Y}}_j\|} \right) \|\tilde{\mathbf{a}}\|. \quad (7)$$

This motivates the following iterative procedure to solve (3) and (5) (*standard-vector normalization*):

$$\hat{\mathbf{y}}_j^{(0)} = \tilde{\mathbf{y}}_j, \quad (8)$$

$$\hat{\mathbf{y}}_j^{(t)} = \hat{\mathbf{y}}_j^{(t-1)} - \hat{\mathbf{b}}^{(t-1)}, \quad \text{for } t \geq 1, \quad (9)$$

$$\hat{\mathbf{b}}^{(t)} = \frac{\sum_{j=1}^g \frac{\hat{\mathbf{y}}_j^{(t)}}{\|\hat{\mathbf{y}}_j^{(t)}\|}}{\sum_{j=1}^g \frac{1}{\|\hat{\mathbf{y}}_j^{(t)}\|}}, \quad \text{for } t \geq 0. \quad (10)$$

At convergence, $\lim_{t \rightarrow \infty} \hat{\mathbf{b}}^{(t)} = \mathbf{0}$, which implies $\lim_{t \rightarrow \infty} \hat{\mathbf{y}}_j^{(t)} = \tilde{\mathbf{x}}_j$ and $\sum_{t=0}^{\infty} \hat{\mathbf{b}}^{(t)} = \tilde{\mathbf{a}}$. Convergence is faster the more symmetric the empirical distribution of $\tilde{\mathbf{x}}_j/\|\tilde{\mathbf{x}}_j\|$ is on the unit $(n-2)$ -sphere. Convergence is optimal with spherically symmetric distributions, such as the Gaussian distribution, because in that case

$$\mathbb{E} \left(\frac{\tilde{\mathbf{Y}}_j}{\|\tilde{\mathbf{Y}}_j\|} \right) = \lambda \tilde{\mathbf{a}}, \quad \text{with } 0 < \lambda < \mathbb{E} \left(\frac{1}{\|\tilde{\mathbf{Y}}_j\|} \right). \quad (11)$$

Assuming no correlation between genes, an approximation of the statistical error at step t can be obtained with

$$\mathbb{E} \left(\|\hat{\mathbf{b}}^{(t)}\| \right) \approx \frac{\sqrt{g}}{\sum_{j=1}^g \frac{1}{\|\hat{\mathbf{y}}_j^{(t)}\|}}. \quad (12)$$

This statistical error is compared with the numerical error to assess convergence.

See Supplementary Material for a detailed exposition of the mathematical methods, and Supplementary Movies 1–5 for an illustration.

Supplementary Tables

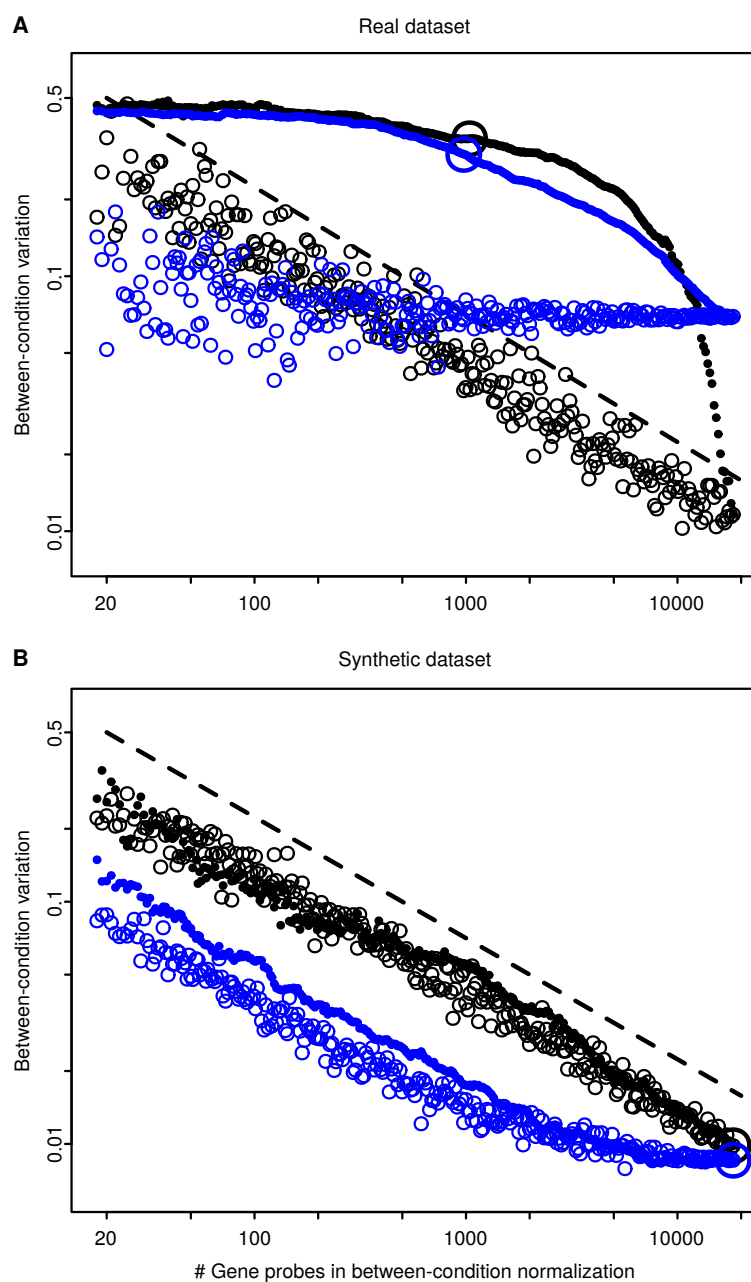
Supplementary Table 1: Experimental conditions of the toxicity experiment on *E. cryp-
ticus*, listed in the same order as they appear in each panel of Figure 1, from left to
right.

Condition number	Condition ID	Condition description
1	Ag.AgNO3.EC20.3d	AgNO3 EC20 3 days
2	Ag.AgNO3.EC20.7d	AgNO3 EC20 7 days
3	Ag.AgNO3.EC50.3d	AgNO3 EC50 3 days
4	Ag.AgNO3.EC50.7d	AgNO3 EC50 7 days
5	Ag.Coated.EC20.3d	Ag-NPs PVP-Coated EC20 3 days
6	Ag.Coated.EC20.7d	Ag-NPs PVP-Coated EC20 7 days
7	Ag.Coated.EC50.3d	Ag-NPs PVP-Coated EC50 3 days
8	Ag.Coated.EC50.7d	Ag-NPs PVP-Coated EC50 7 days
9	Ag.NC.EC20.3d	Ag-NPs Non-Coated EC20 3 days
10	Ag.NC.EC20.7d	Ag-NPs Non-Coated EC20 7 days
11	Ag.NC.EC50.3d	Ag-NPs Non-Coated EC50 3 days
12	Ag.NC.EC50.7d	Ag-NPs Non-Coated EC50 7 days
13	Ag.NM300K.EC20.3d	Ag NM300K EC20 3 days
14	Ag.NM300K.EC20.7d	Ag NM300K EC20 7 days
15	Ag.NM300K.EC50.3d	Ag NM300K EC50 3 days
16	Ag.NM300K.EC50.7d	Ag NM300K EC50 7 days
17	Ag.CT.3d	Ag Control 3 days
18	Ag.CT.7d	Ag Control 7 days
19	Ag.CTD.3d	Ag Control Dispersant 3 days
20	Ag.CTD.7d	Ag Control Dispersant 7 days
21	Cu.CuNO3.EC20.3d	CuNO3 EC20 3 days
22	Cu.CuNO3.EC20.7d	CuNO3 EC20 7 days
23	Cu.CuNO3.EC50.3d	CuNO3 EC50 3 days
24	Cu.CuNO3.EC50.7d	CuNO3 EC50 7 days
25	Cu.Cu.NPs.EC20.3d	Cu-NPs EC20 3 days
26	Cu.Cu.NPs.EC20.7d	Cu-NPs EC20 7 days
27	Cu.Cu.NPs.EC50.3d	Cu-NPs EC50 3 days
28	Cu.Cu.NPs.EC50.7d	Cu-NPs EC50 7 days
29	Cu.Cu.Nwires.EC20.3d	Cu-NWires EC20 3 days
30	Cu.Cu.Nwires.EC20.7d	Cu-NWires EC20 7 days
31	Cu.Cu.Nwires.EC50.3d	Cu-NWires EC50 3 days
32	Cu.Cu.Nwires.EC50.7d	Cu-NWires EC50 7 days
33	Cu.Cu.field.EC20.3d	Cu-Field EC20 3 days
34	Cu.Cu.field.EC20.7d	Cu-Field EC20 7 days
35	Cu.Cu.field.EC50.3d	Cu-Field EC50 3 days
36	Cu.Cu.field.EC50.7d	Cu-Field EC50 7 days
37	Cu.CT.3d	Cu Control 3 days
38	Cu.CT.7d	Cu Control 7 days
39	Ni.NiNO3.EC20.3d	NiNO3 EC20 3 days
40	Ni.NiNO3.EC20.7d	NiNO3 EC20 7 days
41	Ni.NiNO3.EC50.3d	NiNO3 EC50 3 days
42	Ni.NiNO3.EC50.7d	NiNO3 EC50 7 days
43	Ni.Ni.NPs.EC20.3d	Ni-NPs EC20 3 days
44	Ni.Ni.NPs.EC20.7d	Ni-NPs EC20 7 days
45	Ni.Ni.NPs.EC50.3d	Ni-NPs EC50 3 days
46	Ni.Ni.NPs.EC50.7d	Ni-NPs EC50 7 days
47	Ni.CT.3d	Ni Control 3 days
48	Ni.CT.7d	Ni Control 7 days
49	Uv.UV.D1.5d	UV Dose 1
50	Uv.UV.D2.5d	UV Dose 2
51	Uv.CT.5d	UV Control

Supplementary Table 2: Treatment vs control comparisons, listed in increasing number of differentially expressed gene probes (DEGP) obtained with standard-vector condition-decomposition normalization and limma statistical analysis. This is the same order as in Figure 3, from left to right.

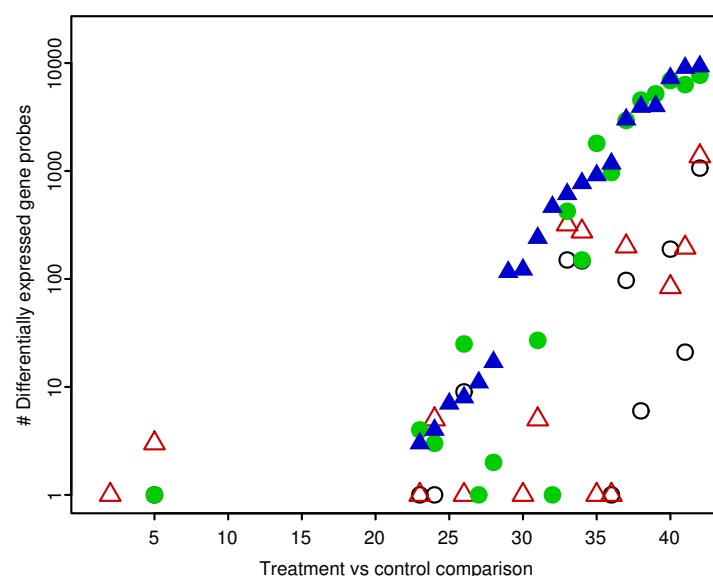
Comparison number	Treatment ID	Control ID	Treatment description	Number of DEGP
1	Ag.NM300K.EC20.7d	Ag.CTD.7d	Ag NM300K EC20 7 days	2
2	Ni.Ni.NPs.EC20.7d	Ni.CT.7d	Ni-NPs EC20 7 days	7
3	Cu.Cu.NO3.EC20.3d	Cu.CT.3d	CuNO3 EC20 3 days	27
4	Cu.Cu.NO3.EC50.7d	Cu.CT.7d	CuNO3 EC50 7 days	31
5	Cu.Cu.NPs.EC20.3d	Cu.CT.3d	Cu-NPs EC20 3 days	31
6	Ag.AgNO3.EC50.7d	Ag.CT.7d	AgNO3 EC50 7 days	33
7	Cu.Cu.NPs.EC20.7d	Cu.CT.7d	Cu-NPs EC20 7 days	33
8	Ag.NM300K.EC20.3d	Ag.CTD.3d	Ag NM300K EC20 3 days	38
9	Ag.NM300K.EC50.3d	Ag.CTD.3d	Ag NM300K EC50 3 days	52
10	Ni.NiNO3.EC20.3d	Ni.CT.3d	NiNO3 EC20 3 days	74
11	Ni.NiNO3.EC20.7d	Ni.CT.7d	NiNO3 EC20 7 days	79
12	Ag.NC.EC20.7d	Ag.CT.7d	Ag-NPs Non-Coated EC20 7 days	106
13	Ni.Ni.NPs.EC50.7d	Ni.CT.7d	Ni-NPs EC50 7 days	107
14	Ag.AgNO3.EC20.7d	Ag.CT.7d	AgNO3 EC20 7 days	111
15	Ag.NC.EC50.7d	Ag.CT.7d	Ag-NPs Non-Coated EC50 7 days	163
16	Ag.NC.EC20.3d	Ag.CT.3d	Ag-NPs Non-Coated EC20 3 days	244
17	Ag.AgNO3.EC50.3d	Ag.CT.3d	AgNO3 EC50 3 days	254
18	Ag.Coated.EC20.7d	Ag.CT.7d	Ag-NPs PVP-Coated EC20 7 days	261
19	Ni.NiNO3.EC50.7d	Ni.CT.7d	NiNO3 EC50 7 days	329
20	Cu.Cu.NPs.EC50.7d	Cu.CT.7d	Cu-NPs EC50 7 days	343
21	Ag.Coated.EC50.7d	Ag.CT.7d	Ag-NPs PVP-Coated EC50 7 days	346
22	Cu.Cu.Nwires.EC50.7d	Cu.CT.7d	Cu-NWires EC50 7 days	387
23	Cu.Cu.NO3.EC20.7d	Cu.CT.7d	CuNO3 EC20 7 days	393
24	Cu.Cu.Nwires.EC20.7d	Cu.CT.7d	Cu-NWires EC20 7 days	478
25	Cu.Cu.NO3.EC50.3d	Cu.CT.3d	CuNO3 EC50 3 days	522
26	Ag.AgNO3.EC20.3d	Ag.CT.3d	AgNO3 EC20 3 days	911
27	Ag.Coated.EC20.3d	Ag.CT.3d	Ag-NPs PVP-Coated EC20 3 days	930
28	Ag.NM300K.EC50.7d	Ag.CTD.7d	Ag NM300K EC50 7 days	1,264
29	Ni.Ni.NPs.EC20.3d	Ni.CT.3d	Ni-NPs EC20 3 days	1,460
30	Cu.Cu.field.EC20.7d	Cu.CT.7d	Cu-Field EC20 7 days	1,627
31	Ni.NiNO3.EC50.3d	Ni.CT.3d	NiNO3 EC50 3 days	1,649
32	Uv.UV.D2.5d	Uv.CT.5d	UV Dose 2	1,864
33	Ni.Ni.NPs.EC50.3d	Ni.CT.3d	Ni-NPs EC50 3 days	2,341
34	Cu.Cu.field.EC50.3d	Cu.CT.3d	Cu-Field EC50 3 days	3,578
35	Cu.Cu.field.EC50.7d	Cu.CT.7d	Cu-Field EC50 7 days	4,412
36	Uv.UV.D1.5d	Uv.CT.5d	UV Dose 1	4,746
37	Cu.Cu.NPs.EC50.3d	Cu.CT.3d	Cu-NPs EC50 3 days	5,993
38	Cu.Cu.field.EC20.3d	Cu.CT.3d	Cu-Field EC20 3 days	9,225
39	Ag.Coated.EC50.3d	Ag.CT.3d	Ag-NPs PVP-Coated EC50 3 days	9,474
40	Cu.Cu.Nwires.EC20.3d	Cu.CT.3d	Cu-NWires EC20 3 days	9,753
41	Ag.NC.EC50.3d	Ag.CT.3d	Ag-NPs Non-Coated EC50 3 days	9,876
42	Cu.Cu.Nwires.EC50.3d	Cu.CT.3d	Cu-NWires EC50 3 days	10,287

Supplementary Figures



Supplementary Figure S1

Supplementary Figure 1: Representing between-condition variation as the standard deviation of the within-condition median averages (averages of sample median expression levels, for all samples under the condition) yields similar results to those obtained with within-condition mean averages (Fig. 2). The panels show the detected variation as a function of the number of gene probes used in the between-condition normalization of the real dataset (A) and synthetic dataset (B). Labeling is the same as in Figure 2. Each point in either of the panels indicates the variation obtained with one complete normalization (black circles, median condition-decomposition normalization; blue circles, standard-vector condition-decomposition normalization). Gene probes were selected in two ways: randomly (empty circles) or in decreasing order of p -values (filled circles). Big circles show the working points of the algorithms whose results are depicted in Figures 1G–J. Black dashed lines show references for $n^{-1/2}$ decays, with the same values in both panels.



Supplementary Figure 2: With t-tests, the condition-decomposition normalizations also detected much larger numbers of differentially expressed gene probes (DEGP). The figure shows the number of DEGP obtained with a statistical analysis based on t-tests instead of limma (Fig. 3A). Labeling is the same as in Figure 3A (empty black circles, median normalization; empty red triangles, quantile normalization; filled green circles, median condition-decomposition normalization; filled blue triangles, standard-vector condition-decomposition normalization). Treatment vs control comparisons are ordered according to the number of DEGP identified with standard-vector condition-decomposition normalization, increasing from left to right. This order (not shown) was similar but not exactly the same as in Figure 3A.

Supplementary Mathematical Methods

Contents

SM1 Vectorial representation of sample data	35
SM2 Linear decomposition of the normalization problem	37
SM3 Permutation invariance of multivariate data	42
SM4 Standard-vector normalization	48
SM5 Identification of non-differentially expressed genes	50

SM1 Vectorial representation of sample data

Let x_1, \dots, x_n be the samples of n independent and identically distributed random variables X_1, \dots, X_n . Let us represent the samples x_1, \dots, x_n with the \mathbb{R}^n column vector $\mathbf{x} = (x_1, \dots, x_n)'$, and let us denote the sample mean by $\bar{x} = \sum_{i=1}^n x_i/n$.

Let us define the $\mathbb{R}^n \rightarrow \mathbb{R}^n$ vectorial operators mean ($\bar{\cdot}$) and residual ($\tilde{\cdot}$), respectively, as

$$\bar{\mathbf{x}} = (\bar{x}, \dots, \bar{x})' = \bar{x}\mathbf{1}, \quad (13)$$

$$\tilde{\mathbf{x}} = \mathbf{x} - \bar{\mathbf{x}} = \mathbf{x} - \bar{x}\mathbf{1}, \quad (14)$$

$\mathbf{1}$ being the all-ones column vector of dimension n .

Thus, any sample vector $\mathbf{x} \in \mathbb{R}^n$ can be decomposed as

$$\mathbf{x} = \bar{\mathbf{x}} + \tilde{\mathbf{x}}. \quad (15)$$

The mean vector $\bar{\mathbf{x}}$ contains the sample mean, while the residual vector $\tilde{\mathbf{x}}$ carries the sample variation around the mean.

The vectorial operators mean (13) and residual (14) are linear.

Proposition. For any two sample vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and any two numbers $\alpha, \beta \in \mathbb{R}$,

$$\overline{\alpha\mathbf{x} + \beta\mathbf{y}} = \alpha\bar{\mathbf{x}} + \beta\bar{\mathbf{y}}, \quad (16)$$

$$\widetilde{\alpha\mathbf{x} + \beta\mathbf{y}} = \alpha\tilde{\mathbf{x}} + \beta\tilde{\mathbf{y}}. \quad (17)$$

Proof. Let us denote $\mathbf{x} = (x_1, \dots, x_n)'$ and $\mathbf{y} = (y_1, \dots, y_n)'$.

$$\begin{aligned} \overline{\alpha\mathbf{x} + \beta\mathbf{y}} &= \frac{\sum_{i=1}^n (\alpha x_i + \beta y_i)}{n} \mathbf{1} = \alpha \frac{\sum_{i=1}^n x_i}{n} \mathbf{1} + \beta \frac{\sum_{i=1}^n y_i}{n} \mathbf{1} = \alpha\bar{\mathbf{x}} + \beta\bar{\mathbf{y}}, \\ \widetilde{\alpha\mathbf{x} + \beta\mathbf{y}} &= \alpha\mathbf{x} + \beta\mathbf{y} - \overline{\alpha\mathbf{x} + \beta\mathbf{y}} = \alpha\mathbf{x} + \beta\mathbf{y} - (\alpha\bar{\mathbf{x}} + \beta\bar{\mathbf{y}}), \\ &= \alpha(\mathbf{x} - \bar{\mathbf{x}}) + \beta(\mathbf{y} - \bar{\mathbf{y}}) = \alpha\tilde{\mathbf{x}} + \beta\tilde{\mathbf{y}}. \quad \square \end{aligned}$$

An essential property of the mean and residual vectors is that they belong to subspaces that are orthogonal complements (Eaton, 2007). Hence, for any sample vector $\mathbf{x} \in \mathbb{R}^n$, the mean vector $\bar{\mathbf{x}}$ belongs to the subspace of dimension 1 spanned by the unit vector $\hat{\mathbf{1}} = \mathbf{1}/\sqrt{n}$, while the residual vector $\tilde{\mathbf{x}}$ belongs to the $(n-1)$ -dimensional hyperplane orthogonal to $\hat{\mathbf{1}}$.

The lengths of the mean vector and residual vector are equal, up to a scaling factor, to the sample mean and sample standard deviation, respectively. For a set of samples x_1, \dots, x_n , where $n \geq 2$, let us denote the sample mean as before by $\bar{x} = \sum_{i=1}^n x_i/n$, and the sample variance as $s_x^2 = \sum_{i=1}^n (x_i - \bar{x})^2/(n-1)$. Then, the lengths of the mean and residual vectors obtained from the sample vector $\mathbf{x} = (x_1, \dots, x_n)'$ are

$$\|\bar{\mathbf{x}}\| = \sqrt{n \bar{x}^2} = \sqrt{n} |\bar{x}|, \quad (18)$$

$$\|\tilde{\mathbf{x}}\| = \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{n-1} s_x. \quad (19)$$

Finally, let us define the standard vector of the sample vector $\mathbf{x} = (x_1, \dots, x_n)'$ ($n \geq 2$), as

$$\text{stdvec}(\mathbf{x}) = \sqrt{n-1} \frac{\tilde{\mathbf{x}}}{\|\tilde{\mathbf{x}}\|}, \quad (20)$$

whenever $\tilde{\mathbf{x}} \neq \mathbf{0}$, or otherwise as $\text{stdvec}(\mathbf{x}) = \mathbf{0}$. $\mathbf{0}$ is the all-zeros column vector of dimension n .

For a given number of samples n , all the non-zero standard vectors belong to the $(n-2)$ -sphere of radius $\sqrt{n-1}$, embedded in the $(n-1)$ -dimensional hyperplane perpendicular to $\hat{\mathbf{1}}$. Besides, all the components of a standard vector are equal to the corresponding standardized samples,

$$\sqrt{n-1} \frac{\tilde{x}_i}{\|\tilde{\mathbf{x}}\|} = \frac{x_i - \bar{x}}{s_x}. \quad (21)$$

For the degenerate case of having only two samples ($n=2$), the only possible values of a non-zero standard vector are $\pm(1/\sqrt{2}, -1/\sqrt{2})'$.

SM2 Linear decomposition of the normalization problem

Let us consider a gene expression dataset, with g genes and c experimental conditions. Each condition k has s_k samples. The total number of samples is $s = \sum_{k=1}^c s_k$.

Let us denote the *observed* expression level of gene j in the sample i of condition k by $y_{ij}^{(k)}$. We assume that the observed level $y_{ij}^{(k)}$ is equal, in the usual \log_2 -scale, to the addition of the normalization factor $a_i^{(k)}$ to the *true* gene expression level $x_{ij}^{(k)}$,

$$y_{ij}^{(k)} = x_{ij}^{(k)} + a_i^{(k)}. \quad (22)$$

Solving the *normalization problem* amounts to finding the normalization factors $a_i^{(k)}$ from the observed values $y_{ij}^{(k)}$. The normalization factors can be understood as sample-wide changes in the concentration of mRNA molecules by multiplicative factors equal to $2^{a_i^{(k)}}$. These changes are caused by technical reasons in the assay and are independent of the biological variation in the true levels $x_{ij}^{(k)}$.

Let us represent the true and observed expression levels, $x_{ij}^{(k)}$ and $y_{ij}^{(k)}$, of gene j in the

samples $i = 1 \dots s_k$ of condition k , by the s_k -dimensional vectors

$$\mathbf{x}_j^{(k)} = (x_{1j}^{(k)}, \dots, x_{s_k j}^{(k)})', \quad (23)$$

$$\mathbf{y}_j^{(k)} = (y_{1j}^{(k)}, \dots, y_{s_k j}^{(k)})'. \quad (24)$$

Let us also represent the unknown normalization factors of condition k by the s_k -dimensional vector

$$\mathbf{a}^{(k)} = (a_1^{(k)}, \dots, a_{s_k}^{(k)})'. \quad (25)$$

From (22)–(25), the normalization problem can be written in vectorial form as

$$\mathbf{y}_j^{(k)} = \mathbf{x}_j^{(k)} + \mathbf{a}^{(k)}. \quad (26)$$

Applying the vectorial operators mean (13) and residual (14), we obtain

$$\bar{\mathbf{y}}_j^{(k)} = \bar{\mathbf{x}}_j^{(k)} + \bar{\mathbf{a}}^{(k)}, \quad (27)$$

$$\tilde{\mathbf{y}}_j^{(k)} = \tilde{\mathbf{x}}_j^{(k)} + \tilde{\mathbf{a}}^{(k)}. \quad (28)$$

The residual-vector equations (28) correspond to the c within-condition normalizations. Each within-condition normalization uses the equations (28) particular to a condition k , for the subset of genes $\mathcal{G}_k \subseteq \{1, \dots, g\}$ that have expression level available and of enough quality in that experimental condition.

Let us denote the condition means for each gene as

$$\bar{x}_j^{(k)} = \frac{\sum_{i=1}^{s_k} x_{ij}^{(k)}}{s_k}, \quad (29)$$

$$\bar{y}_j^{(k)} = \frac{\sum_{i=1}^{s_k} y_{ij}^{(k)}}{s_k}, \quad (30)$$

$$\bar{a}^{(k)} = \frac{\sum_{i=1}^{s_k} a_i^{(k)}}{s_k}, \quad (31)$$

so that

$$\bar{\mathbf{x}}_j^{(k)} = \bar{x}_j^{(k)} \mathbf{1}_{s_k}, \quad (32)$$

$$\bar{\mathbf{y}}_j^{(k)} = \bar{y}_j^{(k)} \mathbf{1}_{s_k}, \quad (33)$$

$$\bar{\mathbf{a}}^{(k)} = \bar{a}^{(k)} \mathbf{1}_{s_k}, \quad (34)$$

$\mathbf{1}_{s_k}$ being the all-ones column vector of dimension s_k .

Then, the mean-vector equations (27) can be written as

$$\bar{y}_j^{(k)} \mathbf{1}_{s_k} = \bar{x}_j^{(k)} \mathbf{1}_{s_k} + \bar{a}^{(k)} \mathbf{1}_{s_k}, \quad (35)$$

so they reduce to the scalar equations

$$\bar{y}_j^{(k)} = \bar{x}_j^{(k)} + \bar{a}^{(k)}. \quad (36)$$

Let us define the vectors of conditions means as

$$\mathbf{x}_j^* = (\bar{x}_j^{(1)}, \dots, \bar{x}_j^{(c)})', \quad (37)$$

$$\mathbf{y}_j^* = (\bar{y}_j^{(1)}, \dots, \bar{y}_j^{(c)})', \quad (38)$$

$$\mathbf{a}^* = (\bar{a}^{(1)}, \dots, \bar{a}^{(c)})', \quad (39)$$

and let us express the condition-mean equations in vectorial form as

$$\mathbf{y}_j^* = \mathbf{x}_j^* + \mathbf{a}^*. \quad (40)$$

Applying again the mean and variance operators, we obtain

$$\bar{\mathbf{y}}_j^* = \bar{\mathbf{x}}_j^* + \bar{\mathbf{a}}^*, \quad (41)$$

$$\tilde{\mathbf{y}}_j^* = \tilde{\mathbf{x}}_j^* + \tilde{\mathbf{a}}^*. \quad (42)$$

The residual-vector equations on condition means (42) correspond to the single between-condition normalization, in a similar way as (28) do for the each of the within-condition normalizations. There is one equation (42) per gene. The only equations used in the between-condition normalization are those of the subset of genes $\mathcal{G}^* \subseteq \{1, \dots, g\}$ that show no evidence of variation across experimental conditions, according to a statistical test.

Given that $\bar{\mathbf{a}}^* = \bar{a}^* \mathbf{1}_c$, (41) has the only unknown \bar{a}^* . The meaning of \bar{a}^* is a conversion factor between the scale the true and observed expression levels. This factor depends on

the technology used to measure the expression levels and finding it is out of the scope of the normalization problem. Therefore, without loss of generality, we assume $\bar{\mathbf{a}}^* = \mathbf{0}$, so

$$\bar{\mathbf{a}}^* = \mathbf{0}_c, \quad (43)$$

$$\mathbf{a}^* = \tilde{\mathbf{a}}^*. \quad (44)$$

The solution of the between-condition normalization, $\tilde{\mathbf{a}}^*$, allows to find the mean vectors of the normalization factors $\bar{\mathbf{a}}^{(k)}$, via (34), (39) and (44). The within-condition normalizations yield the residual vectors $\tilde{\mathbf{a}}^{(k)}$. The complete solution to the normalization problem is finally obtained, for each condition k , with

$$\mathbf{a}^{(k)} = \bar{\mathbf{a}}^{(k)} + \tilde{\mathbf{a}}^{(k)}. \quad (45)$$

Thus, the original normalization problem (26) has been divided in $c+1$ normalization sub-problems on residual vectors, stated by (28) and (42). In fact, this linear decomposition is possible for any partition of the set of s samples. The choice of the partition as the one defined by the experimental conditions is motivated by the need to control the biological variation among the genes used in each normalization. All the $c+1$ normalizations face the same kind of *normalization of residuals problem*, which we define in general as follows.

Normalization of Residuals Problem. Let y_{ij} be the i -th observed value of feature j , in a dataset with $n \geq 2$ observations for each of the m features. The observed values y_{ij} are equal to the true values x_{ij} plus the normalization factors a_i , which are constant across features. In vectorial form, there are m equations

$$\mathbf{y}_j = \mathbf{x}_j + \mathbf{a}, \quad (46)$$

where the vectors belong to \mathbb{R}^n . As a consequence

$$\tilde{\mathbf{y}}_j = \tilde{\mathbf{x}}_j + \tilde{\mathbf{a}}. \quad (47)$$

Solving the *normalization of residuals problem* amounts to finding the residual vector of normalization factors $\tilde{\mathbf{a}}$ from the observed residual vectors $\tilde{\mathbf{y}}_j$. In the within-condition

normalizations, the features are gene expression levels, with one observation per sample of the corresponding experimental condition. In the between-condition normalization, the features are means of gene expression levels, with one observation per condition.

There is, however, an additional requirement imposed by the methods with which we propose to solve the between-condition normalization. We would like to consider the condition means $\bar{x}_j^{(k)}$ in (36) as sample data across conditions. This only holds when all the conditions have the same number of samples. Otherwise, we balance the condition means so that they result from the same number of samples in all conditions, according to the procedure described in the following.

Let s^* be the minimum number of samples across conditions, $s^* = \min\{s_1, \dots, s_c\}$. Let $\mathcal{S}_j^{(k)}$ be independent random samples (without replacement) of size s^* from the set of indexes $\{1, \dots, s_k\}$, with one sample per gene j and condition k . Then, the balanced condition means are defined as

$$\bar{x}_j^{(k)*} = \frac{\sum_{i \in \mathcal{S}_j^{(k)}} x_{ij}^{(k)}}{s^*}, \quad (48)$$

$$\bar{y}_j^{(k)*} = \frac{\sum_{i \in \mathcal{S}_j^{(k)}} y_{ij}^{(k)}}{s^*}, \quad (49)$$

$$\bar{a}_j^{(k)*} = \frac{\sum_{i \in \mathcal{S}_j^{(k)}} a_i^{(k)}}{s^*}. \quad (50)$$

From (22), the balanced condition means verify a relationship similar to (36),

$$\bar{y}_j^{(k)*} = \bar{x}_j^{(k)*} + \bar{a}_j^{(k)*}. \quad (51)$$

Moreover, the average of $\bar{a}_j^{(k)*}$ across the sampling subsets $\mathcal{S}_j^{(k)}$ is equal to the unknown $\bar{a}^{(k)}$. This implies that (51) are, on average, equivalent to (36). Hence, we use the following vectors of balanced conditions means

$$\mathbf{x}_j^* = (\bar{x}_j^{(1)*}, \dots, \bar{x}_j^{(c)*}), \quad (52)$$

$$\mathbf{y}_j^* = (\bar{y}_j^{(1)*}, \dots, \bar{y}_j^{(c)*}), \quad (53)$$

instead of (37), (38), in order to build the condition-mean equations (40). This balancing of the condition means is only required when the experimental conditions have different number of samples.

SM3 Permutation invariance of multivariate data

Let x_{ij} and y_{ij} be, respectively, the true and observed values of a dataset with n observations of m features, as defined in the *normalization of residuals problem* above.

We have assumed that the n true values x_{1j}, \dots, x_{nj} of feature j are samples of independent and identically distributed random variables X_{1j}, \dots, X_{nj} . These random variables can be represented with the random vector $\mathbf{X}_j = (X_{1j}, \dots, X_{nj})'$, carried by the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and with induced space $(\mathbb{R}^n, \mathbb{B}^n, \mathbb{P})$. Let us define the random vectors $\bar{\mathbf{X}}_j$ and $\tilde{\mathbf{X}}_j$ with the vectorial operators mean (13) and residual (14), respectively,

$$\bar{X}_j = \sum_{i=1}^n \frac{X_{ij}}{n}, \quad (54)$$

$$\bar{\mathbf{X}}_j = (\bar{X}_j, \dots, \bar{X}_j)' = \bar{X}_j \mathbf{1}, \quad (55)$$

$$\tilde{\mathbf{X}}_j = \mathbf{X}_j - \bar{\mathbf{X}}_j = \mathbf{X}_j - \bar{X}_j \mathbf{1}. \quad (56)$$

$\mathbf{X}_j = \bar{\mathbf{X}}_j + \tilde{\mathbf{X}}_j$ holds for any random vector \mathbf{X}_j , as well as the other properties presented above. Let us assume that $E(\|\mathbf{X}_j\|) < \infty$ and that $P(\|\tilde{\mathbf{X}}_j\| = 0) = 0$, which imply that $\tilde{\mathbf{X}}_j/\|\tilde{\mathbf{X}}_j\|$ has length 1 almost surely.

The standard random vector $\sqrt{n-1} \tilde{\mathbf{X}}_j/\|\tilde{\mathbf{X}}_j\|$ is a pivotal quantity, where the location (mean) and scale (standard deviation) of feature j have been removed. The probability distribution of $\tilde{\mathbf{X}}_j/\|\tilde{\mathbf{X}}_j\|$ across the remaining degrees of freedom over the unit $(n-2)$ -sphere is governed by the parametric family of the random variables X_{1j}, \dots, X_{nj} . Moreover, the independence and identity of distribution across the n observations implies that the distribution of \mathbf{X}_j is *exchangeable*, i.e. invariant with respect to permutations of the observation labels. As a result, $\tilde{\mathbf{X}}_j/\|\tilde{\mathbf{X}}_j\|$ is also permutation invariant, which geometrically corresponds to symmetries with respect to the $n!$ permutations of the axes

in the n -dimensional space of random vectors, projected onto the $(n - 1)$ -dimensional hyperplane of residual vectors.

Residual vectors and standard vectors have been widely studied, especially in relation to elliptically symmetric distributions and linear models (Fang et al., 1990; Gupta et al., 2013), and to the invariances of probability distributions (Kallenberg, 2005). Here, we consider these vectors from the viewpoint of the problem of normalizing multivariate data, and its relationship with permutation invariance.

It is well known that, for a multivariate distribution with independent and identically distributed components, the expected value of the standard vector is zero (Eaton, 2007), given that it is so for each component. We prove this here for completeness, and to show that it is also a necessary consequence of the permutation invariance of the distribution.

Proposition. The expected value of any true (i.e. without normalization issues) standard vector is zero. If the $n \geq 2$ samples of feature j are independent and identically distributed, then

$$\mathbb{E} \left(\sqrt{n-1} \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \right) = \mathbf{0}. \quad (57)$$

Proof. Let \mathcal{P}_n be the set of all the permutation matrices in $\mathbb{R}^{n \times n}$. Then, for any $P \in \mathcal{P}_n$, $\tilde{\mathbf{X}}_j / \|\tilde{\mathbf{X}}_j\|$ is equal in distribution to $P \tilde{\mathbf{X}}_j / \|\tilde{\mathbf{X}}_j\|$. This implies that

$$\mathbb{E} \left(\frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \right) = \mathbb{E} \left(P \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \right) = P \mathbb{E} \left(\frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \right).$$

The only vectors that are invariant with respect to all possible permutations are those that have all components identical. Therefore, $\mathbb{E}(\tilde{\mathbf{X}}_j / \|\tilde{\mathbf{X}}_j\|) = \alpha \hat{\mathbf{1}}$, with $\alpha \in \mathbb{R}$. However, $\tilde{\mathbf{X}}_j' \hat{\mathbf{1}} = 0$, so that $\alpha = \mathbb{E}(\tilde{\mathbf{X}}_j / \|\tilde{\mathbf{X}}_j\|)' \hat{\mathbf{1}} = 0$. Hence $\mathbb{E}(\tilde{\mathbf{X}}_j / \|\tilde{\mathbf{X}}_j\|) = \mathbf{0}$. \square

For each true random vector \mathbf{X}_j , there is an observed random vector $\mathbf{Y}_j = \mathbf{X}_j + \mathbf{A}$, where \mathbf{A} is the random vector of normalization factors. The random vectors \mathbf{X}_j and \mathbf{A} are independent, representing biological and technical variation, respectively. Therefore, and without loss of generality, we assume in what follows a fixed vector of normalization

factors \mathbf{a} , i.e. we condition on the event $\{\mathbf{A} = \mathbf{a}\}$. We also assume that $P(\|\tilde{\mathbf{Y}}_j\| = 0) = 0$, which implies that $\tilde{\mathbf{Y}}_j/\|\tilde{\mathbf{Y}}_j\|$ has length 1 almost surely.

In contrast to the true standard vector $\sqrt{n-1}\tilde{\mathbf{X}}_j/\|\tilde{\mathbf{X}}_j\|$, the observed standard vector $\sqrt{n-1}\tilde{\mathbf{Y}}_j/\|\tilde{\mathbf{Y}}_j\|$ is biased toward the direction of $\tilde{\mathbf{a}}$, with the result that the expected value is not zero.

Proposition. If the $n \geq 2$ samples of feature j are independent and identically distributed, whenever $\tilde{\mathbf{a}} \neq \mathbf{0}$,

$$E\left(\sqrt{n-1}\frac{\tilde{\mathbf{Y}}_j}{\|\tilde{\mathbf{Y}}_j\|}\right) \neq \mathbf{0}. \quad (58)$$

When $n = 2$, there is the additional requirement that $P(\|\tilde{\mathbf{X}}_i\| < \|\tilde{\mathbf{a}}\|) > 0$. This threshold of detection only occurs for the degenerate case of $n = 2$.

Proof. Let us consider the projection of $\tilde{\mathbf{Y}}_j/\|\tilde{\mathbf{Y}}_j\|$ on $\tilde{\mathbf{a}}$, compared to the projection of $\tilde{\mathbf{X}}_j/\|\tilde{\mathbf{X}}_j\|$.

When the vectors $\tilde{\mathbf{X}}_j$ and $\tilde{\mathbf{a}}$ are collinear,

$$\frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} = \pm 1, \quad \text{and} \quad \frac{\tilde{\mathbf{Y}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{Y}}_j\| \|\tilde{\mathbf{a}}\|} = \pm 1,$$

with

$$\frac{\tilde{\mathbf{Y}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{Y}}_j\| \|\tilde{\mathbf{a}}\|} \geq \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|}.$$

This is the only case when $n = 2$. The additional requirement ensures that, for $n = 2$,

$$P\left(\frac{\tilde{\mathbf{Y}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{Y}}_j\| \|\tilde{\mathbf{a}}\|} > \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|}\right) > 0,$$

which implies

$$E\left(\frac{\tilde{\mathbf{Y}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{Y}}_j\| \|\tilde{\mathbf{a}}\|}\right) > E\left(\frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|}\right).$$

Otherwise, when $n > 2$ and the vectors $\tilde{\mathbf{X}}_j$ and $\tilde{\mathbf{a}}$ are not collinear, they lie on a plane. The vector $\tilde{\mathbf{Y}}_j = \tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}$ is the diagonal of the parallelogram defined by $\tilde{\mathbf{X}}_j$ and $\tilde{\mathbf{a}}$. Hence

the angle between $\tilde{\mathbf{Y}}_j$ and $\tilde{\mathbf{a}}$ is strictly less than the angle between $\tilde{\mathbf{X}}_j$ and $\tilde{\mathbf{a}}$, so the cosine of the angle is strictly greater. Thus,

$$\frac{\tilde{\mathbf{Y}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{Y}}_j\| \|\tilde{\mathbf{a}}\|} > \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|}.$$

Due to the permutation symmetries in the distribution of $\tilde{\mathbf{X}}_j/\|\tilde{\mathbf{X}}_j\|$, when $n > 2$ the vector $\tilde{\mathbf{X}}_j$ has non-zero probability of being not collinear with $\tilde{\mathbf{a}}$, i.e. $P(|\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}| < 1) > 0$.

Therefore,

$$P\left(\frac{\tilde{\mathbf{Y}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{Y}}_j\| \|\tilde{\mathbf{a}}\|} > \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|}\right) > 0,$$

which again implies

$$E\left(\frac{\tilde{\mathbf{Y}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{Y}}_j\| \|\tilde{\mathbf{a}}\|}\right) > E\left(\frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|}\right).$$

Finally,

$$\left\| E\left(\frac{\tilde{\mathbf{Y}}_j}{\|\tilde{\mathbf{Y}}_j\|}\right) \right\| \geq E\left(\frac{\tilde{\mathbf{Y}}_j}{\|\tilde{\mathbf{Y}}_j\|}\right)' \frac{\tilde{\mathbf{a}}}{\|\tilde{\mathbf{a}}\|} > E\left(\frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|}\right)' \frac{\tilde{\mathbf{a}}}{\|\tilde{\mathbf{a}}\|} = 0. \quad \square$$

As a consequence, the *normalization of residuals problem* may be restated as the problem of finding the normalization factors $\tilde{\mathbf{a}}$ from the observed vectors $\tilde{\mathbf{y}}_j$, such that the standard vectors $\sqrt{n-1}(\tilde{\mathbf{y}}_j - \tilde{\mathbf{a}})/\|\tilde{\mathbf{y}}_j - \tilde{\mathbf{a}}\|$ are invariant against permutations of the observation labels. Or equivalently, such that the standard vectors $\sqrt{n-1}(\tilde{\mathbf{y}}_j - \tilde{\mathbf{a}})/\|\tilde{\mathbf{y}}_j - \tilde{\mathbf{a}}\|$ have zero mean. The following property provides an approach to the solution.

Proposition. Whenever $\tilde{\mathbf{a}} \neq \mathbf{0}$, the component of the expected value of $\tilde{\mathbf{Y}}_j/\|\tilde{\mathbf{Y}}_j\|$ parallel to $\tilde{\mathbf{a}}$ verifies

$$0 < E\left(\frac{\tilde{\mathbf{Y}}_j}{\|\tilde{\mathbf{Y}}_j\|}\right)' \frac{\tilde{\mathbf{a}}}{\|\tilde{\mathbf{a}}\|} < E\left(\frac{1}{\|\tilde{\mathbf{Y}}_j\|}\right) \|\tilde{\mathbf{a}}\|. \quad (59)$$

As in (58), when $n = 2$ we also assume that $P(\|\tilde{\mathbf{X}}_j\| < \|\tilde{\mathbf{a}}\|) > 0$.

Proof. The first inequality holds from the previous proof. Concerning the second inequality, let us consider

$$\frac{\tilde{\mathbf{Y}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{Y}}_j\| \|\tilde{\mathbf{a}}\|} = \frac{(\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}})' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\| \|\tilde{\mathbf{a}}\|} = \frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} + \frac{\|\tilde{\mathbf{a}}\|}{\|\tilde{\mathbf{Y}}_j\|}.$$

We need to prove that the first term on the RHS has negative expected value. Let us decompose this term into the positive and negative parts,

$$\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} = \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+ - \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^-,$$

where $X^+ = \max(X, 0)$ and $X^- = -\min(X, 0)$.

Because $\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|^2 = \|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2 + 2\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}$,

$$\begin{aligned} \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+ &\leq \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+, \\ \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^- &\geq \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^-. \end{aligned}$$

These inequalities are identities when $\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}$ is of opposite sign to $(\cdot)^\pm$, or when $\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}} = 0$.

Because of the permutation symmetries of $\tilde{\mathbf{X}}_j/\|\tilde{\mathbf{X}}_j\|$, it follows that $P(\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}} \neq 0) > 0$,

which implies

$$\begin{aligned} P \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+ < \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+ \right) &> 0, \\ P \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^- > \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^- \right) &> 0, \end{aligned}$$

and hence

$$\begin{aligned} E \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+ \right) &< E \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+ \right), \\ E \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^- \right) &> E \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^- \right). \end{aligned}$$

For any permutation matrix $P \in \mathcal{P}_n$,

$$\begin{aligned} \frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} &= \frac{\|P \tilde{\mathbf{X}}_j\|}{\sqrt{\|P \tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \quad \text{surely,} \\ \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} &= P \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \quad \text{in distribution,} \end{aligned}$$

so that

$$\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} = P \frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \quad \text{in distribution,}$$

which together with

$$\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \right)' \hat{\mathbf{1}} = 0 \quad \text{surely,}$$

implies, as in (57), that

$$\mathbb{E} \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j}{\|\tilde{\mathbf{X}}_j\|} \right) = \mathbf{0}.$$

Therefore,

$$\mathbb{E} \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+ \right) = \mathbb{E} \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\sqrt{\|\tilde{\mathbf{X}}_j\|^2 + \|\tilde{\mathbf{a}}\|^2}} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^- \right).$$

Back to the initial expected values, it follows that

$$\mathbb{E} \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^+ \right) < \mathbb{E} \left(\left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right)^- \right),$$

which implies

$$\mathbb{E} \left(\frac{\|\tilde{\mathbf{X}}_j\|}{\|\tilde{\mathbf{X}}_j + \tilde{\mathbf{a}}\|} \frac{\tilde{\mathbf{X}}_j' \tilde{\mathbf{a}}}{\|\tilde{\mathbf{X}}_j\| \|\tilde{\mathbf{a}}\|} \right) < 0. \quad \square$$

The Gaussian multivariate distribution, among others, has spherical symmetry besides permutation symmetry. For parametric families with spherical symmetry, the true standard vector $\sqrt{n-1} \tilde{\mathbf{X}}_j / \|\tilde{\mathbf{X}}_j\|$ has uniform distribution over the $(n-2)$ -sphere. As a result, the components of $\tilde{\mathbf{Y}}_j / \|\tilde{\mathbf{Y}}_j\|$ perpendicular to $\tilde{\mathbf{a}}$ are antisymmetric with respect to the direction of $\tilde{\mathbf{a}}$, so that they cancel out in expectation. That is, for parametric families with spherical symmetry, and as long as $\tilde{\mathbf{a}} \neq \mathbf{0}$,

$$\mathbb{E} \left(\frac{\tilde{\mathbf{Y}}_j}{\|\tilde{\mathbf{Y}}_j\|} \right) = \lambda \tilde{\mathbf{a}}, \quad \text{with} \quad 0 < \lambda < \mathbb{E} \left(\frac{1}{\|\tilde{\mathbf{Y}}_j\|} \right). \quad (60)$$

SM4 Standard-vector normalization

The properties (59), (60) suggest the use of

$$\hat{\mathbf{b}} = \frac{\sum_{j=1}^m \frac{\tilde{\mathbf{y}}_j}{\|\tilde{\mathbf{y}}_j\|}}{\sum_{j=1}^m \frac{1}{\|\tilde{\mathbf{y}}_j\|}} \quad (61)$$

to approximate the unknown residual vector of normalization factors $\tilde{\mathbf{a}}$. The following iterative method implements this approach to solve the *normalization of residuals problem*.

Let us define the following recursive sequence, where each step t comprises m vectors $\hat{\mathbf{y}}_j^{(t)}$ ($j \in \{1, \dots, m\}$) and one vector $\hat{\mathbf{b}}^{(t)}$,

$$\hat{\mathbf{y}}_j^{(0)} = \tilde{\mathbf{y}}_j, \quad (62)$$

$$\hat{\mathbf{y}}_j^{(t)} = \hat{\mathbf{y}}_j^{(t-1)} - \hat{\mathbf{b}}^{(t-1)}, \quad \text{for } t \geq 1, \quad (63)$$

$$\hat{\mathbf{b}}^{(t)} = \frac{\sum_{j=1}^m \frac{\hat{\mathbf{y}}_j^{(t)}}{\|\hat{\mathbf{y}}_j^{(t)}\|}}{\sum_{j=1}^m \frac{1}{\|\hat{\mathbf{y}}_j^{(t)}\|}}, \quad \text{for } t \geq 0. \quad (64)$$

We assume that $\hat{\mathbf{y}}_j^{(t)} \neq \mathbf{0}_n$, for all $j \in \{1, \dots, m\}$ and all $t \geq 0$. Nonetheless, an implementation of this algorithm benefits from trimming out a small fraction (e.g. 1%) of the features with lesser $\|\hat{\mathbf{y}}_j^{(t)}\|$ in (64), in order to avoid numerical singularities.

Let us write $\hat{\mathbf{y}}_j^{(t)}$ as a function of the unknowns $\tilde{\mathbf{x}}_j$ and $\tilde{\mathbf{a}}$. For any $t \geq 1$,

$$\hat{\mathbf{y}}_j^{(t)} = \hat{\mathbf{y}}_j^{(t-1)} - \hat{\mathbf{b}}^{(t-1)}, \quad (65)$$

$$= \hat{\mathbf{y}}_j^{(t-2)} - \hat{\mathbf{b}}^{(t-2)} - \hat{\mathbf{b}}^{(t-1)}, \quad (66)$$

$$\vdots \quad (67)$$

$$= \hat{\mathbf{y}}_j^{(0)} - \sum_{r=0}^{t-1} \hat{\mathbf{b}}^{(r)}, \quad (68)$$

$$= \tilde{\mathbf{y}}_j - \sum_{r=0}^{t-1} \hat{\mathbf{b}}^{(r)}, \quad (69)$$

$$= \tilde{\mathbf{x}}_j + \tilde{\mathbf{a}} - \sum_{r=0}^{t-1} \hat{\mathbf{b}}^{(r)}. \quad (70)$$

Note that (70) is also valid for $t = 0$.

Let us also define the vectors $\hat{\mathbf{a}}^{(t)}$, for $t \geq 0$, which describe the vector of normalization factors still to be removed at step t ,

$$\hat{\mathbf{a}}^{(t)} = \tilde{\mathbf{a}} - \sum_{r=0}^{t-1} \hat{\mathbf{b}}^{(r)}, \quad (71)$$

so that, by (70), for $t \geq 0$,

$$\hat{\mathbf{y}}_j^{(t)} = \tilde{\mathbf{x}}_j + \hat{\mathbf{a}}^{(t)}. \quad (72)$$

Therefore, the recursive sequence (62)–(64) faces a new, weaker *normalization of residuals problem* at each step t , with true residual vectors $\tilde{\mathbf{x}}_j$, observed residual vectors $\hat{\mathbf{y}}_j^{(t)}$ and unknown normalization factors $\hat{\mathbf{a}}^{(t)}$. The step t results in the estimation of normalization factors $\hat{\mathbf{b}}^{(t)}$, which are removed from $\hat{\mathbf{y}}_j^{(t)}$, generating the next step. At the beginning, $\hat{\mathbf{y}}_j^{(0)} = \tilde{\mathbf{y}}_j$ and $\hat{\mathbf{a}}^{(0)} = \tilde{\mathbf{a}}$.

At convergence, $\lim_{t \rightarrow \infty} \hat{\mathbf{b}}^{(t)} = \mathbf{0}$. Equations (57), (58), (64) imply that, in such a case, $\lim_{t \rightarrow \infty} \hat{\mathbf{y}}_j^{(t)} = \tilde{\mathbf{x}}_j$ and $\sum_{t=0}^{\infty} \hat{\mathbf{b}}^{(t)} = \tilde{\mathbf{a}}$. Convergence is optimal when the parametric family of the m features has spherical symmetry, Gaussian being the most prominent case. Otherwise, the more uniform the distribution of standard vectors $\sqrt{n-1} \tilde{\mathbf{x}}_j / \|\tilde{\mathbf{x}}_j\|$ is on the $(n-2)$ -sphere, the faster the sequence (62)–(64) converges. See examples of convergence in Supplementary Movies 1–3.

SM5 Identification of non-differentially expressed genes

Let us consider a gene expression dataset, with g genes and c experimental conditions. Each condition k has s_k samples. The total number of samples is $s = \sum_{k=1}^c s_k$. Let us assume that $c \geq 2$ and that $s_k \geq 2$, for all conditions $k \in \{1, \dots, c\}$. Let us also assume that, among the g genes, there is a fraction π_0 of non-differentially expressed (non-DE) genes, with $0 \leq \pi_0 \leq 1$, while the remaining fraction $1 - \pi_0$ comprises the differentially expressed (DE) genes (Storey and Tibshirani, 2003).

Let us consider the usual ANOVA test comparing average expression levels across conditions, gene-by-gene. Under the null hypothesis of a non-differentially expressed gene, the corresponding F -statistic follows the F -distribution with $c - 1$ and $s - c$ degrees of freedom. The test of this hypothesis yields a p -value p_j for each gene $j \in \{1, \dots, g\}$. The obtained p -values p_j follow a probability distribution that can be considered as the mixture of two probability distributions, F_0 and F_1 , for the non-DE genes and the DE genes, respectively (Storey, 2003). The fraction π_0 of non-DE genes follows the uniform distribution on the interval $[0, 1]$,

$$F_0(p) = p, \quad (73)$$

while the fraction $1 - \pi_0$ of DE genes follows a distribution that verifies, for any $p \in (0, 1)$,

$$F_1(p) > p, \quad (74)$$

and the mixture distribution is

$$F(p) = \pi_0 F_0(p) + (1 - \pi_0) F_1(p). \quad (75)$$

Let us further assume that there exists a p^* , with $0 < p^* < 1$, such that $F_1(p) = 1$ for every $p \geq p^*$. In other words, all DE genes have p -value p_j from the ANOVA test such that $p_j \leq p^*$, while only some genes among the non-DE genes have p -value with $p_j > p^*$. This implies that the mixture distribution of p -values is uniform on the interval $[p^*, 1]$,

$$F(p) = \pi_0 p + 1 - \pi_0, \quad \text{for } p^* \leq p \leq 1, \quad (76)$$

$$f(p) = \pi_0, \quad \text{for } p^* < p < 1. \quad (77)$$

On the other hand, for any set of n samples $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ obtained from n independent and identically distributed uniform random variables on the interval $[a, b]$, all the distances between consecutive ordered samples (including boundaries), $x_{(1)} - a, x_{(2)} - x_{(1)}, \dots, x_{(n)} - x_{(n-1)}, b - x_{(n)}$, obey the same distribution (Feller, 1971). Then, it can be realized that, for any j such that $2 \leq j \leq n - 1$, the two subsets of samples $x_{(1)}, \dots, x_{(j-1)}$ and $x_{(j+1)}, \dots, x_{(n)}$ follow uniform distributions on the intervals $[a, x_{(j)}]$ and $[x_{(j)}, b]$, respectively.

Based on these facts, to identify non-DE genes we propose finding the minimum $p_{(j)}$, from the ordered sequence of p -values $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(g)}$, such that a goodness-of-fit test for the uniform distribution on the interval $[p_{(j)}, 1]$, performed on $p_{(j+1)}, \dots, p_{(g)}$, is not rejected. As a result, the genes corresponding to the p -values $p_{(j)}, p_{(j+1)}, \dots, p_{(g)}$ are considered as non-DE genes.

Given the concavity of $F(p)$, the goodness-of-fit test used is the one-sided Kolmogorov-Smirnov test on positive deviations of the empirical distribution function.

See Supplementary Movies 4–5 for examples of this approach to identifying non-differentially expressed genes.