

<https://www.researchgate.net/publication/229099731>—“Logistic-Normal—“verb—“Distributions—“Some—“Prop
some-uses-andfl—proptfl—“verb-ies.pdf:pdf—“endverbfl—“field-journaltitle”—Biometrika”fl—“field-year”—1980”fl—“wa
Vidal”,fl—family’i=B“bibinitperiod-V“bibinitperiod”,fl—given=C.”,fl—given’i=C“bibinitperiod”,fl—“”%fl—
Fern—“a”ndeZ”,fl—family’i=M“bibinitperiod-F“bibinitperiod”,fl—given=J.“bibnamedelima-A.”,fl—given’i=—
Glahn”,fl—family’i=P“bibinitperiod-G“bibinitperiod”,fl—given=V.”,fl—given’i=V“bibinitperiod”,fl—“”%fl—
based-distance-measures-infl—particular; the misstatements that logratio methods destroy distancefl—structures and are denom
8121”fl—“field-issn”—08828121”fl—“field-number”—3”fl—“field-pages”—271“bibrangedash-275”fl—“field-title”—Logratio-an
Plenum-Publishers”%fl—“fl—“strng-namehash”—AJ1”fl—“strng-fullhash”—AJ2”fl—“field-labelnamesource”—author”fl—“f
8121”fl—“field-number”—4”fl—“field-pages”—365“bibrangedash-379”fl—“field-title”—On-criteria-for-measures-of-compositio
1992--On-criteria-for-measures-offl—compositfl—“verb-ional-difference.pdf:pdf—“endverbfl—“field-journaltitle”—Mathemat
9876.00275fl—“endverbfl—“field-issn”—0035-9254”fl—“field-number”—4”fl—“field-pages”—375“bibrangedash-392”fl—“field
9876.00275fl—“endverbfl—“field-volume”—51”fl—“field-journaltitle”—Journal-of-the-Royal-Statistical-Society:Series-Cfl—(Ap
2010-11-10-r106fl—“endverbfl—“field-isbn”—1474-7596”fl—“field-issn”—1465-6906”fl—“field-number”—10”fl—“field-pages”—
2010-11-10-r106.pdf—“endverbfl—“field-volume”—11”fl—“verb-file”fl—“verb:C\$“backslash\$:/Users/dlaroche/AppData/I
2010--Differential-expressionfl—analysisfl—“verb-for-sequence-count-data.pdf:pdf—“endverbfl—“field-journaltitle”—Genom
wide-scale. However, the data produced by thefl—thousands-of-microarray-studies-published-annually-are-confounded-by—”bato
6203”fl—“field-issn”—19326203”fl—“field-number”—2”fl—“field-title”—Removing-batch-effects-in-analysis-of-expression-mic
An-Evaluation-of-Six-Batch-Adjustment-Methods.pdf:pdf—“endverbfl—“field-journaltitle”—PLoS-ONE”fl—“field-year”—2
seq-experiments. The voom method estimates the mean-variancefl—relationship-of-the-log-counts, generates a precision weight
seq-analysts-to-a-large-body-offl—methodology-developed-for-microarrays. Simulation studies show that voomfl—performs-as-we
based-RNA-seq-methods-even-when-thefl—data-are-generated-according-to-the-assumptions-of-the-earlier-methods. Twofl—case
2014-15-2-r29fl—“endverbfl—“field-issn”—1465-6914”fl—“field-number”—2”fl—“field-pages”—R29”fl—“field-title”—voom:1
seq-read-counts.””fl—“verb-url”fl—“verbfl—http://genomebiology.biomedcentral.com/articles/10.1186/gb-2014-
15-2-r2fl—“verb-9fl—“endverbfl—“field-volume”—15”fl—“verb-file”fl—“verb:C\$“backslash\$:/Users/dlaroche/AppData/Lo
2014--voom-Precision-weights-unlockfl—lineafl—“verb-r-model-analysis-tools-for-RNA-seq-read-counts.pdf:pdf—“endverbfl—
throughput-genomic-analysis. We introducedfl—surrogate-variable-analysis-(sva)-for-estimating-these-artifacts-by-(i)fl—identify
based-data-andfl—FPKM-based-data. These updates are available through the sva Bioconductorfl—package-and-I-have-made-fu
4962-(Electronic)\$“backslash\$r0305-1048-(Linking)”fl—“field-issn”—13624962”fl—“field-number”—21”fl—“field-pages”—1“bi
removing-batch-effects-and-other-unwanted-noise-fromfl—seqfl—“verb-encing-data..pdf:pdf—“endverbfl—“field-journaltitl
throughput-technologies-are-widely-used, for-example-to-assay-geneticfl—variants, gene-and-protein-expression, and-epigenetic
0064-(Electronic)\$“backslash\$r1471-0056-(Linking)”fl—“field-issn”—1471-0056”fl—“field-number”—10”fl—“field-pages”—733

throughput data." "fl" "verb-url" fl" "verb-http://dx.doi.org/10.1038/nrg2825 fl" "endverb fl" "field-volume" -11 fl" "verb-throughput data..pdf:pdf fl" "endverb fl" "field-journaltitle" -Nature reviews. Genetics fl" "field-eprinttype" -arXiv fl" "field-proportions, Challenge posed by omics data to fl" "compositional ana, Compositional data, Compositional data analysis- roots in fl" "geosciences, Current raft of nucleotide counting sequencing tec, Impact of fl" "compositional constraints in the omics- b, Impact of compositional fl" "constraints on correlation, Impact of compositional constraints on fl" "multivariate, PPM- whether molecular biosciences treat composi- tio, The Omics Imp fl" "and bioscience experiment paradigms-, Under- and over-expression used in gene fl" "expression, and biological systems with sources of variability, being scarce, in fl" "comparison v- principled fl" "approaches, interpreting measurements of mineral content of ro, molecular fl" "accountant par excellence, percentages A05-Lovell.pdf fl" "verb-f:pdf fl" "endverb fl" "field-year" -2011 fl" "endentry fl fl" "entry-Lovell2015" -article" - fl" "name-author" -Glahn", fl" "family" i=-P" "bibinitperiod-G" "bibinitperiod", fl" "given=-Vera", fl" "given" i=-V" "bibinitperiod", fl" " " "% or fl" "compositional data, differential expression needs careful interpretation, and fl" "correlation a statistical workhorse for analy- is fl" "an inappropriate measure of association. Using yeast gene expression data we fl" "show how correlation can be misleading and expression networks fl" "and clustered heatmaps. While the main aim of this study is to present fl" "proportionality as a means to 7358" fl" "field-number" -3 fl" "field-pages" -e1004075 fl" "field-title" -Proportionality: A Valid Alternative to Correlation fl" "verb-proportionality-PLOSOne2015.pdf:pdf fl" "endverb fl" "field-journaltitle" -PLoS computational biology fl" "field-Fernandez1998" -article" - fl" "name-author" -6" -" "% fl" "hash=MFJA" -" "% fl" "family=-Mart-" -"i" "n- Fern-" -"a" "nde- z", fl" "family" i=-M" "bibinitperiod-F" "bibinitperiod", fl" "given=-J" "bibnamedelima-A", fl" "given" i=-J Vidal", fl" "family" i=-B" "bibinitperiod-V" "bibinitperiod", fl" "given=-C", fl" "given" i=-C", fl" " " "% fl" "hash=PGV Glahn", fl" "family" i=-P" "bibinitperiod-G" "bibinitperiod", fl" "given=-V", fl" "given" i=-V", fl" " " "% fl" "hash=BA 900308-2-8" fl" "field-number" -1 fl" "field-pages" -526" "bibrangedash-531" fl" "field-title" -Measures of difference for comp Fern-" -"a" "nde- z et al. -1998- Measures fl" "of fl" "verb-f difference for compositional data and hierarchical clustering fl" "method On a fl" "Form of Spurious Correlation Which May Arise When Indices Are Used in the fl" "Measurement of Organs": Pearson, K. 4644" fl" "field-number" -2 fl" "field-pages" -321" "bibrangedash-332" fl" "field-title" -Small sample estimation of negative b- based transcriptome surveys suggests fl" "that RNA-seq is likely to become the platform of choice for interrogating fl" "steady-state 2010-11-3-r25 fl" "endverb fl" "field-isbn" -1465-6914 (Electronic) \$" "backslash\$r1465-6906 (Linking)" fl" "field-issn" -1465- 6906 fl" "field-number" -3 fl" "field-pages" -R25 fl" "field-title" -A scaling normalization method for differential expressi- seq data." " fl" "field-volume" -11 fl" "verb-file" fl" "verb-fl: C\$" "backslash\$: /Classes/Dissertation/Proportionality/Backgrou- III Consortium2014" -article" - fl" "name-author" -1" -" "% fl" "hash=S" -" "% fl" "family=-SEQC/MAQC- III Consortium" ", fl" "family" i=-S" "bibinitperiod", fl" " " "% fl" "list-publisher" -1" -" "% fl" "Nature Publishing Group in controls, we fl" "assess RNA sequencing (RNA-seq) performance for junction discovery and fl" "differential expression profiling- exon junctions, with fl" "textgreater" 80- " "% validated by qPCR. We find that measurements of relative fl" "expression are accu-

seq and microarrays do not provide accurate absolute measurements, and gene-specific biases are observed for all examined level profiling. The complete SEQC data sets, comprising >100 billion reads (10Tb), provide unique resources for RNA-seq analyses for clinical and regulatory settings.

doi:10.1038/nbt.2957
 Nat Biotech 34:1696–1703 (2016)
 A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium

doi:10.1038/nbt.2702
 Nat Biotech 32:1546–1554 (2014)
 A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium

doi:10.1038/nbt.3017
 Nat Biotech 33:1696–1703 (2015)
 A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium

doi:10.1038/nbt.3017
 Nat Biotech 33:1696–1703 (2015)
 A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium

RNA-Seq as a Measure of Relative Abundance: oportunities afforded by a compositional analysis framework.

Dominic LaRoche Dean Billheimer Shripad Sinari Kurt Michels
Bonnie LaFleur

September 21, 2016

1 Abstract

2 Introduction

The rapid rise in the use of RNA sequencing technology (RNA-seq) for scientific discovery has led to its consideration as a clinical diagnostic tool. However, as a new technology the analytical accuracy and reproducibility of RNA-seq must be established before it can realize its full clinical utility (**SEQC/MAQC-III Consortium 2014**, **VanKeuren-Jensen 2014**). Recent studies evaluating RNA-seq have found generally high intra-platform and inter-platform congruence across multiple laboratories (**Li 2013**; **tHoen 2013**; **SEQC/MAQC-III Consortium 2014**). Despite these promising results, there is still a need to establish reliable diagnostics, quality control metrics and improve the reproducibility of RNA-seq data. Understanding, and capatilizing on, the relative frequency nature of RNA-Seq data provides tools for identifying batch effects, creating quality control metrics, and improving reproducibility.

Relative frequency measures (hereafter referred to as *compositional data* for consistency with other disciplines) are characterzied as a vector of proportions of some whole. These proportions are necessarily positive and sum to a constant. The positivity and summation constraint complicate the analysis of compositions. For example, John Aitchison (Aitchison [1986](#)) identified the difficulty of interpreting the covariance matrix of a composition that results from the dependency in the data induced by the sum constraint. As early as 1896 Karl Pearson (**Pearson 1896**) identified the spurious correlation problem associated with compositions.

NGS-based RNA-Seq methods are inherently compositional because high-throughput RNA-Seq instruments have a maximum number of reads available per run. For example, the Roche 454 GS Junior ^(TM) claims approximately 100,000 reads per run for shotgun sequencing and 70,000 reads per run for amplicon sequencing. The Illumina Mi-Seq, with shorter read lengths, is limited to 25 million reads per sequencing run. These reads are distributed across all of the samples included in a sequencing run and, therefore, impose a total sum constraint on the data. This constraint cascades down to each probe or tag within a sample which is, in turn, constrained by the total number of reads allocated to the sample. Previous authors have identified the relative nature of RNA-Seq data (**Robinson2007**; **Anders2010**; **Robinson2010**; **Law2014**; **Lovell2015**). For example, Robinson and Smyth (2007) consider counts of RNA tags as relative abundances in their development of a model for estimating differential gene expression implemented in the Bioconductor package edgeR. Similarly, Robinson and Oshlack (2010) explicitly acknowledge the mapped read constraint when developing their widely used Trimmed-Mean of M-values (TMM) normalization method for RNA-Seq data.

Ignoring the sum constraint can lead to unexpected results and erroneous inference (**Pearson1896**; **Lovell2011**; Aitchison [1986](#)). Despite the evidence that RNA-Seq data are compositional in nature, few researchers have extended the broad set of compositional data analysis theory and tools for use in RNA-Seq analysis problems. We extend existing compositional data methodology to include statistical diagnostic tests for the identification of sample outliers and batch effects. We also show how compositional properties can be exploited to improve exploratory analyses and improve reproducibility.

3 Methods

3.1 Compositional Data

We begin with a brief introduction to compositional data, its properties, and some established analytical methods. Compositional data is defined as any data in which all elements are non-negative and sum to a fixed constant (Aitchison [1986](#)). The total sum constraint is common in biological sampling. For example, if a 1 ml sample of blood is taken this sample could be divided into several components such as plasma, red blood cells, white blood cells, and platelets. If the amount of any 1 component were to increase some other

component (or all the other components) must decrease due to the fixed volume of the sample.

For RNA-seq data, the total sum constraint is imposed by the limited number of available reads in each sequencing platform. Since this total differs between platforms we will refer to the total number of available reads as \mathbb{T} . These reads are distributed among the D samples in a sequencing run such that:

$$\sum_{i=1}^D t_i = \mathbb{T} \quad (1)$$

where t_i represents the total reads for sample i . Because of the total sum constraint, the vector \mathbf{t} is completely determined by $D-1$ elements since the D^{th} element of \mathbf{t} can be determined from the other $d = D-1$ elements and the total \mathbb{T} :

$$t_D = \mathbb{T} - \sum_{i=1}^d t_i \quad (2)$$

In 2, any of the elements can be chosen for t_D with the remaining elements labeled $1, \dots, d$ in any order (Aitchison 1986).

From equations 1 and 2 it is clear that the D samples represent a $D-1 = d$ dimensional simplex (S^d). This leads to a difficulty in interpreting the traditional $D \times D$ covariance structure. In particular, it is clear that for a D-part composition \mathbf{x} , $\text{cov}(x_1, x_1 + \dots + x_D) = 0$ since $x_1 + \dots + x_D$ is a constant. Moreover, the sum constraint induces negativity in the covariance matrix,

$$\text{cov}(x_1, x_2) + \dots + \text{cov}(x_1, x_D) = -\text{var}(x_1), \quad (3)$$

which means at least one element of each row of the covariance matrix must be negative. Aitchison refers to this as the "negative bias difficulty" (although 'bias' is not used in the traditional sense; Aitchison 1986, p. 53). The existence of these negative values creates problems for the interpretation of the covariance matrix since values are no longer free to take values between 0 and 1.

Similarly, the compositional geometry must be accounted for when measuring the distance between two compositions or finding the center of a group of compositions (Aitchison2000). Aitchison (Aitchison1992) outlined several properties for any compositional difference metric which must be met: scale invariance, permutation invariance, perturbation invariance (similar to translation invariance for Euclidean distance), and

subcompositional dominance (similar to subspace dominance of Euclidean distance). The scale invariance requirement is ignorable if the difference metric is applied to data on the same scale (which is generally not satisfied in raw RNA-seq data). The permutation invariance is generally satisfied by existing methods (**Martin-Fernandez1998**). However, the perturbation invariance and subcompositional dominance are not generally satisfied.

Because of the difficulties outlined above, standard statistical methodology is not always appropriate (Aitchison 1986) and can produce misleading results (Lovell2015). To overcome these obstacles, Aitchison (Aitchison and Shen 1980) proposed working in ratios of components. We focus on the Centered Log-Ratio (CLR) which treats the parts of the composition symmetrically and provides an informative covariance structure. The CLR transformation is defined for a D -part composition \mathbf{t} as:

$$y_i = \text{CLR}(x_i) = \log \left(\frac{x_i}{g(\mathbf{x})} \right), \quad (4)$$

where $g(\mathbf{t})$ is the geometric mean of \mathbf{t} . The $D \times D$ covariance matrix is then defined as:

$$\Gamma = [\text{cov}(y_i, y_j) : i, j = 1, \dots, D] \quad (5)$$

Although the CLR transformation gives equal treatment to every element of \mathbf{t} , the resulting covariance matrix, Γ , is singular. Therefore, care should be taken when using general multivariate methods on CLR transformed data.

Aitchison (**Aitchison1992**; Aitchison 1986) suggests using the sum of squares of all log-ratio differences. Billheimer, Guttorp, and Fagan (2001) use the geometry of compositions to define a norm which, along with the perturbation operator defined by Aitchison (Aitchison 1986), allow the interpretation of differences in compositions (**Billheimer2001**). Briefly, denote the elementwise multiplication of two positive k -vectors \mathbf{u} and \mathbf{v} by

$$\mathbf{u} \cdot \mathbf{v} \equiv (u_1 v_1, u_2 v_2, \dots, u_k v_k)'.$$

Further define the perturbation operator for composition \mathbf{x} and perturbation $\alpha \in S^d$ as

$$z = \mathbf{x} \oplus \alpha = C(\mathbf{x}\alpha)$$

$$||x|| =$$

Martin-Fernandez et al. (1998) showed that applying either Euclidean distance or Mahalanobis distance metric to CLR transformed data satisfies all the requirements of a compositional distance metric. Euclidean distance on CLR transformed compositions is referred to as Aitchison distance:

$$d_A(x_i, x_j) = \left[\sum_{k=1}^D \left(\log \left(\frac{x_{ik}}{g(x_i)} \right) - \log \left(\frac{x_{jk}}{g(x_j)} \right) \right)^2 \right]^{\frac{1}{2}}$$

or

$$d_A(x_i, x_j) = \left[\sum_{k=1}^D (clr(x_{ik}) - clr(x_{jk}))^2 \right]^{\frac{1}{2}}.$$

3.2 Outlier Detection

Problems with RNA isolation, library preparation, or sequencing may result in a low number of reads for the sample. There is currently no objective way to evaluate sample quality based on the total number of reads attributed to a sample. We develop a method grounded in the compositional nature of RNA-Seq data for objectively identifying samples with potentially poor quality.

For most experimental designs we expect the number of reads in each sample, t_i , to be equivalent notwithstanding random variation. Since these reads are part of a composition it is natural to view them as arising from a Multinomial distribution with equal probabilities. Since each cell has the same probability we test for outlying values using the Binomial distribution with probability $1/D$ and size n = total available reads.

3.3 Batch Effects and Normalization

Batch effects arising from differing laboratory conditions or operator differences have been identified as a problem in high-throughput measurement systems (leek2010; chen2011). Identifying and controlling for batch effects is a critical step in the transition of RNA-Seq from the lab to the clinic. Batch effects are typically identified with a hierarchical clustering method or principal components analysis (PCA) and re-

moved through various normalization methods (**Robinson2007**; **Anders2010**; **Robinson2010**; **Law2014**; **leek2014**).

The compositional nature of RNA-Seq data has important implications for the detection of batch effects because of the difficulty of interpreting the covariance matrix (Aitchison [1986](#)) and the incompatibility with standard measures of distance (**Martin-Fernandez1998**). The CLR transformation facilitates both batch effect detection and normalization. The CLR transformed covariance matrix is suitable for exploration through PCA (**Aitchison2002**) or hierarchical clustering using standard Euclidean distance (**Martin-Fernandez1998**).

References

blx@hook@bibinit

Aitchison, J (1986). *The statistical analysis of compositional data*. Chapman & Hall, Ltd. ISBN: 0-412-28060-4.

URL: <http://dl.acm.org/citation.cfm?id=17272> (cit. on pp. 4–7, 9).