

# clrDV: A differential variability test for RNA-Seq data based on the skew-normal distribution (#82667)

1

First submission

## Guidance from your Editor

Please submit by **6 Apr 2023** for the benefit of the authors (and your token reward) .



### Structure and Criteria

Please read the 'Structure and Criteria' page for general guidance.



### Author notes

Have you read the author notes on the [guidance page](#)?



### Raw data check

Review the raw data.



### Image check

Check that figures and images have not been inappropriately manipulated.

If this article is published your review will be made public. You can choose whether to sign your review. If uploading a PDF please remove any identifiable information (if you want to remain anonymous).

## Files

Download and review all files from the [materials page](#).

30 Figure file(s)

2 Latex file(s)

4 Table file(s)

1 Raw data file(s)

1 Other file(s)



# Structure and Criteria

## Structure your review

The review form is divided into 5 sections. Please consider these when composing your review:

1. BASIC REPORTING
2. EXPERIMENTAL DESIGN
3. VALIDITY OF THE FINDINGS
4. General comments
5. Confidential notes to the editor

You can also annotate this PDF and upload it as part of your review

When ready [submit online](#).

## Editorial Criteria

Use these criteria points to structure your review. The full detailed editorial criteria is on your [guidance page](#).

### BASIC REPORTING

- Clear, unambiguous, professional English language used throughout.
- Intro & background to show context. Literature well referenced & relevant.
- Structure conforms to [Peerj standards](#), discipline norm, or improved for clarity.
- Figures are relevant, high quality, well labelled & described.
- Raw data supplied (see [Peerj policy](#)).

### EXPERIMENTAL DESIGN

- Original primary research within [Scope of the journal](#).
- Research question well defined, relevant & meaningful. It is stated how the research fills an identified knowledge gap.
- Rigorous investigation performed to a high technical & ethical standard.
- Methods described with sufficient detail & information to replicate.

### VALIDITY OF THE FINDINGS

- Impact and novelty not assessed. *Meaningful* replication encouraged where rationale & benefit to literature is clearly stated.
- All underlying data have been provided; they are robust, statistically sound, & controlled.
- Conclusions are well stated, linked to original research question & limited to supporting results.



The best reviewers use these techniques

## Tip

## Example

**Support criticisms with evidence from the text or from other sources**

*Smith et al (J of Methodology, 2005, V3, pp 123) have shown that the analysis you use in Lines 241-250 is not the most appropriate for this situation. Please explain why you used this method.*

**Give specific suggestions on how to improve the manuscript**

*Your introduction needs more detail. I suggest that you improve the description at lines 57- 86 to provide more justification for your study (specifically, you should expand upon the knowledge gap being filled).*

**Comment on language and grammar issues**

*The English language should be improved to ensure that an international audience can clearly understand your text. Some examples where the language could be improved include lines 23, 77, 121, 128 - the current phrasing makes comprehension difficult. I suggest you have a colleague who is proficient in English and familiar with the subject matter review your manuscript, or contact a professional editing service.*

**Organize by importance of the issues, and number your points**

- 1. Your most important issue*
- 2. The next most important item*
- 3. ...*
- 4. The least important points*

**Please provide constructive criticism, and avoid personal opinions**

*I thank you for providing the raw data, however your supplemental files need more descriptive metadata identifiers to be useful to future readers. Although your results are compelling, the data analysis should be improved in the following ways: AA, BB, CC*

**Comment on strengths (as well as weaknesses) of the manuscript**

*I commend the authors for their extensive data set, compiled over many years of detailed fieldwork. In addition, the manuscript is clearly written in professional, unambiguous language. If there is a weakness, it is in the statistical analysis (as I have noted above) which should be improved upon before Acceptance.*

# clrDV: A differential variability test for RNA-Seq data based on the skew-normal distribution

Hongxiang Li<sup>1</sup>, Tsung Fei Khang<sup>Corresp. 1, 2</sup>

<sup>1</sup> Institute of Mathematical Sciences, Universiti Malaya, Kuala Lumpur, Malaysia

<sup>2</sup> Universiti Malaya Centre for Data Analytics, Universiti Malaya, Kuala Lumpur, Malaysia

Corresponding Author: Tsung Fei Khang  
Email address: tfkhang@um.edu.my

**Background.** Pathological conditions may result in certain genes having expression variance that differs markedly from control's. Finding such genes from gene expression data can provide invaluable candidates for therapeutic intervention. Under the dominant paradigm for modeling RNA-Seq gene counts using the negative binomial model, tests of differential variability are challenging to develop, owing to dependence of the variance on the mean.

**Methods.** Here, we describe clrDV, a statistical method for detecting genes that show differential variability between two populations. We present the skew-normal distribution for modeling gene-wise null distribution of centered log-ratio transformation of compositional RNA-seq data.

**Results.** Simulation results show that clrDV has false discovery rate and probability of Type II error that are on par with or superior to existing methodologies. In addition, its run time is faster than the closest competitor's, and remains relatively constant for increasing sample size per group. Analysis of a large neurodegenerative disease RNA-Seq dataset using clrDV successfully recovers multiple gene candidates that have been reported to be associated with Alzheimer's disease. Additionally, we find that most of the genes with differential variability have smaller relative gene expression variance in the Alzheimer's disease population compared to the control population.

# 1 clrDV: A Differential Variability Test for 2 RNA-Seq Data Based on the Skew-normal 3 Distribution

4 Hongxiang Li<sup>1</sup> and Tsung Fei Khang<sup>1,2</sup>

5 <sup>1</sup>Institute of Mathematical Sciences, Faculty of Science, Universiti Malaya, 50603 Kuala  
6 Lumpur, Malaysia

7 <sup>2</sup>Universiti Malaya Centre for Data Analytics, Universiti Malaya, 50603 Kuala Lumpur,  
8 Malaysia

9 Corresponding author:

10 Tsung Fei Khang<sup>2</sup>

11 Email address: tfkhang@um.edu.my

## 12 ABSTRACT

13 **Background.** Pathological conditions may result in certain genes having expression  
14 variance that differs markedly from control's. Finding such genes from gene expression  
15 data can provide invaluable candidates for therapeutic intervention. Under the dominant  
16 paradigm for modeling RNA-Seq gene counts using the negative binomial model,  
17 tests of differential variability are challenging to develop, owing to dependence of the  
18 variance on the mean.

19 **Methods.** Here, we describe clrDV, a statistical method for detecting genes that  
20 show differential variability between two populations. We present the skew-normal  
21 distribution for modeling gene-wise null distribution of centered log-ratio transformation  
22 of compositional RNA-seq data.

23 **Results.** Simulation results show that clrDV has false discovery rate and probability  
24 of Type II error that are on par with or superior to existing methodologies. In addition,  
25 its run time is faster than the closest competitor's, and remains relatively constant  
26 for increasing sample size per group. Analysis of a large neurodegenerative disease  
27 RNA-Seq dataset using clrDV successfully recovers multiple gene candidates that  
28 have been reported to be associated with Alzheimer's disease. Additionally, we find  
29 that most of the genes with differential variability have smaller relative gene expression  
30 variance in the Alzheimer's disease population compared to the control population.

## 31 1 INTRODUCTION

### 32 1.1 Background

33 Finding patterns of gene expression variation that are associated with a biological  
34 condition of interest is the first step towards elucidating the molecular basis underlying  
35 a biological process. Currently, bulk tissue mRNA collected under specific biological  
36 conditions through RNA-sequencing (RNA-Seq) technologies remains an important  
37 approach for studying gene expression patterns. Typically, genes that show statistically  
38 and biologically meaningful difference in mean expression between conditions are of

39 interest. Indeed, pathological conditions frequently manifest as gene sets with altered  
40 mean mRNA expression levels. The identification of these genes is important for  
41 understanding how the functions of normal molecular pathways are perturbed (Van den  
42 Berge et al., 2019). Hence, detecting genes that are differentially expressed is a routine  
43 and main use of RNA-Seq data (Stark et al., 2019). To analyse differential gene  
44 expression, a multitude of statistical tests have been developed throughout the years.  
45 Methods such as edgeR (Robinson et al., 2010), DESeq2 (Love et al., 2014) and voom  
46 (Law et al., 2014) have become established, go-to methods for differential expression  
47 (DE) analysis.

48 To obtain a more complete picture of patterns of gene expression variation, we need  
49 to look beyond genes with significantly different mean expression (DE genes) between  
50 conditions (Gorlov et al., 2012). Genes that show differential variability (DV genes)  
51 are likely to be important as well because many biological phenomena are explained  
52 by changes in the variance, rather than the mean, of the distribution of gene expression  
53 level (de Jong et al., 2019). For example, genes that show differential variability  
54 between undifferentiated and differentiating states have been found to be related to  
55 body axis development, neuronal movement, and transcriptional regulation during the  
56 neural differentiation process (Ando et al., 2015). In cancer biology, DV genes are  
57 useful as biomarkers for predicting tumor progression and prognosis (Dinalankara and  
58 Corrada Bravo, 2015), and patient survival (Strbenac et al., 2016). Gorlov et al. (2012)  
59 found that genes with larger expression variance in tumors compared to normal cells  
60 show stronger association with clinically important features. In network biology, genes  
61 with high variability in expression correlate with their positions within the signaling  
62 network hierarchy (Komurov and Ram, 2010). Finally, increased gene expression  
63 variability is a common outcome of aging (Bahar et al., 2006; Stegeman and Weake,  
64 2017). Standard DE analyses are likely to miss DV gene candidates, since they are not  
65 optimized for detecting differences in expression variability.

66 To date, only a few methods are available for finding DV genes using RNA-Seq data.  
67 In contrast, even in 2015, there were at least 20 methods for detecting DE genes (Khang  
68 and Lau, 2015). For testing differential variability of genes between two populations  
69 using RNA-Seq data, initial methods co-opted techniques from microarray data analysis.  
70 DiffVar (Phipson and Oshlack, 2014) is an empirical Bayes method that depends on  
71 the limma (Smyth, 2005) framework. Subsequently, negative binomial models became  
72 popular. MDSeq (Ran and Daye, 2017) uses the coefficient of dispersion ( $\sigma^2/\mu$ ) from  
73 the negative binomial 1 (NB1) generalized linear model as a measure for variability. The  
74 variance from the NB1 model is a function of the mean  $\mu$  and the dispersion  $\phi$  parameter  
75 ( $\sigma^2 = \phi\mu$ ). The parameter  $\mu$  is treated as a technical component, whereas  $\phi$  is treated  
76 as a biological component and interpreted as a parameter for gene expression variability.  
77 de Jong et al. (2019) proposed a DV test that uses the generalized additive models for  
78 location, scale and shape (GAMLSS; (Rigby and Stasinopoulos, 2005)) framework for  
79 quantifying expression variability. GAMLSS is based on the negative binomial 2 (NB2)  
80 model, whereby the mean and the variance are related quadratically as  $\sigma^2 = \mu + \phi\mu^2$ .  
81 Recently, Roberts et al. (2022) developed DiffDist, a hierarchical Bayesian model based  
82 on the NB2 model. In their work, gene expression variability is measured using the  
83 dispersion parameter  $\phi$ , which is treated as a log-normal prior. Subsequently, test of  
84 difference in dispersion between two conditions is based on the posterior distribution

85 simulated using Markov Chain Monte Carlo (MCMC).

86 In this paper, we wish to propose **clrDV** - a novel method for detecting DV genes  
87 between two conditions in RNA-Seq data that is based on a compositional data analysis  
88 framework. The method involves a log-ratio transformation of the raw gene counts,  
89 which results in a continuous variable. We show that the skew-normal distribution  
90 with centered parameters (Azzalini, 1985) is an appropriate model for the null distribu-  
91 tion. Subsequently, we construct a Wald test statistic for testing differential variability.  
92 Through simulations, we show how well clr-DV performs compared to existing methods.  
93 Finally, we demonstrate the applied value of clrDV by using it to identify biologically  
94 meaningful genes in the analysis of a large RNA-seq dataset from a neurodegenerative  
95 disease study.

## 96 **1.2 Motivation**

97 The general idea of conducting a test of differential variability for RNA-Seq data  
98 involves testing the equality of variances (equivalently, standard deviations) between two  
99 populations. The variance parameter is embedded in some probability distribution that  
100 approximates the distribution of gene (more generally, transcript) counts, assuming the  
101 null hypothesis is correct. The standard approach models RNA- Seq data as a discrete  
102 random variable.

103 Before modeling can be done, the raw count data need to be normalized to account  
104 for variation in the sequencing depth of each sample. Commonly used methods include  
105 the trimmed mean of M values (TMM) (Robinson and Oshlack, 2010), the median-of-  
106 ratios method (Anders and Huber, 2010; Love et al., 2014), upper-quartile (Bullard  
107 et al., 2010), conditional quartile normalization (Hansen et al., 2012), etc. After this,  
108 a model that accounts for overdispersion commonly seen in RNA-Seq data (e.g. the  
109 NB distribution) is used, but alternative models are possible (Esnaola et al., 2013).  
110 Statistical tests of differential variability can then be based on estimators of suitable  
111 model parameters for representing expression variability.

112 In recent years, there has been an increasing call towards adopting a **compositional**  
113 **data analysis** (CoDA) framework for improving the analysis of RNA-Seq data. Indeed,  
114 in the closely related field of microbiome data analysis, CoDA forms the main theo-  
115 retical framework of data analysis and differential abundance methods (Gloor et al.,  
116 2017). Nevertheless, the diffusion of CoDA approach into RNA-Seq data analysis is  
117 slow, possibly because established protocols for routine analyses such as differential  
118 expression analysis (e.g. DESeq2, edgeR) are all based on discrete count models such as  
119 the NB model. Quinn et al. (2018b) argued that next-generation sequencing abundance  
120 data should be viewed as inherently compositional because only a portion of genes may  
121 be sampled by sequencers, and cells are likely to be constrained in their capacity for  
122 mRNA production. Furthermore, Quinn et al. (2018a) showed the feasibility of applying  
123 ALDEx2 (Fernandes et al., 2014), a tool developed for differential abundance analysis  
124 in microbiome studies under a CoDA framework, to differential expression analysis  
125 using RNA-Seq data. Encouragingly, they reported that ALDEx2 shows superior perfor-  
126 mance with respect to precision and recall when compared against edgeR and DESeq2.  
127 By removing the need to rely on assumptions that justify normalization protocols in  
128 standard count-based approaches, log-ratio based transformations of RNA-Seq data in  
129 compositional form is potentially more attractive and effective for differential expression

130 analyses (Quinn et al., 2019). More recently, McGee et al. (2019) developed absSimSeq -  
 131 a novel simulation protocol for generating realistic RNA-Seq data using a compositional  
 132 data framework.

The key step in processing compositional data involves log-ratio transformation, for which several variants are available. The simplest is the centered log-ratio (CLR) transformation, first proposed by Aitchison (1986). After CLR- transformation, the simplex space of the compositional data is transformed into the Euclidean space. It is then convenient to view CLR-transformed values as realizations of a continuous random variable. To be concrete, let  $X_{gi}$  be the read count for gene  $g$  and sample  $i$ , where  $g = 1, 2, \dots, G$  and  $i = 1, 2, \dots, n$ . For a  $G$ -component composition  $\{x_{1i}, x_{2i}, \dots, x_{Gi}\}$ , the CLR- transformation of  $X_{gi}$  is given by

$$\text{CLR}(X_{gi}) = \log \left\{ \frac{x_{gi}}{(\prod_{g'=1}^G x_{g'i})^{1/G}} \right\} = \log(x_{gi}) - \frac{1}{G} \sum_{g'=1}^G \log(x_{g'i}),$$

133 for  $g' = 1, 2, \dots, G$ . We call  $\text{CLR}(X_{gi})$  the relative gene expression, or CLR-transformed  
 134 count, of gene  $g$  and sample  $i$ . A pseudo-value 0.5 is added if  $x_{gi} = 0$  for any  $i$ . Thus,  
 135 the main challenge for using CLR-transformed data to develop a test for differential vari-  
 136 ability is modeling them using a tractable probability distribution for which estimation  
 137 of the variance parameter is practical.

## 138 2 MATERIALS AND METHODS

### 139 2.1 The skew-normal model for CLR-transformed data

We show that the null distribution of CLR-transformed count data approximately follows the skew-normal distribution (Azzalini, 1985; Azzalini and Capitanio, 2014) (see Supplementary Material S1). Denote the relative gene expression from gene  $g$  in sample  $i$  by  $Y_{gi}$ . Thus,  $Y_{gi}$  has a skew-normal distribution with centered parameters (CP), that is,  $Y_{gi} \sim \text{SN}_C(\mu_g, \sigma_g, \gamma_g)$ , where  $\mu_g$  is the mean,  $\sigma_g$  is the standard deviation, and  $\gamma_g$  is the skewness parameter,  $g = 1, 2, \dots, G$  and  $i = 1, 2, \dots, n$ . The parameter vector  $\theta_g^{(C)} = (\mu_g, \sigma_g, \gamma_g)$  has parameter space  $\mathbb{R} \times \mathbb{R}^+ \times (-k, k)$ , where  $k = \sqrt{2}(4 - \pi)/(\pi - 2)^{3/2} \approx 0.9953$ . The special case of  $\gamma_g = 0$  results in a normal distribution with mean  $\mu_g$  and variance  $\sigma_g^2$ . The probability density function of a skew-normal distribution with direct parameters (DP) is given by

$$f(y_{gi}; \xi_g, \omega_g, \alpha_g) = \frac{2}{\omega_g} \phi\left(\frac{y_{gi} - \xi_g}{\omega_g}\right) \Phi\left(\alpha_g \frac{y_{gi} - \xi_g}{\omega_g}\right),$$

with location parameter  $\xi_g \in \mathbb{R}$ , scale parameter  $\omega_g \in \mathbb{R}^+$ , and skewness parameter  $\alpha_g \in \mathbb{R}$ ;  $\phi(\cdot)$  and  $\Phi(\cdot)$  are the probability density function and the cumulative distribution function of the standard normal distribution, respectively. The skew-normal distribution with CP is derived from the DP form via the mapping (Azzalini and Capitanio, 2014)

$$\mu_g = \xi_g + b\omega_g\delta_g, \quad \sigma_g = \omega_g\sqrt{1 - b^2\delta_g^2}, \quad \gamma_g = \frac{4 - \pi}{2} \frac{b^3\alpha_g^3}{\{1 + (1 - b^2)\alpha_g^2\}^{3/2}}; \quad (1)$$



and the inverse mapping is provided by

$$\xi_g = \mu_g - b\omega_g\delta_g, \quad \omega_g = \frac{\sigma_g}{\sqrt{1 - b^2\sigma_g^2}}, \quad \alpha_g = \frac{R}{\sqrt{b^2 - (1 - b^2)R^2}}, \quad (2)$$

140 where  $b = \sqrt{2/\pi}$ ,  $\delta_g = \alpha_g/\sqrt{1 + \alpha_g^2}$ , and  $R = \sqrt[3]{2\gamma_g/(4 - \pi)}$ .

For a single sample, the log-likelihood function for  $\boldsymbol{\theta}_g^{(D)} = (\xi_g, \omega_g, \alpha_g)^T$  is given by

$$\ell_1 = \log L(\boldsymbol{\theta}_g^{(D)}; y_{gi}) = c - \log \omega_g - \frac{(y_{gi} - \xi_g)^2}{2\omega_g^2} + \zeta_0\left(\alpha_g \frac{y_{gi} - \xi_g}{\omega_g}\right),$$

where  $c$  is a constant and  $\zeta_0(\cdot) = \log\{2\Phi(\cdot)\}$ . Taking  $z_{gi} = (y_{gi} - \xi_g)/\omega_g$ , we obtain the partial derivatives of  $\ell_1$ :

$$\frac{\partial \ell_1}{\partial \xi_g} = \frac{z_{gi}}{\omega_g} - \frac{\alpha_g}{\omega_g} \zeta_1(\alpha_g z_{gi}), \quad \frac{\partial \ell_1}{\partial \omega_g} = -\frac{1}{\omega_g} + \frac{z_{gi}^2}{\omega_g} - \frac{\alpha_g}{\omega_g} \zeta_1(\alpha_g z_{gi})z_{gi}, \quad \frac{\partial \ell_1}{\partial \alpha_g} = \zeta_1(\alpha_g z_{gi})z_{gi};$$

thus the likelihood equations for a sample of size  $n$  are given by

$$\sum_{i=1}^n z_{gi} - \alpha_g \sum_{i=1}^n \zeta_1(\alpha_g z_{gi}) = 0, \quad \sum_{i=1}^n z_{gi}^2 - \alpha_g \sum_{i=1}^n z_{gi} \zeta_1(\alpha_g z_{gi}) = n, \quad \sum_{i=1}^n z_{gi} \zeta_1(\alpha_g z_{gi}) = 0, \quad (3)$$

where  $\zeta_1(\cdot) = \phi(\cdot)/\Phi(\cdot)$ . Numerical methods are necessary to solve these equations. Azzalini and Capitanio (2014) suggested that a sample size up to about 50 may be necessary for the skew-normal distribution. To initialize the search, method of moments (MM) estimates are chosen as starting points for the CP components in Equation (1). The MM estimators for the centered parameters are given by

$$\tilde{\mu}_g = \bar{Y}_g, \quad \tilde{\sigma}_g = S_g, \quad \tilde{\gamma}_g = \frac{M_{g,3}}{S_g^3}, \quad (4)$$

141 respectively, where  $\bar{Y}_g$  is the sample mean,  $S_g$  is the sample standard deviation, and  $M_{g,3}$   
 142 is the sample third central moment. By estimating the CP components in Equation (1)  
 143 using Equation (4), and then converting them to DP components using Equation (2),  
 144 we obtain the MM estimators of the DP components:  $\tilde{\xi}_g$ ,  $\tilde{\omega}_g$  and  $\tilde{\alpha}_g$ . Subsequently, a  
 145 search of the DP space where Equation (3) holds is done. Once  $\hat{\boldsymbol{\theta}}_g^{(D)} = (\hat{\xi}_g, \hat{\omega}_g, \hat{\alpha}_g)$  is  
 146 obtained, it is mapped to Equation (1) to get  $\hat{\boldsymbol{\theta}}_g^{(C)} = (\hat{\mu}_g, \hat{\sigma}_g, \hat{\gamma}_g)$ , the maximum likelihood  
 147 estimators of the centered parameters.

Under regular maximum likelihood estimation, certain data values can produce a divergent  $\hat{\alpha}_g$ . To overcome this problem, Azzalini and Arellano-Valle (2013) proposed a maximum penalized likelihood estimation (“Qpenalty”) approach. A non-negative penalty term  $Q$  that penalizes the divergence of the skewness parameter  $\alpha_g$  is formulated as  $Q = c_1 \log(1 + c_2 \alpha_g^2)$ , where  $c_1 \approx 0.87591$  and  $c_2 \approx 0.85625$  (Azzalini and Arellano-Valle, 2013; Azzalini and Capitanio, 2014). Then, the maximum penalized likelihood for  $\boldsymbol{\theta}_g^{(D)}$  is the penalized log-likelihood

$$\ell_p(\boldsymbol{\theta}_g^{(D)}) = \ell(\boldsymbol{\theta}_g^{(D)}; \mathbf{y}_g) - Q, \quad (5)$$

where  $\mathbf{y}_g = (y_{g1}, y_{g2}, \dots, y_{gn})$ ,  $\ell(\boldsymbol{\theta}_g^{(D)}; \mathbf{y}_g)$  is the log-likelihood function with respect to the parameter vector  $\boldsymbol{\theta}_g^{(D)}$ :

$$\ell(\boldsymbol{\theta}_g^{(D)}; \mathbf{y}_g) = \text{constant} - n \log \omega_g - \sum_{i=1}^n \frac{(y_{gi} - \xi_g)^2}{2\omega_g^2} + \sum_{i=1}^n \zeta_0\left(\alpha_g \frac{y_{gi} - \xi_g}{\omega_g}\right).$$

148 The maximum penalized likelihood estimator (MPLE),  $\tilde{\boldsymbol{\theta}}_g^{(D)}$ , is a finite point that maxi-  
 149 mizes  $\ell_p(\boldsymbol{\theta}_g^{(D)})$ . The standard errors of  $\tilde{\boldsymbol{\theta}}_g^{(D)}$  can be approximated from the correspond-  
 150 ing penalized information matrix as  $\text{Var}(\tilde{\boldsymbol{\theta}}_g^{(D)}) \approx -\ell_p''(\tilde{\boldsymbol{\theta}}_g^{(D)})^{-1}$ .

The ‘‘MPpenalty’’ approach (Azzalini and Capitanio, 2014) defines the penalty function  $Q$  in Equation (5) as  $-\log \pi_m(\alpha_g)$ , where  $\pi_m$  is a prior distribution for the skewness parameter  $\alpha_g$ . The matching prior (Cabras et al., 2012) for  $\alpha_g$ , allowing for the presence of  $\boldsymbol{\psi} = (\xi_g, \omega_g)$ , is given by

$$\pi_m(\alpha_g) \propto \left( I_{\alpha_g \alpha_g}(\hat{\boldsymbol{\psi}}, \alpha_g) - I_{\alpha_g \boldsymbol{\psi}}(\hat{\boldsymbol{\psi}}, \alpha_g) I_{\boldsymbol{\psi} \boldsymbol{\psi}}(\hat{\boldsymbol{\psi}}, \alpha_g)^{-1} I_{\boldsymbol{\psi} \alpha_g}(\hat{\boldsymbol{\psi}}, \alpha_g) \right)^{1/2},$$

151 where the terms involved are specific blocks of the Fisher information matrix  $\mathbf{I}$  of  $\boldsymbol{\theta}_g^{(D)}$   
 152 (see Supplementary Material S1 for details). Since  $\pi_m(0) = 0$ , the matching prior penalty  
 153 effectively penalizes  $\alpha_g = 0$  with  $Q = \infty$ .

154 To perform parameter estimation and carry out related numerical tasks involving  
 155 the skew-normal distribution, we used the `sn` (Azzalini, 2022) R package. Regular  
 156 maximum likelihood estimation of parameters of the skew-normal model was first  
 157 done using the function `selm()`. If NA values were returned, we used the maximum  
 158 penalized likelihood estimation as implemented using the `Qpenalty` option. If NA  
 159 values persisted, the `MPpenalty` option was used.

For RNA-Seq experiments comparing two populations, testing for differential variability is equivalent to testing the equality of the standard deviation of relative gene expressions in two populations, that is,  $\sigma_{g,1} = \sigma_{g,2}$ . For this purpose, we can use the Wald statistic

$$Z_g = \frac{\hat{\sigma}_{g,2} - \hat{\sigma}_{g,1}}{\sqrt{\text{Var}(\hat{\sigma}_{g,2}) + \text{Var}(\hat{\sigma}_{g,1})}},$$

160 for  $g = 1, 2, \dots, G$ , where  $\hat{\sigma}_{g,j}$ ,  $j = 1, 2$  are the maximum likelihood estimators of the  
 161 standard deviation of the skew-normal distribution with centered parameters for popula-  
 162 tion 1 and population 2, and  $\text{Var}(\hat{\sigma}_{g,j})$ ,  $j = 1, 2$  are the corresponding diagonal elements  
 163 of the estimated Fisher information matrix of centered parameters  $\boldsymbol{\theta}_g^{(C)} = (\mu_g, \sigma_g, \gamma_g)$ .  
 164 To control the false discovery rate (FDR) as a result of conducting multiple independent  
 165 hypothesis tests across genes, we applied the Benjamini-Yekutieli procedure (Benjamini  
 166 and Yekutieli, 2001). Note that in the context of samples, FDR is estimated as the  
 167 sample proportion of false discoveries.

## 168 2.2 Data Description

169 In order to study the performance of `clrDV` and other existing methods with respect to  
 170 FDR and probability of Type II error, it is necessary to simulate the null distribution with

171 realistic parameter values. For this purpose, we used two real RNA-Seq datasets. The  
172 first dataset (GEO accession number: GSE123658) contains whole blood RNA-Seq data  
173 from from 39 Type 1 diabetes patients and 43 healthy donors (Leal Valentim et al., 2020),  
174 with 16,785 transcripts. The second dataset (GEO accession number: GSE150318)  
175 contains longitudinal gene expression data from 114 short-lived killfish *Nothobranchius*  
176 *furzeri* measured at 10 weeks and 20 weeks of age (Kelmer Sacramento et al., 2020),  
177 with 26,739 transcripts. Hereafter, we call these two datasets the “Valentim dataset” and  
178 the “Kelmer dataset”.

179 For empirical assessment, we used the Mayo RNASeq dataset (Allen et al., 2016),  
180 which consists of 278 samples and 64,253 transcripts. In this study, RNA was isolated  
181 from the temporal cortex of brains of patients with four biological conditions: control  
182 ( $n = 80$ ), Alzheimer’s disease (AD;  $n = 84$ ), progressive supranuclear palsy (PSP;  
183  $n = 84$ ) and pathologic aging ( $n = 30$ ). We chose to compare the control group against  
184 the AD and the PSP group respectively, since the sample sizes in these groups are  
185 reasonably large and balanced.

### 186 2.3 Simulation study

187 Only transcripts that satisfy two conditions in each group were used for simulation:  
188 (i) average count-per- million (CPM) above 0.5; and (ii) less than 85% of samples  
189 have zero count. Then, 2000 of the filtered genes were randomly selected. For each  
190 gene, an NB2 model was fitted. We simulated 10% of the genes to be DV genes by  
191 multiplying their size parameter ( $1/\phi$ ) with a random value  $x$ , where  $x \in (0.25, 0.5) \cup$   
192  $(2, 4)$ . Counts were then simulated based on the fitted NB2 model, for six sample sizes  
193 (50, 100, 125, 150, 200, 250) using the `polyester` (Frazee et al., 2015) R package. A  
194 total of 30 instances were thus simulated. Genes with BY-adjusted  $p$ -value  $< 0.05$  were  
195 flagged as having differential variability.

196 The performance of `clrDV` against `MDSeq`, `diffVar`, and `GAMLSS` (Benjamini-  
197 Hochberg (BH) and Benjamini-Yekutieli (BY) variants) was evaluated by considering  
198 their FDR and probability of Type II error. Additionally, we also recorded the run time  
199 of each method. `DiffDist` was excluded from the evaluation since it needs to perform  
200 MCMC simulations to generate the posterior distribution. As such, it is computationally  
201 expensive to implement and difficult to justify as a choice for routine application. Indeed,  
202 running `DiffDist` on an RNA-Seq dataset with 43 samples per group and 23,416 tran-  
203 scripts, Roberts et al. (2022) reported that `DiffDist` took about three hours to complete,  
204 compared to 12 minutes for `GAMLSS` and 4 minutes for `MDSeq`.

### 205 2.4 Empirical assessment

206 We applied `clrDV` to the Mayo RNA-Seq dataset to assess its capacity for detecting  
207 DV genes that are contextually meaningful. Analysis using `MDSeq` and `GAMLSS`  
208 (BH and and BY variants) were also done. We dropped `diffVar` because this method  
209 performed poorly during the simulation stage. Volcano plots were used to inspect the  
210 biological effect size and statistical significance of all genes tested. Venn diagrams were  
211 used to identify sets of genes that are identically recovered by all three methods, by  
212 combinations of two methods, or uniquely recovered by a single method. Violin plots of  
213 selected DV genes were made to verify computational results.

## 214 **2.5 Tools and computing environment**

215 Computational tasks were done in a computer with a 1.80 GHz i5-8265U CPU and an  
216 8GB RAM processor. R (R Core Team, 2022) (version 4.2.1) operating in Windows 10  
217 was used. The complete list of R packages used is given in the Supplementary Material  
218 S2. ENSEMBL gene ID to gene symbol conversion was done using the application  
219 programming interface of the BioTools.fr website (Saurin, 2022).

## 220 **3 RESULTS**

### 221 **3.1 Simulation study**

222 We found that the skew-normal distribution with centered parameters fit the CLR-  
223 transformed count data well. Two examples are given in Figure 1. Additional examples  
224 can be readily inspected using the R codes provided. Figure 2 shows the scatter plots of  
225 probability of Type II error against FDR for analysis of the simulated Valentim dataset,  
226 for each of the six sample size per group scenarios. diffVar is clearly uniformly inferior  
227 to all other methods (mean probability of Type II error  $> 0.05$  and FDR  $> 0.17$ , for all  
228 sample sizes).

229 For sample size of 50, all methods show relatively larger mean probability of Type  
230 II error ( $> 0.2$ ); additionally, diffVar and GAMLSS-BH show high mean FDR ( $> 0.05$ ).  
231 Against MDSeq, clrDV is uniformly superior with respect to mean FDR and mean  
232 probability of Type II error; against GAMLSS-BH, clrDV has uniformly superior mean  
233 FDR; against GAMLSS-BY, clrDV gives approximately similar mean FDR and mean  
234 probability of Type II error. When sample size is very large (250), clrDV, MDSeq and  
235 GAMLSS-BY give similar performance. With respect to computing speed, clrDV is  
236 substantially faster than GAMLSS (both BH and BY variants) as sample size increases  
237 (Supplementary Material Table S1). For the analysis of simulated data from the Kelmer  
238 dataset, we find clrDV to have comparable mean FDR and mean probability of Type  
239 II error (Figure 3) as MDSeq and GAMLSS-BY. However, clrDV computing time  
240 remains almost constant across the six sample sizes, whereas MDSeq and GAMLSS  
241 have computing times that increase with sample size (Supplementary Material Table S2).  
242 diffVar and GAMLSS-BH are inferior in controlling FDR across all six sample sizes.

### 243 **3.2 Analysis of the Mayo RNA-Seq dataset**

244 After filtering, sample sizes of the control, the AD and the PSP groups were 78, 82,  
245 and 84, respectively. For the AD and the control group comparison, a total of 18,664  
246 transcripts were left; for the PSP and control comparison, 18,636 transcripts were left.  
247 For MDSeq and GAMLSS, we normalized the raw counts using TMM normalization.

#### 248 **3.2.1 Detection of genes with differential variability**

249 Applying the procedure described in Section 2, we estimated the standard deviation  
250 of the CLR- transformed data, computed the Wald statistic and subsequently the BY-  
251 adjusted p-value for each tested gene. For the control vs. AD comparison, we detected a  
252 set of 4754 DV genes (see Supplementary Table S3 for complete list); for the control  
253 vs. PSP comparison, 4859 DV genes were detected (see Supplementary Table S4 for  
254 complete list). For the majority of DV genes, the estimated standard deviation in the  
255 control group is larger than the one in the treatment group (Figure 4). This observation

256 suggests that genes with decreased expression variability among patients with AD are  
257 far more common than those that show increased variability.

258 Figure 5 shows the number of significant DV genes identified by clrDV, MDSeq,  
259 GAMLSS-BH and GAMLSS- BY for the control vs. AD comparison (see Supplemen-  
260 tary Table S5 for complete list). GAMLSS-BH detected the most DV genes (9926),  
261 followed by MDSeq (6924), and clrDV (4754). The high confidence gene set, defined  
262 as the intersection of DV genes from each method, contains genes with estimated  
263  $\log_2(\text{SD ratio})$  that is relatively large ( $> 0.5$ ). About 99.8% (4743/4754) of DV genes  
264 detected by clrDV are also identified by MDSeq or GAMLSS-BH; 92.0% (4374/4754)  
265 are detected by both MDSeq and GAMLSS-BH; about 0.2% (11/4754) are uniquely  
266 identified by clrDV. GAMLSS-BH identified very large numbers of DV genes in this  
267 dataset, but the majority of these are probably false positives, given its relatively poorer  
268 control of FDR as shown in the results of the simulation studies. Moreover, these DV  
269 genes have estimated  $\log_2(\text{SD ratio})$  with relatively small magnitude, as indicated by  
270 the violin plots (Figure 5(c)).

271 Using GAMLSS-BY, only 6079 DV genes were detected, compared to 9926 using  
272 GAMLSS-BH. Thus, GAMLSS-BY primarily helps improve FDR by reducing the  
273 number of DV genes called. Between 61.7% (4271/6924) and 89.8% (4271/4754)  
274 of the DV genes detected by one method are detected by all three. About 97.0%  
275 (4613/4754) of DV genes detected by clrDV are identified by one of other two methods,  
276 and 3.0% (141/4754) of DV genes detected by clrDV are unique.

277 The result of the control vs. PSP comparison is similar (Figure 6; Supplementary  
278 Material Table S6). GAMLSS-BH also detected the most number of DV genes (9707),  
279 followed by MDSeq (6894), and clrDV (4859). Up to 99.4% (4831/4859) of DV  
280 genes identified by clrDV are detected by MDSeq or GAMLSS-BH; about 89.1%  
281 (4329/4859) are detected by both MDSeq and GAMLSS-BH; about 0.6% (28/4859)  
282 are uniquely identified by clrDV. Using GAMLSS-BY, only 6024 DV genes were flagged.  
283 Approximately 95.9% (4658/4859) of DV genes identified by clrDV are also identified  
284 by MDSeq or GAMLSS-BY; about 86.1% (4186/4859) are detected by both MDSeq  
285 and GAMLSS-BY; about 4.1% (201/4859) are uniquely detected by clrDV.

286 The violin plots (Figure 5 and Figure 6) suggest that the DV genes uniquely called by  
287 clrDV may be more likely to true positives, given that the magnitude of  $\log_2(\text{SD ratio})$   
288 is generally larger than 0.5. For those uniquely called by GAMLSS or MDSeq, the order  
289 of magnitude is generally below 0.5. With respect to run time, for the control vs. AD  
290 comparison, clrDV took about 7.5 minutes, compared to 6 minutes for MDSeq, and 13  
291 minutes for GAMLSS; for the control vs. PSP comparison, clrDV took about 7 minutes,  
292 while MDSeq used 6 minutes, and GAMLSS used 15 minutes.

### 293 **3.2.2 Biological significance of detected differential variability genes**

294 In the control vs. AD comparison, four of the DV genes that have the largest estimated  
295 SD ratio above 1 are LTBP2, SLPI, C2orf40, and SLC47A1 (Figure 7). All four genes  
296 have been reported to be associated with Alzheimer's disease in the literature. The  
297 latent transforming growth factor (TGF)- beta binding proteins (LTBP) are important  
298 components of the extracellular matrix (Robertson et al., 2015). They interact with  
299 fibrillin microfibrils, and are known to be mediators of TGF- $\beta$  functions (Rifkin et al.,  
300 2018), dysfunctions of which have been implicated in Alzheimer's disease (Das and

301 Golde, 2006). Then, the secretory leukocyte protease inhibitor protein (SLPI) is known  
302 to regulate the penetrance of frontotemporal lobar degeneration (FTLD) in patients who  
303 have mutations in the progranulin gene (Ghidoni et al., 2014). Loss of progranulin  
304 function has been found to enhance microglial neuroinflammation, which is implicated  
305 in Alzheimer's disease (Mendsaikhon et al., 2019). Podvin et al. (2016) found that  
306 C2orf40 is a neuroimmune factor in Alzheimer's disease. The SLC47A1 (solute carrier  
307 family 47 member 1) protein is expressed in both the kidney and the brain, and recent  
308 research has suggested a linkage between kidney diseases and Alzheimer's disease (Shi  
309 et al., 2018; Kelly and Rothwell, 2022).

310 We detected 74 genes from the SLC family in the high confidence DV gene set,  
311 including four members of the SLC39 family. Lang et al. (2012) demonstrated the  
312 modulating effect of dZip1, the ortholog of human SLC39 family transporter, on zinc  
313 ion uptake using a *Drosophila* model. Zinc is known to induce amyloid beta formation  
314 (Bush et al., 1994). Inhibition of dZip1 produces substantial reduction of amyloid beta  
315 peptide 42 ( $A\beta_{42}$ ) fibril deposits and less neurodegeneration in  $A\beta_{42}$ -transgenic flies.

316 Two of the DV genes with estimated SD ratio substantially smaller than 1 are PELP1  
317 and GP1BB (Figure 7). PELP1 mediates E2 inhibition of GSK3 $\beta$ , a neurodegenerative  
318 kinase signaling pathway in the brain (Thakkar et al., 2018). GSK3 $\beta$  is implicated  
319 in Alzheimer's disease as a key mediator of cell death (Llorens-Martin et al., 2014).  
320 The GP1BB gene produces glycoprotein 1b-beta (GPIb $\beta$ ), a subunit of the GPIb-IX-V  
321 protein complex on the surface of platelet cells. Amyloid beta peptides are known to  
322 be actively released by platelets (Bush et al., 1990; Casoli et al., 2007). Visconte et al.  
323 (2020) recently reported that recruitment of GPIb-IX-V is required for fibrillar amyloid  
324  $A\beta_{40}$  and  $A\beta_{42}$  to induce platelet aggregation. The study of the role of platelets and  
325 the pathogenesis of Alzheimer's disease is an active topic (Caticala et al., 2012).

326 We note that approximately half of genes in the high confidence gene set from the  
327 control vs. AD comparison (4271 genes) are also found in the high confidence gene  
328 sets from the control vs. PSP comparison (4186 genes). Altogether, 2149 DV genes are  
329 common to both comparisons. This observation is consistent with recent findings that  
330 transcriptomic changes are in AD and PSP relative to control are strongly correlated  
331 (Wang et al., 2022).

## 332 4 DISCUSSION

333 Our present work demonstrates that when analyzing gene expression data using the  
334 CoDA framework, the skew-normal distribution provides a natural way to model CLR-  
335 transformed data. The skew-normal distribution is a tractable model with mature  
336 computational support through the `sn` R package. A test of differential variability can  
337 therefore be based directly on the standard deviation parameter of the skew-normal  
338 distribution. Moreover, a test of differential expression that is based on the mean  
339 parameter can be derived as well. With these tests, it becomes possible to develop  
340 methods for detecting three classes of genes in two-population comparisons: (i) equal  
341 variance, different mean; (ii) equal mean, different variance; (iii) different mean, different  
342 variance. Although `clrDV` cannot differentiate genes of the second and the third type,  
343 inspection of violin plots should be useful for ascertaining whether the DV genes also  
344 appear to differ significantly in the mean of their relative expression level.

345 We observed that in the comparisons between control vs. AD and control vs. PSP, a  
346 majority of the DV genes identified by clrDV (between 86.1% and 92.0%) were already  
347 included in the high-confidence gene set, where the estimated  $\log_2(\text{SD ratio})$  has a  
348 **relatively large magnitude**. Thus, it seems that clrDV alone should be able to recover  
349 most of the DV genes of interest.

350 The relative poorer performance of MDSeq and GAMLSS could be caused by the  
351 choice of normalization. It is known that incorrect normalization leads to inflated  
352 FDR in differential expression analyses (Evans et al., 2018), yet the assumptions that  
353 justify a normalization method are usually not testable. Since existing normalization  
354 methods have been developed for the purpose of finding differentially expressed genes,  
355 the assumptions that justify their use are probably suboptimal for differential variability  
356 tests. Consequently, the performance of existing count-based approaches for DV test is  
357 likely sensitive to the choice of normalization method. However, it is beyond the scope  
358 of the present work to optimize the choice of normalization step for these count-based  
359 methods.

360 On the aspect of practical application, we note that the R codes provided by de Jong  
361 et al. (2019) for GAMLSS are not sufficiently generic and require further user modifica-  
362 tions to be suitable for routine use as a DV test. In addition, GAMLSS uses BH rather  
363 than BY as the default setting for multiple comparisons adjustment. For MDSeq, we  
364 found that it may occasionally encounter difficulties in estimating model parameters.  
365 In our analysis of the Mayo RNA-Seq dataset, we observed that 45 genes returned NA  
366 parameter estimates in the control vs. AD and the control vs. PSP comparisons. **Given**  
367 **these findings, we believe clrDV is currently the most practical and effective method for**  
368 **researchers who wish to conduct differential variability test using RNA- Seq data.**

## 369 5 CONCLUSIONS

370 Variability of gene expression at aberrant levels is one of the hallmarks of disrupted or  
371 dysregulated biological processes. Hence, detection of genes with differential variability  
372 should accompany routine differential expression analysis to expand the pool of potential  
373 therapeutic intervention targets. clrDV offers a novel approach for identifying DV genes  
374 in RNA-seq data. By modeling the null distribution of centered log- ratio transformed  
375 RNA-Seq data using a skew-normal distribution, clrDV can detect genes with expression  
376 variance that differs significantly between two populations. Simulation results demon-  
377 strate that clrDV has a comparable or superior false discovery rate and probability of  
378 Type II error compared to existing methods, while also having a faster run time for larger  
379 sample sizes per group. Applying clrDV to the Mayo RNA-seq dataset, we identified  
380 several genes associated with Alzheimer's disease, many of which had smaller relative  
381 gene expression variance in the Alzheimer's disease population compared to the control  
382 population. Crucially, the compositional data analysis framework used in this work can  
383 be extended to create statistical tests for differential expression and differential skewness  
384 using RNA-seq data. Results from this extension will be reported elsewhere.

## 385 CODE AVAILABILITY

386 We have created an R package called `clrDV` to perform the differential variability test  
387 described here. The R package and codes for reproducing the analyses in this study are

388 available at <https://github.com/Divo-Lee/clrDV>.

## 389 DATA AVAILABILITY

390 The present work did not generate no new datasets. We used datasets published by  
391 other researchers described in Section 3.2. The RNA-Seq datasets with GEO accession  
392 numbers GSE123658 and GSE150318 are freely available in the NCBI Gene Omnibus  
393 Expression database. We obtained permission from AD Knowledge Portal (accessible at  
394 <https://adknowledgeportal.org>) to access and use the Mayo RNASeq dataset for research  
395 purpose.

## 396 ACKNOWLEDGMENTS

397 We wish to thank Dr. C.Y. Ung for helpful discussions. We are grateful to Professor  
398 A. Azzalini for developing the theoretical foundations as well as computational tools  
399 for the skew normal distribution over many years, without which the present work  
400 would not be possible. The results published here are in whole or in part based on data  
401 obtained from the AD Knowledge Portal. The Mayo RNAseq study data was led by  
402 Dr. Nilüfer Ertekin-Taner, Mayo Clinic, Jacksonville, FL as part of the multi-PI U01  
403 AG046139 (MPIs Golde, Ertekin-Taner, Younkin, Price). Samples were provided from  
404 the following sources: The Mayo Clinic Brain Bank. Data collection was supported  
405 through funding by NIA grants P50 AG016574, R01 AG032990, U01 AG046139, R01  
406 AG018023, U01 AG006576, U01 AG006786, R01 AG025711, R01 AG017216, R01  
407 AG003949, NINDS grant R01 NS080820, CurePSP Foundation, and support from Mayo  
408 Foundation. Study data includes samples collected through the Sun Health Research  
409 Institute Brain and Body Donation Program of Sun City, Arizona. The Brain and Body  
410 Donation Program is supported by the National Institute of Neurological Disorders and  
411 Stroke (U24 NS072026 National Brain and Tissue Resource for Parkinsons Disease and  
412 Related Disorders), the National Institute on Aging (P30 AG19610 Arizona Alzheimers  
413 Disease Core Center), the Arizona Department of Health Services (contract 211002,  
414 Arizona Alzheimers Research Center), the Arizona Biomedical Research Commission  
415 (contracts 4001, 0011, 05-901 and 1001 to the Arizona Parkinson's Disease Consortium)  
416 and the Michael J. Fox Foundation for Parkinsons Research.

## 417 Author contributions

418 Conceptualization: Tsung Fei Khang; Methodology: all authors; Formal analysis and  
419 investigation: all authors; Writing - original draft preparation: all authors; Writing -  
420 review and editing: all authors; Supervision: Tsung Fei Khang.

## 421 Financial disclosure

422 None reported.

## 423 Conflict of interest

424 The authors declare no potential conflict of interests.



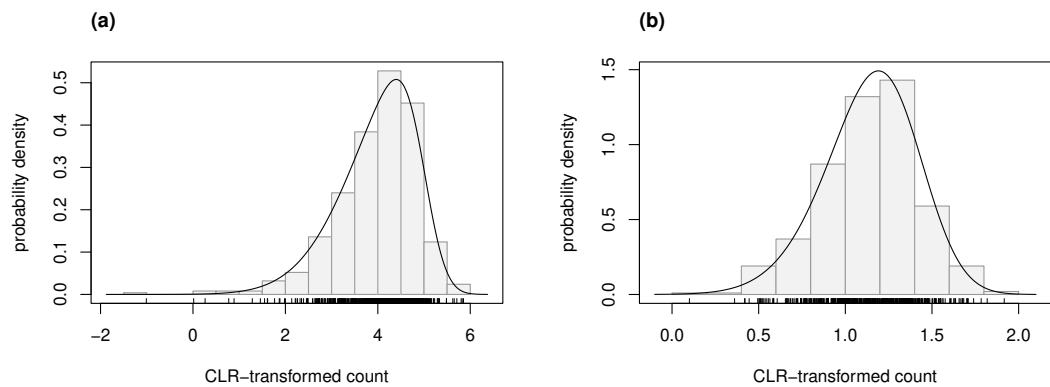
425 **REFERENCES**

- 426 Aitchison, J. (1986). *The Statistical Analysis of Compositional Data*. Chapman & Hall,  
427 London.
- 428 Allen, M., Carrasquillo, M. M., Funk, C., Heavner, B. D., Zou, F., Younkin, C. S.,  
429 Burgess, J. D., Chai, H. S., Crook, J., Eddy, J. A., et al. (2016). Human whole genome  
430 genotype and transcriptome data for Alzheimer's and other neurodegenerative diseases.  
431 *Scientific Data*, 3(89).
- 432 Anders, S. and Huber, W. (2010). Differential expression analysis for sequence count  
433 data. *Genome Biology*, 11:R106.
- 434 Ando, T., Kato, R., and Honda, H. (2015). Differential variability and correlation of gene  
435 expression identifies key genes involved in neuronal differentiation. *BMC Systems  
436 Biology*, 9:Article no. 82.
- 437 Azzalini, A. (1985). A class of distributions which includes the normal ones. *Scandina-  
438 vian Journal of Statistics*, 12(2):171–178.
- 439 Azzalini, A. (2022). *The R package sn: The skew-normal and related distributions such  
440 as the skew-t and the SUN (version 2.1.0)*. Università degli Studi di Padova, Italia.  
441 Available at: <http://azzalini.stat.unipd.it/SN/>.
- 442 Azzalini, A. and Arellano-Valle, R. B. (2013). Maximum penalized likelihood estima-  
443 tion for skew-normal and skew-t distributions. *Journal of Statistical Planning and  
444 Inference*, 143(2):419–433.
- 445 Azzalini, A. and Capitanio, A. (2014). *The Skew-Normal and Related Families*. Cam-  
446 bridge University Press, Cambridge.
- 447 Bahar, R., Hartmann, C. H., Rodriguez, K. A., Denny, A. D., Busuttill, R. A., Dollé,  
448 M. E. T., Calder, R. B., Chisholm, G. B., Pollock, B. H., and Klein, C. A. (2006).  
449 Increased cell-to-cell variation in gene expression in ageing mouse heart. *Nature*,  
450 441:1011–1014.
- 451 Benjamini, Y. and Yekutieli, D. (2001). The control of the false discovery rate in multiple  
452 testing under dependency. *Annals of Statistics*, 29(4):1165–1188.
- 453 Bullard, J. H., Purdom, E., Hansen, K. D., and Dudoit, S. (2010). Evaluation of statistical  
454 methods for normalization and differential expression in mRNA-Seq experiments.  
455 *BMC Bioinformatics*, 11(96).
- 456 Bush, A. I., Martins, R. N., Rumble, B., Moir, R., Fuller, S., Milward, E., Currie, J.,  
457 Ames, D., Weidemann, A., Fischer, P., Multhaup, G., Beyreuther, K., and Masters,  
458 C. L. (1990). The amyloid precursor protein of Alzheimer's disease is released by  
459 human platelets. *Journal of Biological Chemistry*, 265(26):15977–15983.
- 460 Bush, A. I., Pettingell, W. H., Multhaup, G., d'Paradis, M., Vonsattel, J. P., Gusella, J. F.,  
461 Beyreuther, K., Masters, C. L., and Tanzi, R. E. (1994). Rapid induction of Alzheimer  
462 A beta amyloid formation by zinc. *Science*, 265(5177):1464–1467.
- 463 Cabras, S., Racugno, W., Castellanos, M. E., and Ventura, L. (2012). A matching prior  
464 for the shape parameter of the skew-normal distribution. *Scandinavian Journal of  
465 Statistics*, 39(2):236–247.
- 466 Casoli, T., Di Stefano, G., Giorgetti, B., Grossi, Y., Balietti, M., Fattoretti, P., and  
467 Bertoni-Freddari, C. (2007). Release of beta-amyloid from high-density platelets:  
468 implications for Alzheimer's disease pathology. *Annals of the New York Academy of  
469 Sciences*, 1096:170–178.

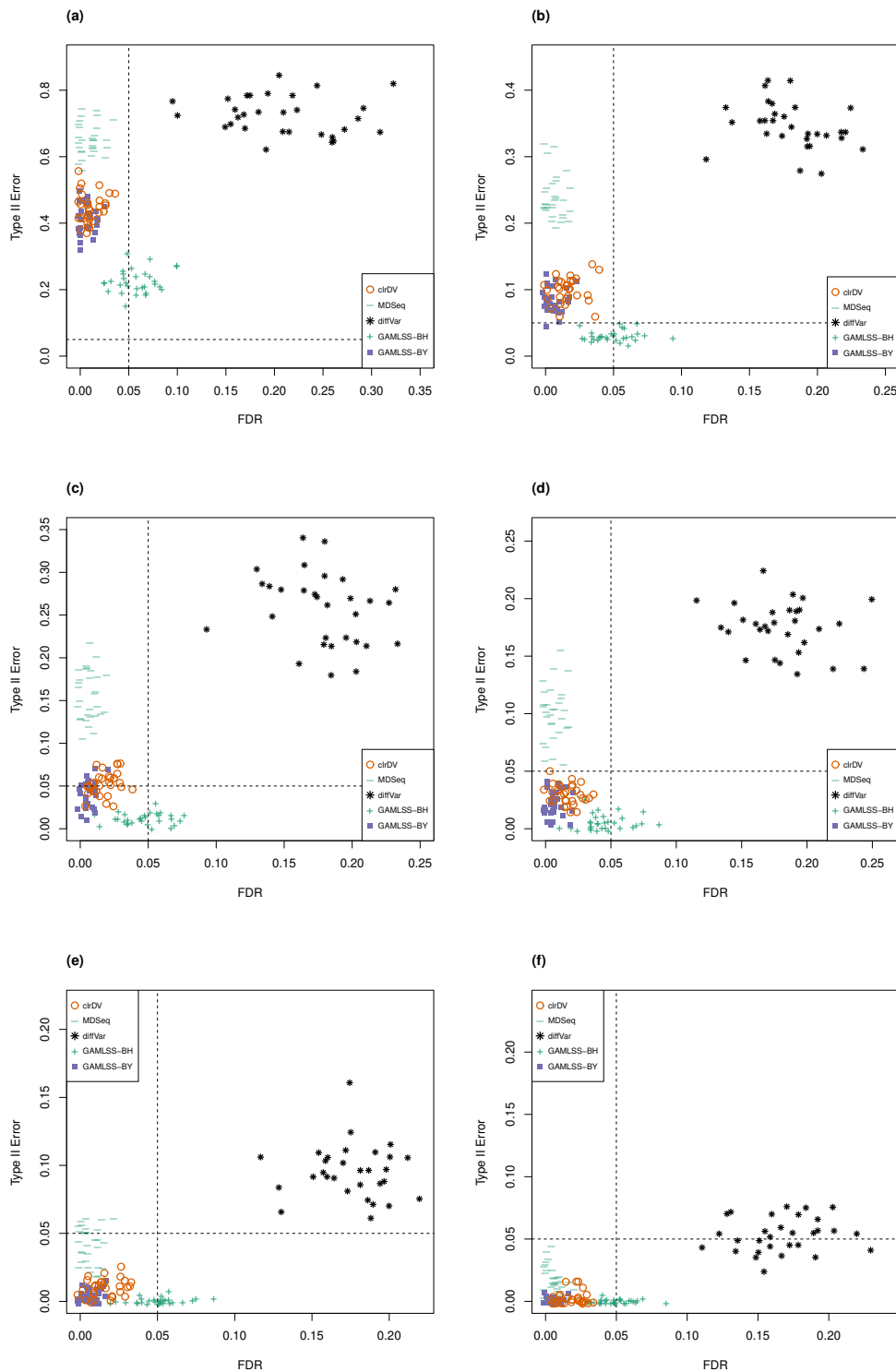
- 470 Catricala, S., Torti, M., and Ricevuti, G. (2012). Alzheimer disease and platelets: how's  
471 that relevant. *Immunity and Ageing*, 9(1):20.
- 472 Das, P. and Golde, T. (2006). Dysfunction of TGF- $\beta$  signaling in Alzheimer's disease.  
473 *Journal of Clinical Investigations*, 116(11):2855–2857.
- 474 de Jong, T. V., Moshkin, Y. M., and Guryev, V. (2019). Gene expression variability: the  
475 other dimension in transcriptome analysis. *Physiological Genomics*, 51(5):145–158.
- 476 Dinalankara, W. and Corrada Bravo, H. (2015). Gene expression signatures based on  
477 variability can robustly predict tumor progression and prognosis. *Cancer Informatics*,  
478 2015:71–81.
- 479 Esnaola, M., Puig, P., Gonzalez, D., Castelo, R., and Gonzalez, J. R. (2013). A  
480 flexible count data model to fit the wide diversity of expression profiles arising from  
481 extensively replicated RNA-seq experiments. *BMC Bioinformatics*, 14(1):1–22.
- 482 Evans, C., Hardin, J., and Stoebel, D. M. (2018). Selecting between-sample RNA-  
483 Seq normalization methods from the perspective of their assumptions. *Briefings in*  
484 *Bioinformatics*, 19(5):776–792.
- 485 Fernandes, A. D., Reid, J. N. S., Macklaim, J. M., McMurrough, T. A., Edgell, D. R.,  
486 and Gloor, G. B. (2014). Unifying the analysis of high-throughput sequencing  
487 datasets: characterizing RNA-seq, 16S rRNA gene sequencing and selective growth  
488 experiments by compositional data analysis. *Microbiome*, 2(1):1–13.
- 489 Frazee, A. C., Jaffe, A. E., Langmead, B., and Leek, J. T. (2015). Polyester: simulating  
490 RNA-seq datasets with differential transcript expression. *Bioinformatics*, 31(17):2778–  
491 2784.
- 492 Ghidoni, R., Flocco, R., Paterlini, A., Glionna, M., Caruana, L., Tonoli, E., Binetti,  
493 G., and Benussi, L. (2014). Secretory leukocyte protease inhibitor protein regulates  
494 the penetrance of frontotemporal lobar degeneration in progranulin mutation carriers.  
495 *Journal of Alzheimer's Disease*, 38:533–539.
- 496 Gloor, G., Macklaim, J., Pawlowsky-Glahn, V., and Egozcue, J. J. (2017). Microbiome  
497 datasets are compositional: and this is not optional. *Frontiers in Microbiology*, 8:2224.
- 498 Gorlov, I. P., Byun, J., Zhao, H., Logothetis, C. J., and Gorlova, O. Y. (2012). Beyond  
499 comparing means: the usefulness of analyzing interindividual variation in gene  
500 expression for identifying genes associated with cancer development. *Journal of*  
501 *Bioinformatics and Computational Biology*, 10:1241013.
- 502 Hansen, K. D., Irizarry, R. A., and Wu, Z. (2012). Removing technical variability in  
503 RNA-seq data using conditional quantile normalization. *Biostatistics*, 13(2):204–216.
- 504 Kelly, D. M. and Rothwell, P. M. (2022). Disentangling the relationship between chronic  
505 kidney disease and cognitive disorders. *Frontiers in Neurology*, 13:Article 830064.
- 506 Kelmer Sacramento, E., Kirkpatrick, J., Mazzetto, M., Baumgart, M., Bartolome, A.,  
507 Di Sanzo, S., Caterino, C., Sanguanini, M., Papaevgeniou, N., Lefaki, M., et al. (2020).  
508 Reduced proteasome activity in the aging brain results in ribosome stoichiometry loss  
509 and aggregation. *Molecular Systems Biology*, 16(6):e9596.
- 510 Khang, T. F. and Lau, C. Y. (2015). Getting the most out of RNA-seq data analysis.  
511 *PeerJ*, 3:e1360.
- 512 Komurov, K. and Ram, P. T. (2010). Patterns of human gene expression variance show  
513 strong associations with signaling network hierarchy. *BMC Systems Biology*, 4:154.
- 514 Lang, M., Wang, L., Fan, Q., Xiao, G., Wang, X., Zhong, Y., and Zhou, B. (2012). Ge-  
515 netic inhibition of solute-linked carrier 39 family transporter 1 ameliorates A $\beta$  pathol-

- ogy in a *Drosophila* model of Alzheimer's disease. *PLoS Genetics*, 8(4):e1002683.
- Law, C. W., Chen, Y., Shi, W., and Smyth, G. K. (2014). voom: Precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biology*, 15(2):R29.
- Leal Valentim, F., Mariotti-Ferrandiz, E., Klatzmann, D., Six, A., and Konza, O. (2020). Transimmunom whole blood RNA-seq data from type 1 diabetic patients and healthy volunteers. Unpublished GEO dataset. GEO accession number: GSE123658.
- Llorens-Martin, M., Jurado, J., Hernandez, F., and Avila, J. (2014). GSK-3 $\beta$ , a pivotal kinase in Alzheimer disease. *Frontiers in Molecular Neuroscience*, 7:46.
- Love, M. I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology*, 15(12):550.
- McGee, W. A., Pimentel, H., Pachter, L., and Wu, J. Y. (2019). Compositional data analysis is necessary for simulating and analyzing RNA-Seq data. *bioRxiv*, 564955:doi: <https://doi.org/10.1101/564955>.
- Mendsaikhan, A., Tooyama, I., and Walker, D. G. (2019). Microglial progranulin: Involvement in Alzheimer's disease and neurodegenerative diseases. *Cells*, 8(3):230.
- Phipson, B. and Oshlack, A. (2014). DiffVar: a new method for detecting differential variability with application to methylation in cancer and aging. *Genome Biology*, 15(9):465.
- Podvin, S., Miller, M. C., Rossi, R., Chukwueke, J., Donahue, J. E., Johanson, C. E., Baird, A., and Stopa, E. G. (2016). The orphan C2orf40 gene is a neuroimmune factor in Alzheimer's disease. *JSM Alzheimer's Disease and Related Dementia*, 3(1):1020.
- Quinn, T. P., Crowley, T. M., and Richardson, M. F. (2018a). Benchmarking differential expression analysis tools for RNA-Seq: normalization-based vs. log-ratio transformation-based methods. *BMC Bioinformatics*, 19(274).
- Quinn, T. P., Erb, I., Gloor, G., Notredame, C., Richardson, M. F., and Crowley, T. M. (2019). A field guide for the compositional analysis of any-omics data. *GigaScience*, 8(9):giz107.
- Quinn, T. P., Erb, I., Richardson, M. F., and Crowley, T. M. (2018b). Understanding sequencing data as compositions: an outlook and review. *Bioinformatics*, 34(16):2870–2878.
- R Core Team (2022). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Ran, D. and Daye, Z. J. (2017). Gene expression variability and the analysis of large-scale RNA-seq studies with the MDSeq. *Nucleic Acids Research*, 45(13):e127–e127.
- Rifkin, D. B., Rifkin, W. J., and Zilberberg, L. (2018). LTBP in biology and medicine; LTBP diseases. *Matrix Biology*, 71-72:90–99.
- Rigby, R. A. and Stasinopoulos, D. M. (2005). Generalized additive models for location, scale and shape. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 54(3):507–554.
- Roberts, A. G. K., Catchpoole, D. R., and Kennedy, P. J. (2022). Identification of differentially distributed gene expression and distinct sets of cancer-related genes identified by changes in mean and variability. *NAR Genomics and Bioinformatics*, 4(1):lqab124.
- Robertson, I. B., Horiguchi, M., Zilberberg, L., Dabovic, B., Hadjiolova, K., and Rifkin, D. B. (2015). Latent TGF- $\beta$ -binding proteins. *Matrix Biology*, 47:44–53.
- Robinson, M. D., McCarthy, D. J., and Smyth, G. K. (2010). edgeR: a Bioconductor

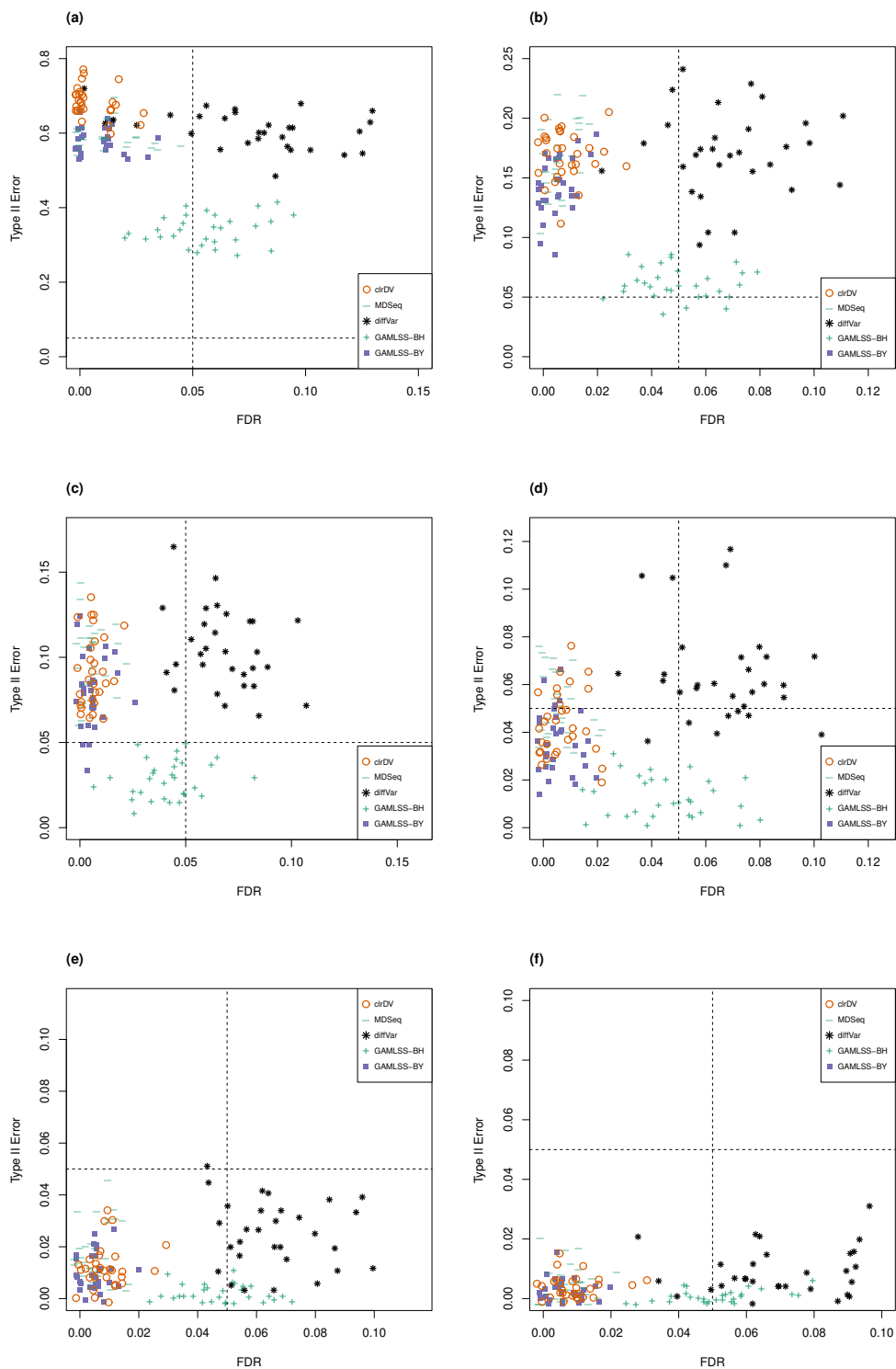
- 562 package for differential expression analysis of digital gene expression data. *Bioinform-*  
563 *atics*, 26(1):139–140.
- 564 Robinson, M. D. and Oshlack, A. (2010). A scaling normalization method for differential  
565 expression analysis of RNA-seq data. *Genome Biology*, 11(3):R25.
- 566 Saurin, A. (2022). Bioinformatics tools for genomics and transcriptomics  
567 analyses: ENSEMBL ID to Gene Symbol Converter. Available at:  
568 [https://www.biotoools.fr/human/ensembl\\_symbol\\_converter](https://www.biotoools.fr/human/ensembl_symbol_converter). Accessed: 31 August  
569 2022.
- 570 Shi, Y., Liu, Z., Shen, Y., and Zhu, H. (2018). A novel perspective linkage between kid-  
571 ney function and Alzheimer’s disease. *Frontiers in Cellular Neuroscience*, 12:Article  
572 384.
- 573 Smyth, G. K. (2005). limma: Linear models for microarray data. In Gentleman,  
574 R., Carey, V. J., Huber, W., Irizarry, R. A., and Dudoit, S., editors, *Bioinformatics  
575 and Computational Biology Solutions Using R and Bioconductor*, pages 397–420.  
576 Springer, NY.
- 577 Stark, R., Grzelak, M., and Hadfield, J. (2019). RNA sequencing: the teenage years.  
578 *Nature Reviews Genetics*, 20:631–656.
- 579 Stegeman, R. and Weake, V. M. (2017). Transcriptional signatures of aging. *Journal of  
580 Molecular Biology*, 429(16):2427–2437.
- 581 Strbenac, D., Mann, G. J., Yang, J. Y. H., and Ormerod, J. T. (2016). Differential distribu-  
582 tion improves gene selection stability and has competitive classification performance  
583 for patient survival. *Nucleic Acids Research*, 44(13):e119.
- 584 Thakkar, R., Sareddy, G. R., Zhang, Q., Wang, R., Vadlamudi, R. K., and Brann, D.  
585 (2018). PELP1: a key mediator of oestrogen signalling and actions in the brain.  
586 *Journal of Neuroendocrinology*, 30(2):e12484.
- 587 Van den Berge, K., Hembach, K. M., Sonesson, C., Tiberi, S., Clement, L., Love, M. I.,  
588 Patro, R., and Robinson, M. D. (2019). RNA sequencing data: Hitchhiker’s guide to  
589 expression analysis. *Annual Review of Biomedical Data Science*, 2:139–173.
- 590 Visconte, C., Canino, J., Vismara, M., Guidetti, G. F., Raimondi, S., Pula, G., Torti,  
591 M., and Canobbio, I. (2020). Fibrillar amyloid peptides promote platelet aggregation  
592 through the coordinated action of ITAM- and ROS-dependent pathways. *Journal of  
593 Thrombosis and Haemostasis*, 18(11):3029–3042.
- 594 Wang, X., Allen, M., İş, O., Reddy, J. S., Tutor-New, F. Q., Casey, M. C., Carrasquillo,  
595 M. M., Oatman, S. R., Min, Y., Asmann, Y. W., Funk, C., Nguyen, T., Ho, C. C. G.,  
596 Malphrus, K. M., Seyfried, N. T., Levey, A. I., Younkin, S. G., Murray, M. E., Dickson,  
597 D. W., Price, N. D., Golde, T. E., and Ertekin-Taner, N. (2022). Alzheimer’s disease  
598 and progressive supranuclear palsy share similar transcriptomic changes in distinct  
599 brain regions. *Journal of Clinical Investigation*, 132(2):e149904.



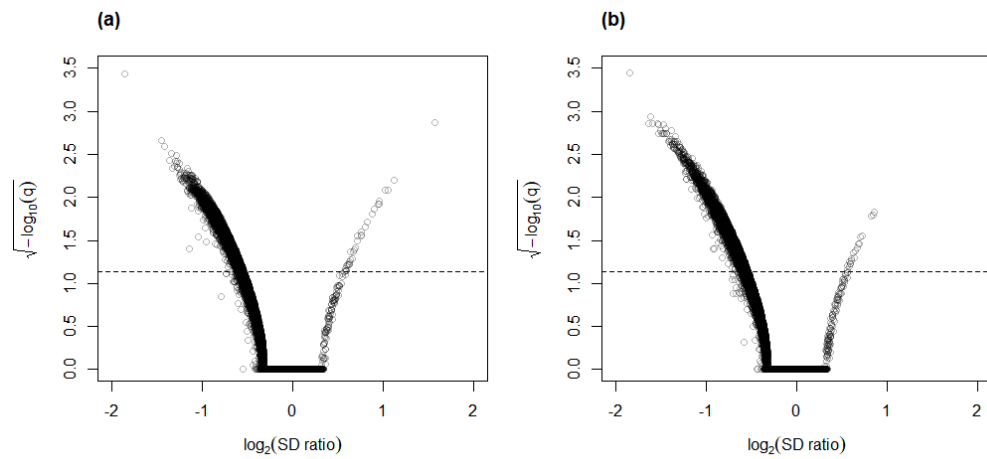
**Figure 1.** Histograms of CLR-transformed counts for two genes with fitted skew-normal curve for (a) the Valentim dataset ( $\hat{\mu} = 3.968$  (s.e. = 0.038),  $\hat{\sigma} = 0.858$ , (s.e. = 0.030) and  $\hat{\gamma} = -0.732$  (s.e. = 0.055)); (b) the Kelmer dataset ( $\hat{\mu} = 1.140$  (s.e. = 0.012),  $\hat{\sigma} = 0.275$  (s.e. = 0.009) and  $\hat{\gamma} = -0.336$  (s.e. = 0.107)).



**Figure 2.** Scatter plots of probability of Type II error vs. FDR for simulation study of the Valentim dataset (30 instances) for samples size per group of (a) 50, (b) 100, (c) 125, (d) 150, (e) 200, and (f) 250. Dashed lines represent probability of Type II error and FDR of 0.05.

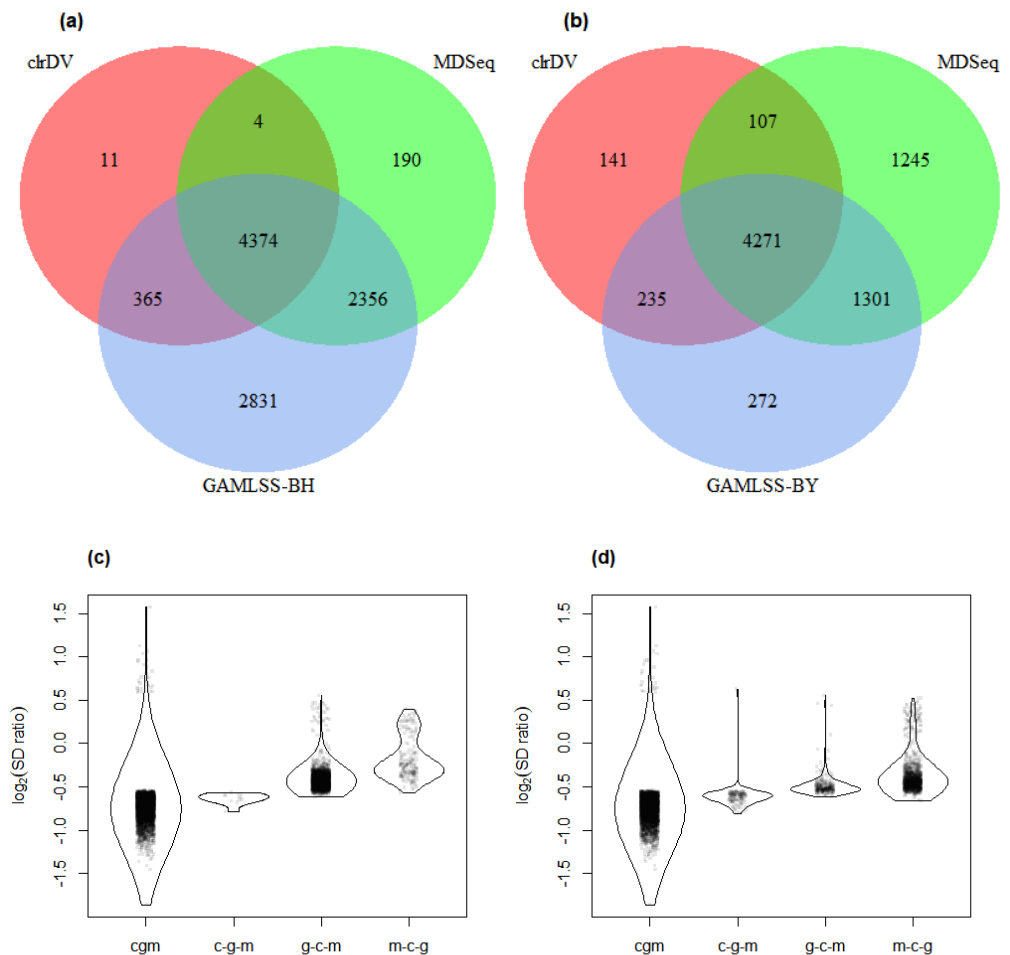


**Figure 3.** Scatter plots of probability of Type II error vs. FDR for simulation study of the Kelmer dataset (30 instances) for samples size per group of (a) 50, (b) 100, (c) 125, (d) 150, (e) 200, and (f) 250. Dashed lines represent probability of Type II error and FDR of 0.05.

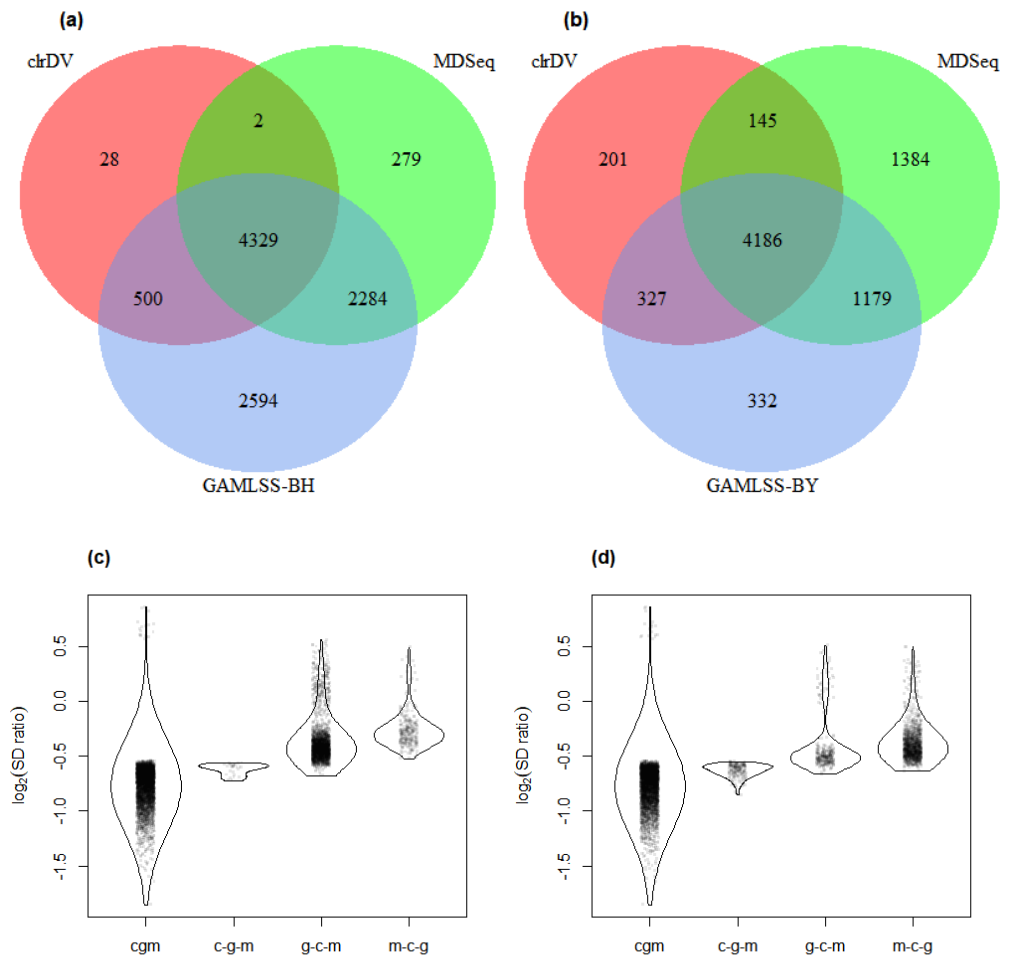


**Figure 4.** Volcano plots for (a) control vs. AD and (b) control vs. PSP comparisons for the Mayo RNA-Seq dataset. Dashed line represents the threshold of BY- adjusted p-value ( $q$ ) at 0.05 for flagging DV genes. The number of DV genes with  $\log_2(\text{SD ratio}) > 0$  and  $\log_2(\text{SD ratio}) < 0$  respectively: (a) 32 and 4722; (b) 19 and 4840.

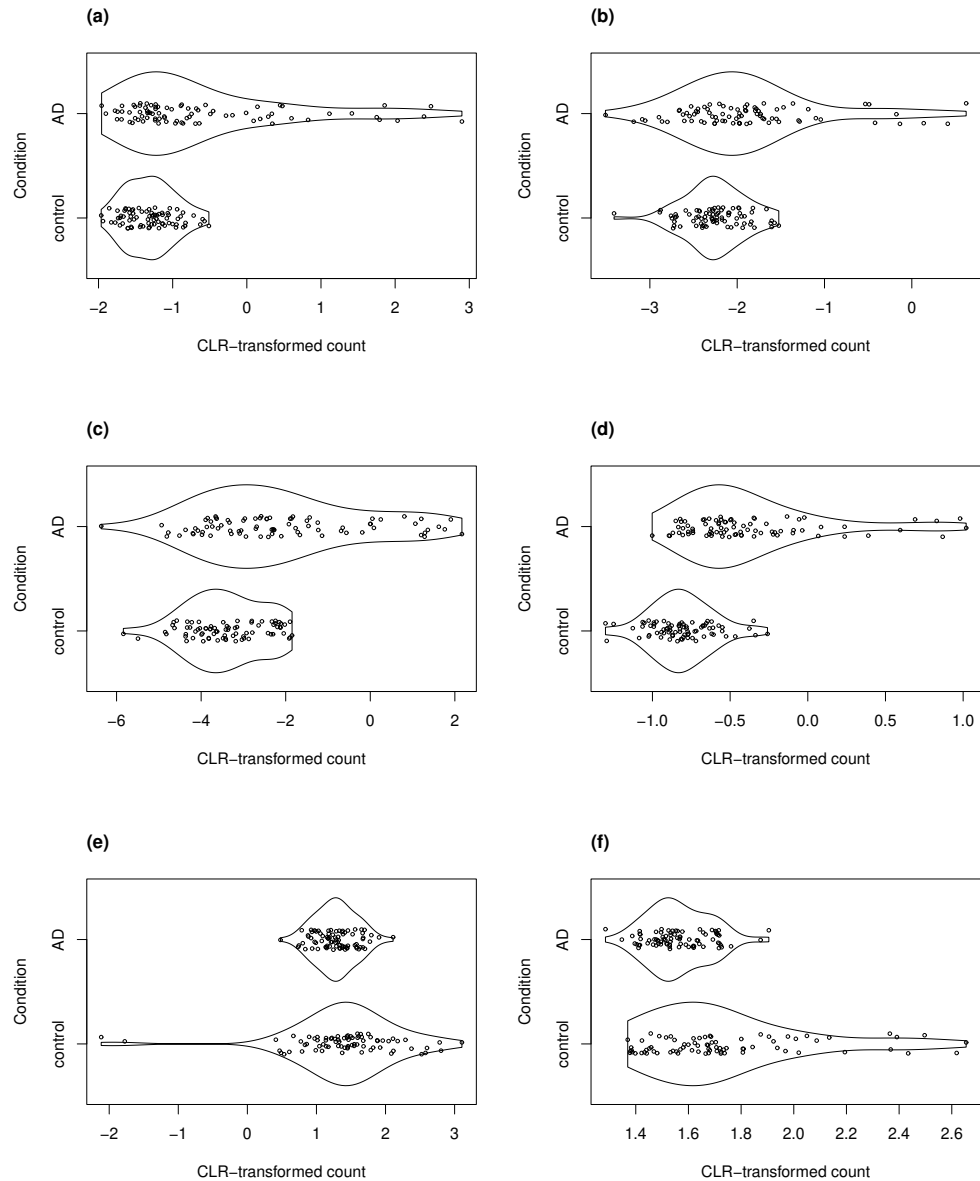




**Figure 5.** Venn diagrams of DV genes detected by clrDV, MDSeq and (a) GAMLSS-BH; (b) GAMLSS-BY for the control vs. AD comparison. Violin plots of the distribution of estimated  $\log_2(\text{SD ratio})$  of the DV genes detected using clrDV, MDSeq and (c) GAMLSS-BH; (d) GAMLSS-BY. Abbreviations: cgm = DV genes detected by clrDV, GAMLSS and MDSeq; c-g-m = DV genes detected by clrDV only; g-c-m = DV genes detected by GAMLSS-BH only; m-c-g = DV genes detected by MDSeq only.



**Figure 6.** Venn diagrams of DV genes detected by clrDV, MDSeq and (a) GAMLSS-BH; (b) GAMLSS-BY for the control vs. PSP comparison. Violin plots of the distribution of estimated  $\log_2(\text{SD ratio})$  of the DV genes detected using clrDV, MDSeq and (c) GAMLSS-BH; (d) GAMLSS-BY. Abbreviations: cgm = DV genes detected by clrDV, GAMLSS and MDSeq; c-g-m = DV genes detected by clrDV only; g-c-m = DV genes detected by GAMLSS-BH only; m-c-g = DV genes detected by MDSeq only.



**Figure 7.** Violin plots of selected DV genes detected in the control vs. AD comparison. (a) SLC47A1, (b) C2orf40, (c) SLPI, (d) LTBP2 have the largest SD ratio ( $> 2$ ); (e) GP1BB and (f) PELP1 have SD ratio about 0.4.