

Global Analysis of RNA-Seq Experiment: Multiple Data Sets & Multiple Genes

By

Bin Zhuo

A DISSERTATION

submitted to

Oregon State University

in partial fulfillment of  
the requirements for the  
degree of

Doctor of Philosophy in Statistics

Presented June 22, 2016  
Commencement June 2016



# AN ABSTRACT OF THE DISSERTATION OF

Bin Zhuo for the degree of Doctor of Philosophy in Statistics presented on  
June 22, 2016.

Title:

Global Analysis of RNA-Seq Experiment: Multiple Data Sets & Multiple Genes

Abstract approved:

---

Yanming Di

This is the abstract for my honors thesis. I'm going to start here.

Key Words: keyword1, keyword2, keyword3

Corresponding e-mail address: zhuob@oregonstate.edu

©Copyright by Bin Zhuo  
June 22, 2016  
All Rights Reserved

Global Analysis of RNA-Seq Experiment: Multiple Data Sets & Multiple Genes

By

Bin Zhuo

A DISSERTATION

submitted to

Oregon State University

in partial fulfillment of  
the requirements for the  
degree of

Doctor of Philosophy in Statistics

Presented June 22, 2016  
Commencement June 2016

Doctor of Philosophy in Statistics dissertation of Bin Zhuo presented on  
June 22, 2016

APPROVED:

---

Major Professor, representing Statistics

---

Chair of the Department of Statistics

---

Dean of the Graduate School

I understand that my project will become part of the permanent collection of Oregon State University libraries. My signature below authorizes release of my dissertation to any reader upon request.

---

Bin Zhuo, Author

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Biological question of interest . . . . .	1
1.1.1	Background . . . . .	1
1.1.2	Statistical issues . . . . .	3
1.1.3	Questions for this thesis . . . . .	6
1.1.3.1	Identifying stably expressed genes . . . . .	6
1.1.3.2	Estimating correlations of test statistics . . . . .	7
1.1.3.3	Adjusting for correlations in competitive gene set test . . . . .	7
1.2	Statistical Methods . . . . .	8
1.2.1	Generalized linear mixed models . . . . .	9
1.2.1.1	Classical linear models . . . . .	9
1.2.1.2	Linear mixed models . . . . .	10
1.2.1.3	Generalized linear models . . . . .	11
1.2.1.4	Generalized linear mixed models . . . . .	12
1.2.1.5	An example—Poisson log-linear mixed-effect model . . . . .	13
1.2.2	Estimation of generalized linear mixed models . . . . .	15
1.2.2.1	Likelihood function approach . . . . .	15
1.2.2.2	Estimation based on linearization . . . . .	17
1.2.2.3	Bayes approach . . . . .	20
1.2.2.4	Example of estimating parameters . . . . .	23
1.3	Multiple hypothesis testing . . . . .	24
1.4	Disertation Objective . . . . .	25
<b>2</b>	<b>Chapter 2</b>	<b>25</b>
<b>3</b>	<b>Chapter 3</b>	<b>25</b>

4	Conclusion	25
5	Chapter 3	25



**List of Figures**

1     Work flow of data preprocessing: from raw reads sequencing data to  
read counts. Raw data in this workflow are retrieved from the NCBI  
respository. Data processing is carried out based on two softwares—the  
SRA Toolkit [41] and the Rsubread aligner [47]. . . . . 4

## List of Tables

# 1. Introduction

## 1.1. Biological question of interest

### 1.1.1. Background

Gene is a piece of DNA that encodes a functional RNA or protein product, and is the basic physical and functional unit of heredity. The process by which genes are used to synthesize functional gene products is called *gene expression*. A gene is considered to be expressed in a cell or group of cells when a gene product is detected. These products can be transcribed messenger RNA (mRNA) and proteins for protein coding genes, or functional RNA species such as transfer RNA (tRNA) or small nuclear RNA (snRNA) for non-protein coding genes. Since the information encoded in a gene is first transcribed into RNA molecules, which is then used to make functional gene products, the RNAs transcribed in a certain condition reflect the current state of the cell.

#### **Why do people do expression analysis?**

In a typical gene expression experiment, researchers are usually interested in comparing expression levels of one or more genes from different sources. Factors for comparison can be *before vs after* effect in a drug treatment, *tumor vs normal* tissues in clinical study, or *wild type vs mutant* strains in plant research. Another important factor is time-course, where cells/tissues at different stages are sampled with the purpose of discovering temporal pattern of gene expression. There are many other types of experiment, each with specific factors of interest to be studied.

#### **What tools do people use to measure gene expression?**

The expression levels of a gene can be measured using techniques such as complementary DNA (cDNA) libraries, microarray analysis, RNA fingerprinting by arbitrary primed PCR (RAP-PCR), expressed sequence tag (EST) sequencing, serial analysis of gene expression (SAGE), and RNA sequencing (RNA-Seq) (see [12] for

a review). RNA-Seq, also known as *whole transcriptome shotgun sequencing* [57], is a next-generation sequencing (NGS) technology used to uncover the presence and quality of RNA in a biological sample. It is rapidly becoming technology of choice for transcriptome profiling over the past few years. The standard procedure of an RNA-Seq experiment runs as follows [21]: first, the RNAs in the biological sample are fragmented and reverse-transcribed into cDNAs; second, the cDNA fragments are amplified and sequenced in a high-throughput sequencing platform (e.g., Illumine 3000, <http://www.illumina.com>) to generate (up to) hundreds of millions of reads; third, those reads are mapped to a reference genome or a reference transcriptome. It is the number of reads aligned to each gene (referred to as “read count”) on the reference genome/transcriptome that quantifies the genes’ expression levels.

### **pros and cons about RNA-Seq**

RNA-Seq technology offers several key advantages over other methods [79], the most important of which are that it does not require prior knowledge of an organism for detecting transcripts, and that it is sensitive to genes expressed at either low or higher levels and thus provides higher dynamic range. The sequencing of RNA allows researchers to study the entire transcriptome of a species using only small amount of RNA. It has been demonstrated that a coordinated effort between RNA-Seq and real time PCR (RT-PCR) is one of the most effective ways to identify new exons [31]. However, one major challenge of this technique is data processing: RNA-Seq experiment produces a huge amount of reads (up to hundreds of millions per sample) and obtaining the expression profiles requires fast read mapping tools as well as a lot of computing resource [40, 44].

### **A workflow of pre-processing RNA-Seq data**

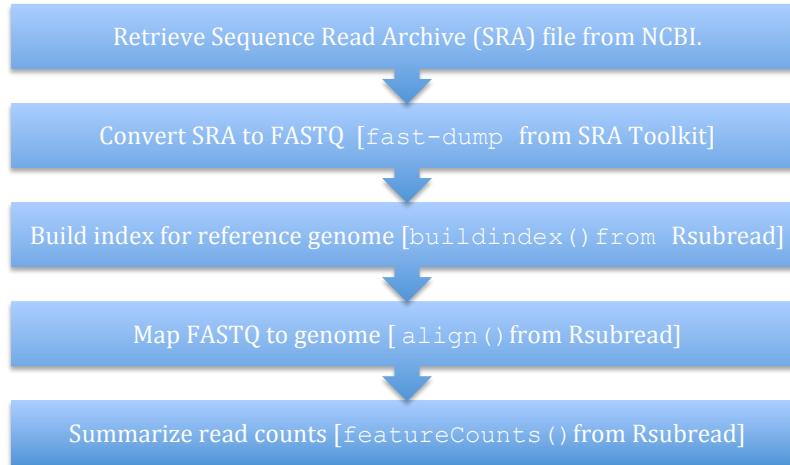
Preprocessing RNA-Seq data consists of two main steps: 1) mapping reads to the reference genome/transcriptome, and 2) summarizing read counts at given genomic feature (e.g., exon, gene or transcript) level. Read mapping is the first computational, and usually, the most time-consuming step in RNA-Seq data analysis. Currently,

there are many alignment tools available, for example, **Bowtie** [39, 40], **BWA** [42, 43], **Subread** [70] and **STAR** [17]. In all situations, an index of either the reference genome/transcriptome or the reads is built at the beginning using hash tables or Burrows-Wheeler transform (BWT) [11]. The index allows fast retrieval of the set of positions in the reference sequence where the reads are more likely to align. Once those positions are decided, alignment is performed in the candidate regions. The precision and speed of the alignment is mainly determined by the algorithm used in the alignment tool (see [30] or [45] for a review). After the reads have been aligned, the numbers of reads mapped to each unit of a specified genomic feature are counted, giving the estimate of the corresponding expression levels. This can be done using **HTSeq** [2] or **featureCounts** [46], among other options. Finally, a read count matrix is obtained in which each row represents a genomic feature unit and each column corresponds to a biological sample.

In this research work, we assemble an in-house pipeline to process RNA-Seq data sets based on the R [62] platform. This pipeline, modified from a standard procedure given by Anders et al. [4], is designed to work for sequencing data available at the *National Center for Biotechnology Information* (NCBI, <http://www.ncbi.nlm.nih.gov/>). It uses the SRA (Sequence Read Archive) Toolkit [41] to convert SRA files to FASTQ files, the **Subread** aligner [47] to map reads to the reference genome, and then the **featureCounts** [46] to summarize counts (see Figure 1 for the work flow). We will use this pipeline to process multiple RNA-Seq data sets that are needed in Chapter ??.

### 1.1.2. Statistical issues

The statistical analysis beginning from the read count matrix consists of three major parts: 1) normalization—adjusting for sources of bias between samples; 2) differential expression (DE) analysis—testing whether the expression levels of a gene are associated with treatment or experimental variables; and 3) gene set test—detecting which



**Figure 1:** Work flow of data preprocessing: from raw reads sequencing data to read counts. Raw data in this workflow are retrieved from the NCBI repository. Data processing is carried out based on two softwares—the SRA Toolkit [41] and the Rsubread aligner [47].

biological pathways are over-represented with DE genes.

## Normalization

Despite the optimistic claim that RNA-Seq does not need sophisticated normalization [79], many works have shown that normalization of count data is highly desirable to account for various sources of bias between samples before accessing differential expression [3, 16, 29, 65, 66, 67]. Normalization is needed for adjusting differences in sequencing depths or library sizes (total number of mapped reads for each biological sample) due to chance variation in sample preparation. In DE analysis, gene expression levels are often estimated from relative read frequencies. Therefore, normalization is also needed to account for the apparent reduction or increase in relative read frequencies of non-differentially expressed genes simply to accommodate the increased or decreased relative frequencies of truly DE genes. Currently there are many normalization methods available, such as the trimmed mean of M-values (TMM) [67], the DESeq normalization [3], and remove unwanted variation (RUV) [65].

## DE analysis

Identification of DE genes is the key task in many gene expression studies. DE analysis uncovers the association between expression levels of a gene and the covariates of interest. The covariates could either be categorical (e.g., treatment/control status, cell types), or continuous (e.g., reagent concentration, time). For example, to understand the effect of a drug, one might ask which genes are *up-regulated* (increased expression levels) or *down-regulated* (decreased expression levels) between treatment and control groups? Finding these genes will help researchers to understand the cause of a disease and to develop effective medicine. In recent years, many statistical tools have been developed for DE detection (methods review can be found in [63, 69, 71]) in RNA-Seq experiments. Most of those approaches are based on Poisson [51, 78] or Negative Binomial (NB) distribution [3, 15, 59, 68, 83] because RNA-Seq expression data are present in the form of counts. The NB distribution based models are more popular for their flexibility to deal with *over-dispersion* (a.k.a. extra-Poisson variation) that are often observed in RNA-Seq expression data.

## Gene set test

DE analysis evaluates each individual gene separately, but it fails to provide insights into biological mechanisms since genes may be correlated and function together. For this reason, *gene set test* is a frequently used technique that enables researchers to examine an ensemble of genes simultaneously and thus improves interpretability of DE results. Gene set test is the assessment of the association between a set of DE genes, which are significantly correlated with treatment or experimental design variables, and a prior set of genes, which are biologically related. Depending on the definition of the null hypothesis, there are two types of gene set test: the *self-contained* test and the *competitive* test [24]. A self-contained test examines a set of genes by a fixed standard without reference to other genes in the genome (see, for example, [26, 25, 35, 77, 81]). A competitive test compares DE genes in the test set to those

not in the test set [76, 82, 84]. The competitive gene set test is much more popular among genomic literatures [23, 24].

### 1.1.3. Questions for this thesis

In this thesis, we focus on three aspects of gene expression analysis: identifying stably expressed genes from multiple RNA-Seq data sets (Chapter ??); estimating correlations between test statistics via sample correlations [NEED TO REVISE] (Chapter 2); and adjusting for correlations in competitive gene set test (Chapter 5).

#### 1.1.3.1 Identifying stably expressed genes

Many of the current normalization methods, for example, TMM [67] and DESeq [3] normalizations, assume that the majority of genes are not DE within the experiment under investigation. However, this assumption could be violated for some experiments where over 50% of the genes' expression levels are altered by the treatments [50, 80]. The consequence with such assumption can be alleviated if one could identify a set of stably expressed genes whose expression levels are stable across different experimental conditions. This motivates us to identify such a set of genes by exploring a large number of existing RNA-Seq data sets.

In microarray studies, there have been many attempts to find reference genes for normalization. Traditionally, the *house-keeping genes* are used as reference genes. However, a number of works have shown that house-keeping genes are not necessarily stably expressed according to numerical stability measures (see, for example, [13, 36]). Another choice, the *spike-in genes*, is not reliable for normalization due to the same issue [65]. A popular approach has been to search from large sets of experiments for reference genes [13, 14, 22, 27, 72] whose expression stability are evaluated by some numerical stability measure. Validation experiments (e.g. reverse transcription-PCR) show that reference genes identified by numerical methods generally outperform house-keeping genes or spike-in genes in terms of expression stability [13, 32]. We will



follow the strategy of quantifying gene expression stability by numerical measures and identify stably expressed genes.

Identifying stably expressed genes not only helps count normalization, but also improves interpretability and comparability of RNA-Seq experiments in integrative analysis. Since genes are measured by relative frequencies, we argue that DE is a relative concept: when a normalization procedure is applied to a single data set, it effectively uses an implicit reference set of genes. Furthermore, making the reference set explicit will be beneficial during DE analysis, because often times biologists compare results from one experiment to those from others experiments whose data are publicly available.

### **1.1.3.2 Estimating correlations of test statistics**

NEED SOMETHING HERE

### **1.1.3.3 Adjusting for correlations in competitive gene set test**

Competitive gene set test compares DE genes in the set against those in its complementary set. A number of statistical methodologies have been developed for this purpose (literature reviews can be found in [33, 37, 56]). Broadly speaking, all of the competitive gene set tests fall into two categories based on whether they assume independence of expression profiles among genes. In earlier literatures, the inter-gene correlations were not taken care of in the enrichment analysis procedure, such as SigPathway [76], PAGE [38], MRSGE [55] or the  $2 \times 2$  contingency-table-based tests [1, 34, 85]. However, it has been argued that such test procedures will result in inflated type I error [20, 23, 24, 82, 84], as genes within a gene set are often co-expressed and function together.

Several approaches have been proposed to address inter-gene correlation problems in competitive gene set test. One attempt is to evaluate the significance of the test set by permuting sample labels [20, 23, 75]. Sample permutation does not require

an explicit understanding of the underlying correlation structure among genes, and is therefore supposed to protect the test against such correlations. One very famous example of this kind is the *gene set enrichment analysis* (GSEA) procedure [75]. Yet, sample permutation method has been criticized for several reasons: first, it cannot be applied to experiments having small number of biological replicates (e.g., three samples each for a two-group comparison, which is common in RNA-Seq experiments); second, it is computationally intensive since tens of thousands of DE tests are involved in each permutation; third, and most importantly, it implicitly alters the null hypothesis being tested and makes the null and alternative difficult to be characterized [24, 37, 82]. Another attempt has been to incorporate the inter-gene correlations into the formulation of gene set test procedure [82, 84]. CAMERA [82] estimates a *variance inflation factor* (VIF) from sample correlation (after the treatment effects removed), and then includes it in its test statistic to assess the significance of the gene set. The same VIF has also been used by QuSAGE [84] to adjust for inter-gene correlations. However, accurate estimation of VIF relies on the assumption that correlation between any two gene-level statistics are almost the same as correlation between their corresponding expression levels. In Chapter 2, we will demonstrate that this assumption is easily violated when differentially expressed genes are present, and as a remedy, we will propose a new gene set test procedure in Chapter 5.

## 1.2. Statistical Methods

We have mentioned earlier that RNA-Seq data are essentially present in the form of count matrices. Therefore it might not be appropriate to impose normal distribution on gene expression profiles, especially when the sample size is small. Generalized linear models (GLMs) are a natural choice for analyzing RNA-Seq data. In addition, to account for random terms in biological experiments, GLMs are sometimes extended to generalized linear mixed models (GLMMs). In this section, we will first describe the formulation GLMMs, and then review commonly used methods for parameter

estimation under this framework.

### 1.2.1. Generalized linear mixed models

GLMMs are a natural generalization of classical linear models. To illustrate this point, we will begin with classical linear models, and discuss how to generalize them to linear mixed models and then to GLMMs by relaxing different layers of assumptions.

#### 1.2.1.1 Classical linear models

In a classical linear model, a vector  $\mathbf{y}$  of  $n$  observations is assumed to be a realization of random variable  $\mathbf{Y}$  whose components are identically distributed with mean  $\boldsymbol{\mu}$ . The systematic part of this model is a specification of the mean  $\boldsymbol{\mu}$  over a few unknown parameters [53]. In the context of classical linear models, the mean is a function of  $p$  covariates  $\mathbf{X}_1, \dots, \mathbf{X}_p$ ,

$$\boldsymbol{\mu} = \beta_0 + \sum_{i=1}^p \beta_i \mathbf{X}_i \quad (1)$$

where  $\beta$ 's are unknown parameters and need to be estimated from data. For  $j$ th observation  $Y_j$ , we specify  $\epsilon_j$ , a random term, to allow for measurement error. Assuming a linear relationship between response  $Y_j$  and predictors  $(x_{1j}, \dots, x_{pj})$ , we present the linear model

$$Y_j = \beta_0 + \beta_1 x_{1j} + \dots + \beta_p x_{pj} + \epsilon_j \quad (2)$$

It is often required that  $\epsilon_i$ 's meet *Gauss-Markov* assumption,

$$E(\epsilon_i) = 0, \text{ Var}[\epsilon_i] = \sigma^2 < \infty, \text{ Cov}[\epsilon_i, \epsilon_j] = 0, \forall i \neq j. \quad (3)$$

In practice, the error term is frequently, if not always, assumed to be normally distributed,

$$\boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I}). \quad (4)$$

### 1.2.1.2 Linear mixed models

The Gauss-Markov assumption in equation (3) is vulnerable in practice, for example, nonconstant variance, or correlated data where  $\text{Cov}[\epsilon_i, \epsilon_j] \neq 0$ . Equation (2) in either case, without loss of generality, can be expressed in matrix form as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad E[\boldsymbol{\epsilon}] = \mathbf{0}, \quad \text{Cov}[\boldsymbol{\epsilon}] = \mathbf{V} \quad (5)$$

where  $\mathbf{V}$  is a known positive definite matrix. Let  $\mathbf{Y}^* = \mathbf{V}^{-1/2}\mathbf{Y} = \mathbf{V}^{-1/2}\mathbf{X}\boldsymbol{\beta} + \mathbf{V}^{-1/2}\boldsymbol{\epsilon}$ . It follows that  $\text{Cov}(\mathbf{Y}^*) = \mathbf{I}$  and the techniques in classical linear models are readily applicable to estimate  $\boldsymbol{\beta}$ . However, this method relies on the assumption that  $\mathbf{V}$  is known which is rarely, if ever, given. On the other hand, the structure of  $\mathbf{V}$ , which depends on experiment setup, can often be specified by a few unknown parameters.

Nonindependence can occur in the form of serial correlation or cluster correlation [64, chapter 17]. Serial correlation usually exists in experiments with repeated measurements—multiple measurements taken from a response variable on the same experimental unit. Several covariance structures are available for implementation (for more details, see Littell et al. [48, chapter 5]). Cluster correlation is present when measurements of a response variable are grouped in some way. In many situations, the covariance of cluster correlated data can be specified using an extension of standard linear model by

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_1\mathbf{u}_1 + \cdots + \mathbf{Z}_q\mathbf{u}_q + \boldsymbol{\epsilon} \quad (6)$$

equation (6) differs from equation (5) only in the  $\mathbf{Z}_i\mathbf{u}_i$  terms, which is the key part of *linear mixed models*. The  $\mathbf{Z}_i$  are known  $n \times p_i$  full rank matrices, usually used to specify membership of predictors in various subgroups. The most important innovation in this model is that instead of estimating  $\mathbf{u}_i$ 's as fixed parameters, we assume them to be unknown random quantities, and  $E[\mathbf{u}_i] = \mathbf{0}$ ,  $\text{Cov}[\mathbf{u}_i] = \sigma_i^2 \mathbf{I}_{p_i}$  for  $i = 1, \dots, q$ . It is, in many cases, reasonable to require that  $\mathbf{u}_i$  are mutually independent, and that  $\mathbf{u}_i$

is independent of  $\epsilon$  for  $i = 1, \dots, q$ . If we further impose normal distribution on the random terms and errors, then equation (6) can be casted in a Bayesian framework,

$$\begin{aligned} \mathbf{y} | \mathbf{u}_1, \dots, \mathbf{u}_q &\sim N_n(\mathbf{X}\boldsymbol{\beta} + \sum_{i=1}^q \mathbf{Z}_i \mathbf{u}_i, \sigma^2 \mathbf{I}_n), \\ \mathbf{u}_i &\sim N_{p_i}(0, \sigma_i^2 \mathbf{I}_{p_i}). \end{aligned} \quad (7)$$

The modeling issues are: (a) estimation of variance components  $\sigma_i^2$  and  $\sigma^2$ ; (b) estimation of random effects  $u_i$  if needed. For the variance component estimation, there are primarily three approaches: (i) procedures based on expected mean squares from analysis of variance (ANOVA); (ii) maximum likelihood (ML); and (iii) restricted/residual maximum likelihood (REML). For more details, see Littell et al. [48, Chapter 1].

### 1.2.1.3 Generalized linear models

We can take a different perspective of classical linear models by arranging equations (1)–(3) into three parts, following the notations of McCullagh and Nelder [53, Chapter 2],

- (i) the *random component*  $Y_j$  has constant variance  $\sigma^2$  and  $E[Y_j] = \mu_j$ .
- (ii) the *systematic component*—the linear predictor  $\eta_j$  is modeled by covariates

$$\begin{aligned} \mathbf{x}_j &=: x_{1j}, \dots, x_{pj}, \\ \eta_j &= \sum_{i=1}^p \beta_i x_{ij} = \mathbf{x}_j \boldsymbol{\beta}. \end{aligned} \quad (8)$$

- (iii) the *link function* relates the random components and the systematic components by

$$\eta_j = g(\mu_j). \quad (9)$$

The classical linear models fits within this framework if we assume that the random components  $Y_j$ 's are independent and normally distributed, and that the link function is identity (i.e.,  $g(\mu_j) = \mu_j$ ).

We can extend part (i)—by allowing  $Y_j$  to come from an exponential family (e.g., Poisson, Gamma or Binomial distribution), and part (iii)—by requiring the link function to be monotonic differentiable (e.g.,  $g(\mu_j) = \log \mu_j$ ). These two extensions lead to the *generalized linear models* (GLMs), a framework that is especially suitable when a normal distribution is no longer appropriate to be assumed on the response.

#### 1.2.1.4 Generalized linear mixed models

Generalized linear mixed models (GLMMs) is a further extension of GLMs that incorporates random components into part (ii), represented in a matrix notation

$$\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \sum_{i=1}^q \mathbf{Z}_i \mathbf{u}_i \quad (10)$$

where  $\mathbf{Z}_i$  and  $\mathbf{u}_i$  are specified in equation (6).

To formally present GLMMs, we start with the conditional distribution of  $\mathbf{y}$  given  $\mathbf{u}$ . It is typical to assume that vector  $\mathbf{y}$  consists of conditionally independent elements, each coming from the exponential family (or similar to the exponential family),

$$\begin{aligned} y_j | \mathbf{u} &\sim \text{indep. } f_{Y_j | \mathbf{u}}(y_j | \mathbf{u}), \\ f_{Y_j | \mathbf{u}}(y_j; \theta, \phi | \mathbf{u}) &= \exp \left[ \frac{y_j \theta_j - b(\theta_j)}{a_j(\phi)} + c(y_j, \phi) \right]. \end{aligned} \quad (11)$$

It can be verified that the conditional mean of  $y_j$  is related to  $\theta_j$  in equation (11) by the identity  $\mu_j = \partial b(\theta_j) / \partial \theta_j$ . The transformation of the mean allows us to model the fixed and the random factors by a linear model

$$\begin{aligned} E[y_j | \mathbf{u}] &= \mu_j, \\ g(\mu_j) = \eta_j &= \mathbf{X}_j \boldsymbol{\beta} + \mathbf{Z}_j \mathbf{u}. \end{aligned} \quad (12)$$

Finally, we assign a distribution to the random effects

$$\mathbf{U} \sim \phi_{\mathbf{U}}(\mathbf{u}), \quad (13)$$

which completes the specification of GLMMs. It is often, if not always, assumed that  $\mathbf{u}$  come from a normal distribution.

#### 1.2.1.5 An example—Poisson log-linear mixed-effect model

We will illustrate one specific type of GLMMs—the Poisson log-linear mixed-effect model in the context of RNA-seq experiments. Suppose we have RNA-Seq expression profiles (in the form of counts) randomly selected from three experiments conducted in three different labs. For each experiment, there are two treatments and two biological replicates for each treatment. We are not interested in the specific levels of treatment, but focus more on the overall variation of treatments. In this sense, the treatment effects are also considered as random. For a single gene, let  $Y_{jkl} \sim \text{Poisson}(\mu_{jkl})$  be the read count for  $j$ th biological sample from  $k$ th treatment of  $l$ th experiment. The link function  $\eta_{jkl} = \log(\mu_{jkl})$  relates the mean  $\mu_{jkl}$  to the linear predictors by equation (12)

$$\log(\mu_{jkl}) = \log(N_{jkl}R_{jkl}) + \xi + a_j + b_{k(j)} + \epsilon_{jkl}, \quad (14)$$

where  $N_{jkl}R_{jkl}$  is the normalized library size (total number of read counts mapped to the genome),  $j = 1, 2, 3$ ,  $k = 1, 2$  and  $l = 1, 2$ ; the random terms  $a_j \sim N(0, \sigma_1^2)$ ,  $b_{k(j)} \sim N(0, \sigma_2^2)$  and  $\epsilon_{jkl} \sim N(0, \sigma_0^2)$  represent the experimental, treatment and sample effects respectively, and are mutually independent. If the observations are sorted by experiment and then by treatment nested in experiment, we can present the model in the form of equation (10), with  $\boldsymbol{\beta} = (\log[N_{111}R_{111}] + \xi, \dots, \log[N_{223}R_{223}] + \xi)^T$ ,  $\mathbf{u} =$





### 1.2.2. Estimation of generalized linear mixed models

There are three general approaches for estimating parameters under GLMM settings [58, Chapter 7]: (i) using numerical method to approximate the integrals for the likelihood functions and obtaining the estimating equations; (ii) linearization of the conditional mean and then iteratively applying linear mixed model techniques to the approximated model; (iii) Bayesian approach.

In the following discussion, we assume conditional distribution of  $\mathbf{Y}$  given  $\mathbf{u}$  is  $f_Y(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u})$ , the link function is  $\boldsymbol{\eta} = g(\boldsymbol{\mu})$ , and  $\boldsymbol{\eta}$  relates the covariates by equation (12). We also assume that the random term  $\mathbf{u}$  have some distribution  $\mathbf{U} \sim \phi_{\mathbf{U}}(\mathbf{u}|\boldsymbol{\Sigma})$ .

#### 1.2.2.1 Likelihood function approach

It is straightforward to write down the likelihood function of  $\mathbf{Y}$  by first obtaining the joint likelihood of  $(\mathbf{Y}, \mathbf{u})$  and then integrating out the random term  $\mathbf{u}$ ,

$$L(\mathbf{Y}|\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \int f(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u})\phi(\mathbf{u}|\boldsymbol{\Sigma})d\mathbf{u}. \quad (15)$$

A major challenge in estimating GLMMs is the integration in equation (15) over the  $n$ -dimensional distribution of  $\mathbf{u}$ . Numerical approximation are usually used in evaluating the integral. In this part we will discuss the *Gauss-Hermite* (GH) quadrature which is recognized as a higher order Laplace approximation [49]. Gauss-Hermite quadrature is used for integrals of the form  $\int_{-\infty}^{\infty} f(x)e^{-x^2}dx$  that can be approximated by a weighted sum of  $f(x)$ :

$$\int_{-\infty}^{\infty} f(x)e^{-x^2}dx \approx \sum_{i=1}^m w_i f(x_i) \quad (16)$$

In equation (16),  $x_i$ 's are the zeros of  $m$ th order Hermite polynomial

$$H_m(x) = (-1)^m \exp\left(\frac{x^2}{2}\right) \frac{d^m}{dx^m} \exp\left(-\frac{x^2}{2}\right)$$

and  $w_i$  are the corresponding weights. For a Hermite polynomial of degree  $m$ ,  $x_i$  and  $w_i$  can be calculated as

$$x_i = i\text{th zero of } H_m(x), \quad w_i = \frac{2^{m-1}m!\sqrt{\pi}}{m^2[H_{m-1}(x_i)]^2}. \quad (17)$$

Equation (16) gives the exact numerical value for all polynomials up to degree of  $2m - 1$ . An improved version of the regular Gauss-Hermite quadrature is to center and scale the quadrature points by the empirical Bayes estimate of the random effects and the Hessian matrix from the Bayes estimate suboptimization [49]. This procedure is called *Adaptive Gauss-Hermite* (AGH) quadrature [60].

The AGH quadrature starts with maximizing the integrand  $h(\mathbf{u}|\mathbf{y}, \boldsymbol{\beta}, \boldsymbol{\Sigma}) := f(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u})\phi(\mathbf{u}|\boldsymbol{\Sigma})$  in equation (15) with respect to the random term  $\mathbf{u}$ . The resulting estimate  $\hat{\mathbf{u}}^{(n)}$  at iteration  $n$  is the joint posterior modes for the random effects. Because  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$  are unknown, they are replaced by the current estimates  $\hat{\boldsymbol{\beta}}^{(n)}$  and  $\hat{\boldsymbol{\Sigma}}^{(n)}$ . The Hessian matrix  $\hat{\mathbf{H}}^{(n)}$  can be obtained by evaluating the second order partial derivatives of  $\log(h(\mathbf{u}|\mathbf{y}, \hat{\boldsymbol{\beta}}^{(n)}, \hat{\boldsymbol{\Sigma}}^{(n)}))$  at  $\hat{\mathbf{u}}^{(n)}$ . Consequently,  $\hat{\boldsymbol{\Omega}}^{(n)} = -\hat{\mathbf{H}}^{(n)}$  is the estimated covariance matrix for the random effects posterior modes. It follows from equation (15) that for the  $i$ th cluster

$$L(\mathbf{Y}_i|\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \int f(\mathbf{y}_i|\boldsymbol{\beta}, \mathbf{u})\phi(\mathbf{u}|\boldsymbol{\Sigma})d\mathbf{u} = \int \frac{f(\mathbf{y}_i|\boldsymbol{\beta}, \mathbf{u})\phi(\mathbf{u}|\boldsymbol{\Sigma})}{\phi(\mathbf{u}|\hat{\mathbf{u}}^{(n)}, \hat{\boldsymbol{\Omega}}^{(n)})}\phi(\mathbf{u}|\hat{\mathbf{u}}^{(n)}, \hat{\boldsymbol{\Omega}}^{(n)})d\mathbf{u} \quad (18)$$

Let  $m$  be the number of quadrature points (i.e., the order of the Hermite polynomial) in each dimension for each random effect term, and  $Q$  the number of random effects. If  $\mathbf{x} = (x_1, \dots, x_m)$  are the nodes for standard Gauss-Hermite quadrature, and  $\mathbf{x}_j^* = (x_{j1}, \dots, x_{jQ})$  is a point on the  $Q$  dimensional quadrature grid, then the centered and scaled nodes are

$$\mathbf{a}_j^* = \hat{\mathbf{u}}^{(n)} + \sqrt{2}[\hat{\boldsymbol{\Omega}}^{(n)}]^{1/2}\mathbf{x}_j^* \quad (19)$$

The centered and scaled nodes, along with the Gauss-Hermite quadrature weights

$\mathbf{w} = (w_1, \dots, w_m)$  are used to construct the  $Q$  dimensional integral of equation (18), approximated by

$$\begin{aligned} L(\mathbf{y}_i | \boldsymbol{\beta}, \boldsymbol{\Sigma}) &\approx \sum_{j_1=1}^m \cdots \sum_{j_Q=1}^m \frac{f(\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{a}_j^*) \phi(\mathbf{a}_j^* | \boldsymbol{\Sigma})}{\phi(\mathbf{a}_j^* | \hat{\mathbf{u}}^{(n)}, \hat{\boldsymbol{\Omega}}^{(n)})} w_{j_1} \cdots w_{j_Q} \\ &= (2)^{Q/2} |\hat{\boldsymbol{\Omega}}^{(n)}|^{1/2} \sum_{j_1=1}^m \cdots \sum_{j_Q=1}^m \left[ f(\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{a}_j^*) \phi(\mathbf{a}_j^* | \boldsymbol{\Sigma}) \prod_{k=1}^Q w_{j_k} \exp(x_{j_k}^2) \right]. \end{aligned} \quad (20)$$

Thus the multidimensional unbounded integral is approximated by a finite summations. Now that the likelihood has the form of equation (20), a number of numerical methods (e.g. Newton-Raphson or Fisher's scoring) can be used to estimate  $(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ .

It should be noted, however, as the number of dimension  $Q$  increases, the computational burden for approximating equation (20) grows exponentially since the total number of nodes is  $m^Q$ . Therefore it is difficult to implement AGH procedure with more than three random effects [8].

### 1.2.2.2 Estimation based on linearization

Under GLMM framework, we have some conditional distribution of  $\mathbf{Y}$  given  $\mathbf{u}$ . Without loss of generality, we assume

$$\begin{aligned} E[\mathbf{Y} | \mathbf{u}] &= \boldsymbol{\mu} = g^{-1}(\boldsymbol{\eta}) = g^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}), \\ \text{Var}[\mathbf{Y} | \mathbf{u}] &= \mathbf{S}, \end{aligned} \quad (21)$$

where  $\mathbf{u} \sim N(\mathbf{0}, \boldsymbol{\Sigma})$ . The linearization is done by Taylor expansion of equation (21) about estimates  $\boldsymbol{\eta}$ . Two approaches proposed by Breslow and Clayton [9]—the *penalized quasi-likelihood* (PQL) and the *marginal quasi-likelihood* (MQL)—may be used for this purpose.

**Penalized Quasi-likelihood** The PQL procedure uses a first order Taylor expansion of  $\boldsymbol{\beta}$  and  $\mathbf{u}$ , at  $\tilde{\boldsymbol{\beta}}$  and  $\tilde{\mathbf{u}}$ , respectively

$$g^{-1}(\boldsymbol{\eta}) \approx g^{-1}(\hat{\boldsymbol{\eta}}) + \tilde{\boldsymbol{\Omega}}_{PQL}(\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}}), \quad (22)$$

where  $\tilde{\boldsymbol{\Omega}}_{PQL}$  is an  $n \times n$  diagonal matrix whose  $(i, i)$  entry is  $\partial g^{-1}(\boldsymbol{\eta}_i) / \partial \boldsymbol{\eta}_i$  evaluated at  $\tilde{\boldsymbol{\eta}} = \mathbf{X}\tilde{\boldsymbol{\beta}} + \mathbf{Z}\tilde{\mathbf{u}}$ . Multiplying both sides by  $\tilde{\boldsymbol{\Omega}}_{PQL}^{-1}$ , equation (22) can be rearranged as

$$\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} \approx \tilde{\boldsymbol{\Omega}}_{PQL}^{-1}[g^{-1}(\boldsymbol{\eta}) - g^{-1}(\tilde{\boldsymbol{\eta}})] + \mathbf{X}\tilde{\boldsymbol{\beta}} + \mathbf{Z}\tilde{\mathbf{u}}. \quad (23)$$

Note that the right hand side of equation (23) is just the expected value, conditioning on  $\tilde{\boldsymbol{\beta}}$  and  $\tilde{\mathbf{u}}$ , of the pseudo-response

$$\tilde{\mathbf{Y}} = \tilde{\boldsymbol{\Omega}}_{PQL}^{-1}[\mathbf{Y} - g^{-1}(\tilde{\boldsymbol{\eta}})] + \mathbf{X}\tilde{\boldsymbol{\beta}} + \mathbf{Z}\tilde{\mathbf{u}}, \quad (24)$$

whose variance-covariance matrix given  $\mathbf{u}$  is

$$\text{Var}[\tilde{\mathbf{Y}}|\mathbf{u}] = \tilde{\boldsymbol{\Omega}}_{PQL}^{-1} \text{Var}[\mathbf{Y}|\mathbf{u}] \tilde{\boldsymbol{\Omega}}_{PQL}^{-1} = \tilde{\boldsymbol{\Omega}}_{PQL}^{-1} \mathbf{S} \tilde{\boldsymbol{\Omega}}_{PQL}^{-1}. \quad (25)$$

Then we can consider the model

$$\tilde{\mathbf{Y}} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon} \quad (26)$$

which is a linear mixed model with pseudo response  $\tilde{\mathbf{Y}}$  with covariance matrix

$$\mathbf{W} = \text{Var}[\tilde{\mathbf{Y}}|\mathbf{u}] = \mathbf{Z}\boldsymbol{\Sigma}\mathbf{Z}' + \tilde{\boldsymbol{\Omega}}_{PQL}^{-1} \mathbf{S} \tilde{\boldsymbol{\Omega}}_{PQL}^{-1}. \quad (27)$$

Model (26) has exactly the same form as the linear mixed models (see Section 1.2.1.2), except that an estimate of  $(\boldsymbol{\beta}, \mathbf{u})$  is needed for calculating the pseudo-response  $\tilde{\mathbf{Y}}$  in equation (24). An iterative procedure can be used to estimate the parameters in

model (26) by substituting raw data  $\mathbf{y}$  for  $\tilde{\mathbf{y}}$  and identity matrix  $\mathbf{I}$  for  $\mathbf{S}$  as starting values. Techniques for fitting LMM such as REML can be readily applied to estimate variance components  $\Sigma$ , upon which  $\hat{\mathbf{W}}$  is calculated. The estimate for  $\beta$  is given by

$$\hat{\beta} = (\mathbf{X}^T \hat{\mathbf{W}}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{W}}^{-1} \mathbf{X} \tilde{\mathbf{y}}, \quad (28)$$

and the estimate for random effect is

$$\hat{\mathbf{u}} = \hat{\Sigma} \mathbf{Z} \hat{\mathbf{W}}^{-1} (\tilde{\mathbf{y}} - \mathbf{X} \hat{\beta}). \quad (29)$$

Then the pseudo-response is updated and the procedure is repeated until convergence is reached for fixed effects and variance components. Note that equation (29) estimates a vector of random effect. For this reason, PQL is also referred to as *subject-specific* estimate procedure.

**Marginal Quasi-likelihood** One of the motivation for MQL is that usually one is more interested in estimating the marginal mean of the response than estimating the conditional mean as is done by equation (29) in PQL. Since  $E[\eta|\mathbf{u}] = \mathbf{X}\beta + \mathbf{Z}\mathbf{u}$ , the unconditional mean is  $E[\eta] = E[E(\eta|\mathbf{u})] = \mathbf{X}\beta$ . A first-order Taylor expansion of  $E[\mathbf{Y}|\mathbf{u}]$  about  $\mathbf{X}\beta$  is given by

$$E[\mathbf{Y}|\mathbf{u}] = g^{-1}(\eta) \approx g^{-1}(\mathbf{X}\beta) + \tilde{\Omega}_{MQL}(\eta - \mathbf{X}\beta) \quad (30)$$

where  $\tilde{\Omega}_{MQL}$  is evaluated at  $\mathbf{X}\beta$  (recall that for PQL,  $\tilde{\Omega}_{PQL}$  is evaluated at  $\mathbf{X}\beta + \mathbf{Z}\mathbf{u}$ ). The unconditional expected value of  $\mathbf{Y}$  is approximately  $g^{-1}(\mathbf{X}\beta)$  by equation (30). The variance of  $\mathbf{Y}$  can then be derived from the relation  $\text{Var}(\mathbf{Y}) = E[\text{Var}(\mathbf{Y}|\mathbf{u})] + \text{Var}[E(\mathbf{Y}|\mathbf{u})]$ , which yields

$$\text{Var}[\mathbf{Y}] = \tilde{\Omega}_{MQL} \mathbf{Z} \Sigma \mathbf{Z}' \tilde{\Omega}_{MQL}' + \mathbf{S}_{\eta_0}. \quad (31)$$

A linearization performed at  $\boldsymbol{\eta}_0 = \mathbf{X}\boldsymbol{\beta}_0$  leads to

$$g^{-1}(\boldsymbol{\eta}) \approx g^{-1}(\mathbf{X}\boldsymbol{\beta}_0) + \tilde{\boldsymbol{\Omega}}_{MQL}(\boldsymbol{\eta} - \mathbf{X}\boldsymbol{\beta}_0), \quad (32)$$

and multiplying both sides by  $\tilde{\boldsymbol{\Omega}}_{MQL}^{-1}$ , equation (32) then can be arranged to

$$\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} \approx \tilde{\boldsymbol{\Omega}}_{MQL}^{-1}[g^{-1}(\boldsymbol{\eta}) - g^{-1}(\boldsymbol{\eta}_0)] + \mathbf{X}\boldsymbol{\beta}_0. \quad (33)$$

Defining the pseudo-response  $\tilde{\mathbf{Y}}_{MQL}$  as

$$\tilde{\mathbf{Y}}_{MQL} = \tilde{\boldsymbol{\Omega}}_{MQL}^{-1}[\mathbf{Y} - g^{-1}(\boldsymbol{\eta}_0)] + \mathbf{X}\boldsymbol{\beta}_0, \quad (34)$$

we next consider the linear mixed model

$$\tilde{\mathbf{Y}}_{MQL} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon},$$

where  $\text{Var}(\boldsymbol{\epsilon})$  is given by equation (31). The estimating procedure for fixed effect parameter  $\boldsymbol{\beta}$  and variance component  $\boldsymbol{\Sigma}$  is the same as those in the PQL approach. Note that the pseudo-response is not a function of  $\mathbf{u}$  any more, so updating this quantity does not require calculating the random effects  $\mathbf{u}$ . Accordingly, the MQL approach is also referred to as *population-averaged* estimate approach.

Breslow and Lin [10] and Pinheiro and Chao [61] point out that PQL approach may lead to asymptotically biased estimates and hence to inconsistency. It is not recommended to use simple PQL method in practice.

### 1.2.2.3 Bayes approach

As mentioned earlier, for models with higher dimensional integrals, it is not practical to evaluate the likelihood function by AGH procedure. For mixed models, a typical strategy is to treat the random effects to be missing data. Following this idea, the

problem of estimating variance components associated with random effects can be simplified. Denote the *complete data* as  $\mathbf{v} = (\mathbf{y}, \mathbf{u})$ , the log-likelihood of  $\mathbf{v}$  can be expressed as

$$\log \pi(\boldsymbol{\beta}, \boldsymbol{\Sigma} | \mathbf{v}) = \log f(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}) + \log \phi(\mathbf{u} | \boldsymbol{\Sigma}) \quad (35)$$

The optimal solution for parameters in equation (35) can be obtained by *expectation-maximization* (EM) algorithm. The EM algorithm consists of two steps, readily implemented as follows:

1. **E-Step.** At  $(k + 1)$ th iteration given  $\boldsymbol{\beta}^{(k)}$  and  $\boldsymbol{\Sigma}^{(k)}$ , calculate

$$\begin{aligned} E_{\boldsymbol{\beta}^{(k)}}[\log f(\boldsymbol{\beta}, \boldsymbol{\Sigma} | \mathbf{v}) | \mathbf{y}] &= Q_1(\boldsymbol{\beta}, \boldsymbol{\beta}^{(k)}), \\ E_{\boldsymbol{\Sigma}^{(k)}}[\log \phi(\boldsymbol{\Sigma} | \mathbf{v}) | \mathbf{y}] &= Q_2(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}^{(k)}). \end{aligned} \quad (36)$$

2. **M-Step.** Maximize  $Q_1$  and  $Q_2$  to update  $\boldsymbol{\beta}^{(k+1)}$  and  $\boldsymbol{\Sigma}^{(k+1)}$ .

The **E** and **M** steps are alternated until convergence is reached. Unfortunately, the expectations in equation (36) cannot be computed in closed form for GLMMs. However, they may be approximated by *Markov chain Monte Carlo* (MCMC). In light of this, McCulloch [54] developed a Monte Carlo EM (MCEM) algorithm.

The Metropolis-Hastings algorithm is used for drawing samples from difficult-to-calculate density functions. For Metropolis algorithm, a proposal distribution  $g(\mathbf{u})$  is selected, from which an initial value of  $\mathbf{u}$  is drawn. The new candidate value  $\mathbf{u}' = (u_1, u_2, \dots, u_{k-1}, u'_k, u_{k+1}, \dots, u_Q)$ , which has all elements the same as previous values except the  $k$ th, is accepted (as opposed to keeping the previous value) with probability

$$A_k(\mathbf{u}', \mathbf{u}) = \min \left\{ 1, \frac{f(\mathbf{u}' | \mathbf{y}, \boldsymbol{\beta}, \boldsymbol{\Sigma})g(\mathbf{u})}{f(\mathbf{u} | \mathbf{y}, \boldsymbol{\beta}, \boldsymbol{\Sigma})g(\mathbf{u}')} \right\}. \quad (37)$$

If we choose  $g(\mathbf{u}) = \phi(\mathbf{u}|\Sigma)$ , the ratio term in equation (37) can be simplified to

$$\begin{aligned}
& \frac{f(\mathbf{u}'|\mathbf{y}, \beta, \Sigma)g(\mathbf{u})}{f(\mathbf{u}|\mathbf{y}, \beta, \Sigma)g(\mathbf{u}')} \\
&= \left[ \frac{f(\mathbf{u}', \mathbf{y}|\beta, \Sigma)}{f(\mathbf{y}|\beta, \Sigma)} \phi(\mathbf{u}|\Sigma) \right] \bigg/ \left[ \frac{f(\mathbf{u}, \mathbf{y}|\beta, \Sigma)}{f(\mathbf{y}|\beta, \Sigma)} \phi(\mathbf{u}'|\Sigma) \right] \\
&= \frac{f(\mathbf{y}|\mathbf{u}', \beta, \Sigma)\phi(\mathbf{u}'|\Sigma)\phi(\mathbf{u}|\Sigma)}{f(\mathbf{y}|\mathbf{u}, \beta, \Sigma)\phi(\mathbf{u}|\Sigma)\phi(\mathbf{u}'|\Sigma)} \\
&= \frac{f(\mathbf{y}|\mathbf{u}', \beta, \Sigma)}{f(\mathbf{y}|\mathbf{u}, \beta, \Sigma)}
\end{aligned} \tag{38}$$

The MCEM procedure combines the EM steps and Metropolis algorithm in estimating the fixed parameters and variance components, summarized as follows:

1. Choose the starting value of  $\beta^{(0)}, \Sigma^{(0)}$ . Set  $b = 0$
2. Generate the sequence  $\mathbf{u}^{(1)}, \mathbf{u}^{(2)}, \dots, \mathbf{u}^{(B)}$  from the conditional distribution of  $\mathbf{u}$  given  $\mathbf{y}$  with Metropolis algorithm.
3. Maximize  $\sum_{b=1}^B \log f(\mathbf{y}|\mathbf{u}^{(b)}, \beta)/B$  and  $\sum_{b=1}^B \log \phi(\mathbf{u}^{(b)}|\Sigma)/B$  to obtain  $\beta^{(m+1)}$  and  $\Sigma^{(m+1)}$
4. Iterate between step 2 and 3 until convergence is reached.

This method can be easily extended to allow for multiple random effects. Yet the advantage comes at a price. A major drawback of MCEM is the computational intensity. First, the convergence of EM algorithm is usually very slow, especially at the neighborhood of maximum of marginal likelihood. Second, the chain in Metropolis algorithm has to run long enough for reliable estimation.

In the Bayes framework, there are other alternatives to estimate the parameters and variance components, for example, *Monte Carlo Newton-Raphson* (MCNR) [54] and MCMC [28].



### 1.2.2.4 Example of estimating parameters

We will demonstrate the estimating procedure with the Poisson log-linear mixed-effect model discussed in Section 1.2.1. The estimation procedure starts from the joint density function of  $\mathbf{Y} = (Y_{jkl})'$  given  $\boldsymbol{\mu} = (\mu_{jkl})'$ ,

$$f(\mathbf{y}|\boldsymbol{\mu}) = \prod_{j,k,l} f(y_{jkl}|\mu_{jkl}) = \prod_{j,k,l} \frac{[\mu_{jkl}]^{y_{jkl}} \exp(-\mu_{jkl})}{y_{jkl}!}. \quad (39)$$

A re-expression of equation (14) in matrix form gives

$$\log \boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_1\mathbf{a} + \mathbf{Z}_2\mathbf{b} + \mathbf{I}_{12}\boldsymbol{\epsilon}.$$

Therefore  $\boldsymbol{\mu} \sim \log N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma})$  where  $\boldsymbol{\mu}_0 = \boldsymbol{\xi} + \log(\mathbf{NR})$  and  $\boldsymbol{\Sigma} = \sigma_1^2 \mathbf{Z}_1 \mathbf{Z}_1' + \sigma_2^2 \mathbf{Z}_2 \mathbf{Z}_2' + \sigma_0^2 \mathbf{I}_{12}$ . The density function of  $\boldsymbol{\mu}$  is then

$$f(\boldsymbol{\mu}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}) = \prod_{j,k,l} \mu_{jkl}^{-1} \cdot \frac{1}{\sqrt{(2\pi)^{12}|\boldsymbol{\Sigma}|}} \exp\left[-\frac{1}{2}(\log \boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\log \boldsymbol{\mu} - \boldsymbol{\mu}_0)\right]. \quad (40)$$

Since  $Y_{jkl} \sim \text{Poisson}(\mu_{jkl})$ , by combining equation (39) and (40), we obtain the joint distribution of  $\mathbf{Y}$  and  $\boldsymbol{\mu}$ ,

$$\begin{aligned} f(\mathbf{y}, \boldsymbol{\mu}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}) \\ = \frac{1}{\sqrt{(2\pi)^{12}|\boldsymbol{\Sigma}|}} \exp\left[-\mathbf{1}^T \boldsymbol{\mu} - \frac{1}{2}(\log \boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\log \boldsymbol{\mu} - \boldsymbol{\mu}_0)\right] \prod_{jkl} \frac{[\mu_{jkl}]^{y_{jkl}-1}}{y_{jkl}!}. \end{aligned} \quad (41)$$

Therefore we can obtain the likelihood function of or the marginal distribution of  $\mathbf{Y}$  by integrating out the random components  $\mathbf{u}$ ,

$$L(\xi, \sigma_1^2, \sigma_2^2, \sigma_0^2|\mathbf{y}) = f(\mathbf{y}|\boldsymbol{\xi}, \boldsymbol{\Sigma}) = \int_{\mathbf{a}, \mathbf{b}, \boldsymbol{\epsilon}} f(\mathbf{y}, \mathbf{a}, \mathbf{b}, \boldsymbol{\epsilon}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}) d\mathbf{a} d\mathbf{b} d\boldsymbol{\epsilon}. \quad (42)$$

The integral in equation (42) can be approximated by adaptive Gaussian-Hermite (AGH) quadrature or MCMC. For AGH quadrature, we first approximate the like-

likelihood by equation (20) and then estimate  $\boldsymbol{\theta} = (\xi, \sigma_0^2, \sigma_1^2, \sigma_2^2)'$  by maximizing the resulting likelihood. The R package `lme4` [5] has an inbuilt function `glmer()` for this purpose. The MCMC procedure has also been implemented in several packages, such as the `Rstan` [73] and the `MCMCPack` [52].

### 1.3. Multiple hypothesis testing

Multiple hypothesis testing procedures deal with type I error rates in a family of tests. The problems arise when we consider a set of statistical inference simultaneously. For each of the individual tests or confidence intervals, there is a type I error which can be controlled by the experimenter. If the family of tests contains one or more true null hypotheses, the probability of rejecting one or more of these true null increases.

While traditional multiple testing procedures focus on modest number of tests, a different set of techniques are needed for large-scale inference, in which tens or even hundreds of thousands of tests are performed simultaneously. For example, in genomics study, expression levels of 50,000 genes for each of 100 individuals can be measured using modern technologies such as microarray or RNA-Sequencing. In testing differential expression (DE), 50,000 tests need to be conducted against the null that there is no DE between treatment/control. This has brought new challenge to the field of multiple hypothesis testing. Benjamini and Hochberg [6] points out that the control of familywise error rate (FWER), i.e. the probability of making one or more false discovery in a set of tests, tends to have substantially less power.

*False discovery rate* (FDR), introduced by Benjamini and Hochberg [6], is the expected proportion of false positives among all significant calls (null rejected). FDR has been studied extensively ([7, 18, 19, 74] and more) over the past two decades. FDR is equivalent to FWER [6] when all hypotheses are true but smaller if there are some true discoveries to be made. We will focus our attention on FDR in this part.

Let  $m$ ,  $m_0$  and  $m_1$  be the number of tests, true nulls and true alternatives respectively. Let also  $F$  and  $T$  be the number of true nulls and true alternatives among

$S$  tests that are declared as significant. Table (??) shows the relation among them.

The FDR is

	Called significance	Called not significant	Total
Null True	F	$m_0 - F$	$m_0$
Alternative true	T	$m_1 - T$	$m_1$
total	S	$m - S$	$m$

#### 1.4. Disertation Objective

## 2. Chapter 2

## 3. Chapter 3

## 4. Conclusion

## 5. Chapter 3

Chapter 3.

## References

- [1] Alexa, A. and Rahnenfuhrer, J. (2010). topgo: enrichment analysis for gene ontology. *R package version*, 2(0).
- [2] Anders, S. (2010). Htseq: Analysing high-throughput sequencing data with python. URL <http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>.
- [3] Anders, S. and Huber, W. (2010). Differential expression analysis for sequence count data. *Genome biology*, 11(10):R106.
- [4] Anders, S., McCarthy, D. J., Chen, Y., Okoniewski, M., Smyth, G. K., Huber, W., and Robinson, M. D. (2013). Count-based differential expression analysis of rna sequencing data using r and bioconductor. *Nature protocols*, 8(9):1765–1786.
- [5] Bates, D., Maechler, M., and Bolker, B. (2012). lme4: Linear mixed-effects models using s4 classes.
- [6] Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 289–300.
- [7] Benjamini, Y. and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Annals of statistics*, pages 1165–1188.
- [8] Bolker, B. M., Brooks, M. E., Clark, C. J., Geange, S. W., Poulsen, J. R., Stevens, M. H. H., and White, J.-S. S. (2009). Generalized linear mixed models: a practical guide for ecology and evolution. *Trends in ecology & evolution*, 24(3):127–135.
- [9] Breslow, N. E. and Clayton, D. G. (1993). Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association*, 88(421):9–25.

- [10] Breslow, N. E. and Lin, X. (1995). Bias correction in generalised linear mixed models with a single component of dispersion. *Biometrika*, 82(1):81–91.
- [11] Burrows, M. and Wheeler, D. (1994). A block-sorting lossless data compression algorithm. In *DIGITAL SRC RESEARCH REPORT*. Citeseer.
- [12] Casassola, A., Brammer, S. P., Chaves, M. S., Ant, J., Grando, M. F., et al. (2013). Gene expression: A review on methods for the study of defense-related gene differential expression in plants. *American Journal of Plant Sciences*, 2013.
- [13] Czechowski, T., Stitt, M., Altmann, T., Udvardi, M. K., and Scheible, W.-R. (2005). Genome-wide identification and testing of superior reference genes for transcript normalization in arabidopsis. *Plant physiology*, 139(1):5–17.
- [14] Dekkers, B. J., Willems, L., Bassel, G. W., van Bolderen-Veldkamp, R. M., Ligterink, W., Hilhorst, H. W., and Bentsink, L. (2012). Identification of reference genes for rt-qpcr expression analysis in arabidopsis and tomato seeds. *Plant and Cell Physiology*, 53(1):28–37.
- [15] Di, Y., Schafer, D. W., Cumbie, J. S., and Chang, J. H. (2011). The NBP negative binomial model for assessing differential gene expression from RNA-Seq. *Stat. Appl. Genet. Mol. Biol.*, 10(1):1–28.
- [16] Dillies, M.-A., Rau, A., Aubert, J., Hennequet-Antier, C., Jeanmougin, M., Servant, N., Keime, C., Marot, G., Castel, D., Estelle, J., et al. (2013). A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. *Brief. Bioinform.*, 14(6):671–683.
- [17] Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and Gingeras, T. R. (2013). Star: ultrafast universal rna-seq aligner. *Bioinformatics*, 29(1):15–21.

- [18] Efron, B. (2004). Large-scale simultaneous hypothesis testing. *Journal of the American Statistical Association*, 99(465).
- [19] Efron, B. (2010). *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*, volume 1. Cambridge University Press.
- [20] Efron, B. and Tibshirani, R. (2007). On testing the significance of sets of genes. *The Annals of Applied Statistics*, pages 107–129.
- [21] Finotello, F. and Di Camillo, B. (2015). Measuring differential gene expression with rna-seq: challenges and strategies for data analysis. *Briefings in functional genomics*, 14(2):130–142.
- [22] Frericks, M. and Esser, C. (2008). A toolbox of novel murine house-keeping genes identified by meta-analysis of large scale gene expression profiles. *BBA-Gene Regul. Mech.*, 1779(12):830–837.
- [23] Gatti, D. M., Barry, W. T., Nobel, A. B., Rusyn, I., and Wright, F. A. (2010). Heading down the wrong pathway: on the influence of correlation within gene sets. *BMC Genomics*, 11(1):574.
- [24] Goeman, J. J. and Bühlmann, P. (2007). Analyzing gene expression data in terms of gene sets: methodological issues. *Bioinformatics*, 23(8):980–987.
- [25] Goeman, J. J., Oosting, J., Cleton-Jansen, A.-M., Anninga, J. K., and Van Houwelingen, H. C. (2005). Testing association of a pathway with survival using gene expression data. *Bioinformatics*, 21(9):1950–1957.
- [26] Goeman, J. J., Van De Geer, S. A., De Kort, F., and Van Houwelingen, H. C. (2004). A global test for groups of genes: testing association with a clinical outcome. *Bioinformatics*, 20(1):93–99.

- [27] Gur-Dedeoglu, B., Konu, O., Bozkurt, B., Ergul, G., Seckin, S., and Yulug, I. G. (2009). Identification of endogenous reference genes for qRT-PCR analysis in normal matched breast tumor tissues. *Oncol. Res.*, 17(8):353–365.
- [28] Hadfield, J. D. et al. (2010). Mcmc methods for multi-response generalized linear mixed models: the mcmcglmm r package. *Journal of Statistical Software*, 33(2):1–22.
- [29] Hansen, K. D., Irizarry, R. A., and Zhijin, W. (2012). Removing technical variability in RNA-Seq data using conditional quantile normalization. *Biostatistics*, 13(2):204–216.
- [30] Hatem, A., Bozdağ, D., Toland, A. E., and Çatalyürek, Ü. V. (2013). Benchmarking short sequence mapping tools. *BMC bioinformatics*, 14(1):1.
- [31] Howald, C., Tanzer, A., Chrast, J., Kokocinski, F., Derrien, T., Walters, N., Gonzalez, J. M., Frankish, A., Aken, B. L., Hourlier, T., et al. (2012). Combining rt-pcr-seq and rna-seq to catalog all genic elements encoded in the human genome. *Genome research*, 22(9):1698–1710.
- [32] Hruz, T., Wyss, M., Docquier, M., Pfaffl, M. W., Masanetz, S., Borghi, L., Verbrugghe, P., Kalaydjieva, L., Bleuler, S., Laule, O., et al. (2011). Refgenes: identification of reliable and condition specific reference genes for rt-qpcr data normalization. *BMC genomics*, 12(1):156.
- [33] Huang, D. W., Sherman, B. T., and Lempicki, R. A. (2009). Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Research*, 37(1):1–13.
- [34] Huang, D. W., Sherman, B. T., Tan, Q., Kir, J., Liu, D., Bryant, D., Guo, Y., Stephens, R., Baseler, M. W., Lane, H. C., et al. (2007). David bioinformatics resources: expanded annotation database and novel algorithms to better extract biology from large gene lists. *Nucleic acids research*, 35(suppl 2):W169–W175.

- [35] Huang, Y.-T. and Lin, X. (2013). Gene set analysis using variance component tests. *BMC Bioinformatics*, 14(1):210.
- [36] Huggett, J., Dheda, K., Bustin, S., and Zumla, A. (2005). Real-time RT-PCR normalisation; strategies and considerations. *Genes Immun.*, 6(4):279–284.
- [37] Khatry, P., Sirota, M., and Butte, A. J. (2012). Ten years of pathway analysis: current approaches and outstanding challenges. *PLoS Computational Biology*, 8(2):e1002375.
- [38] Kim, S.-Y. and Volsky, D. J. (2005). PAGE: parametric analysis of gene set enrichment. *BMC Bioinformatics*, 6(1):144.
- [39] Langmead, B. and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat. Methods*, 9(4):357–359.
- [40] Langmead, B., Trapnell, C., Pop, M., Salzberg, S. L., et al. (2009). Ultrafast and memory-efficient alignment of short dna sequences to the human genome. *Genome biol*, 10(3):R25.
- [41] Leinonen, R., Sugawara, H., and Shumway, M. (2010). The sequence read archive. *Nuc. Acids Res.*, page gkq1019.
- [42] Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with bwa-mem. *arXiv preprint arXiv:1303.3997*.
- [43] Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics*, 25(14):1754–1760.
- [44] Li, H. and Durbin, R. (2010). Fast and accurate long-read alignment with burrows–wheeler transform. *Bioinformatics*, 26(5):589–595.
- [45] Li, H. and Homer, N. (2010). A survey of sequence alignment algorithms for next-generation sequencing. *Briefings in bioinformatics*, 11(5):473–483.



- [46] Liao, Y., Smyth, G. K., and Shi, W. (2013a). featurecounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics*, page btt656.
- [47] Liao, Y., Smyth, G. K., and Shi, W. (2013b). The subread aligner: fast, accurate and scalable read mapping by seed-and-vote. *Nucleic acids research*, 41(10):e108–e108.
- [48] Littell, R. C., Stroup, W. W., Milliken, G. A., Wolfinger, R. D., and Schabenberger, O. (2006). *SAS for mixed models*. SAS institute.
- [49] Liu, Q. and Pierce, D. A. (1994). A note on gauss—hermite quadrature. *Biometrika*, 81(3):624–629.
- [50] Lovén, J., Orlando, D. A., Sigova, A. A., Lin, C. Y., Rahl, P. B., Burge, C. B., Levens, D. L., Lee, T. I., and Young, R. A. (2012). Revisiting global gene expression analysis. *Cell*, 151(3):476–482.
- [51] Marioni, J. C., Mason, C. E., Mane, S. M., Stephens, M., and Gilad, Y. (2008). Rna-seq: an assessment of technical reproducibility and comparison with gene expression arrays. *Genome research*, 18(9):1509–1517.
- [52] Martin, A. D., Quinn, K. M., and Park, J. H. (2011). Mcmcpack: Markov chain monte carlo in r.
- [53] McCullagh, P. and Nelder, J. A. (1989). Generalized linear models.
- [54] McCulloch, C. E. (1997). Maximum likelihood algorithms for generalized linear mixed models. *Journal of the American statistical Association*, 92(437):162–170.
- [55] Michaud, J., Simpson, K. M., Escher, R., Buchet-Poyau, K., Beissbarth, T., Carmichael, C., Ritchie, M. E., Schütz, F., Cannon, P., Liu, M., et al. (2008). Integrative analysis of RUNX1 downstream pathways and target genes. *BMC Genomics*, 9(1):363.

- [56] Mishra, P., Törönen, P., Leino, Y., and Holm, L. (2014). Gene set analysis: limitations in popular existing methods and proposed improvements. *Bioinformatics*, 30(19):2747–2756.
- [57] Morin, R. D., Bainbridge, M., Fejes, A., Hirst, M., Krzywinski, M., Pugh, T. J., McDonald, H., Varhol, R., Jones, S. J., and Marra, M. A. (2008). Profiling the hela s3 transcriptome using randomly primed cdna and massively parallel short-read sequencing. *Biotechniques*, 45(1):81.
- [58] Myers, R. H., Montgomery, D. C., Vining, G. G., and Robinson, T. J. (2012). *Generalized linear models: with applications in engineering and the sciences*, volume 791. John Wiley & Sons.
- [59] Oberg, A. L., Bot, B. M., Grill, D. E., Poland, G. A., and Therneau, T. M. (2012). Technical and biological variance structure in mrna-seq data: life in the real world. *BMC genomics*, 13(1):304.
- [60] Pinheiro, J. C. and Bates, D. M. (1995). Approximations to the log-likelihood function in the nonlinear mixed-effects model. *Journal of computational and Graphical Statistics*, 4(1):12–35.
- [61] Pinheiro, J. C. and Chao, E. C. (2006). Efficient laplacian and adaptive gaussian quadrature algorithms for multilevel generalized linear mixed models. *Journal of Computational and Graphical Statistics*, 15(1).
- [62] R Core Team (2015). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- [63] Rapaport, F., Khanin, R., Liang, Y., Pirun, M., Krek, A., Zumbo, P., Mason, C. E., Socci, N. D., and Betel, D. (2013). Comprehensive evaluation of differential gene expression analysis methods for rna-seq data. *Genome Biol*, 14(9):R95.

- [64] Rencher, A. C. and Schaallje, G. B. (2008). *Linear models in statistics*. John Wiley & Sons.
- [65] 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 Biotech*, 32(9):896–902.
- [66] Risso, D., Schwartz, K., Sherlock, G., and Dudoit, S. (2011). Gc-content normalization for rna-seq data. *BMC bioinformatics*, 12(1):480.
- [67] Robinson, M. D., Oshlack, A., et al. (2010). A scaling normalization method for differential expression analysis of RNA-Seq data. *Genome Biol.*, 11(3):R25.
- [68] Robinson, M. D. and Smyth, G. K. (2007). Moderated statistical tests for assessing differences in tag abundance. *Bioinformatics*, 23(21):2881–2887.
- [69] Seyednasrollah, F., Laiho, A., and Elo, L. L. (2015). Comparison of software packages for detecting differential expression in rna-seq studies. *Briefings in bioinformatics*, 16(1):59–70.
- [70] Shi, W. and Liao, Y. (2013). Subread/rsubread users guide.
- [71] Soneson, C. and Delorenzi, M. (2013). A comparison of methods for differential expression analysis of rna-seq data. *BMC bioinformatics*, 14(1):1.
- [72] Stamova, B. S., Apperson, M., Walker, W. L., Tian, Y., Xu, H., Adamczyk, P., Zhan, X., Liu, D.-Z., Ander, B. P., Liao, I. H., et al. (2009). Identification and validation of suitable endogenous reference genes for gene expression studies in human peripheral blood. *BMC Med. Genom.*, 2(1):49.
- [73] Stan Development Team (2016). *RStan: the R interface to Stan, Version 2.9.0*.
- [74] Storey, J. D. and Tibshirani, R. (2003). Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences*, 100(16):9440–9445.

- [75] Subramanian, A., Tamayo, P., Mootha, V. K., Mukherjee, S., Ebert, B. L., Gillette, M. A., Paulovich, A., Pomeroy, S. L., Golub, T. R., Lander, E. S., et al. (2005). Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United States of America*, 102(43):15545–15550.
- [76] Tian, L., Greenberg, S. A., Kong, S. W., Altschuler, J., Kohane, I. S., and Park, P. J. (2005). Discovering statistically significant pathways in expression profiling studies. *Proceedings of the National Academy of Sciences of the United States of America*, 102(38):13544–13549.
- [77] Tsai, C.-A. and Chen, J. J. (2009). Multivariate analysis of variance test for gene set analysis. *Bioinformatics*, 25(7):897–903.
- [78] Wang, L., Feng, Z., Wang, X., Wang, X., and Zhang, X. (2010). Degseq: an R package for identifying differentially expressed genes from rna-seq data. *Bioinformatics*, 26(1):136–138.
- [79] Wang, Z., Gerstein, M., and Snyder, M. (2009). Rna-seq: a revolutionary tool for transcriptomics. *Nature Reviews Genetics*, 10(1):57–63.
- [80] Wu, D., Hu, Y., Tong, S., Williams, B. R., Smyth, G. K., and Gantier, M. P. (2013a). The use of miRNA microarrays for the analysis of cancer samples with global miRNA decrease. *RNA*, 19(7):876–888.
- [81] Wu, D., Lim, E., Vaillant, F., Asselin-Labat, M.-L., Visvader, J. E., and Smyth, G. K. (2010). ROAST: rotation gene set tests for complex microarray experiments. *Bioinformatics*, 26(17):2176–2182.
- [82] Wu, D. and Smyth, G. K. (2012). Camera: a competitive gene set test accounting for inter-gene correlation. *Nucleic Acids Research*, 40(17):e133–e133.

- [83] Wu, H., Wang, C., and Wu, Z. (2013b). A new shrinkage estimator for dispersion improves differential expression detection in rna-seq data. *Biostatistics*, 14(2):232–243.
- [84] Yaari, G., Bolen, C. R., Thakar, J., and Kleinstein, S. H. (2013). Quantitative set analysis for gene expression: a method to quantify gene set differential expression including gene-gene correlations. *Nucleic Acids Research*, page gkt660.
- [85] Ye, J., Fang, L., Zheng, H., Zhang, Y., Chen, J., Zhang, Z., Wang, J., Li, S., Li, R., Bolund, L., et al. (2006). Wego: a web tool for plotting go annotations. *Nucleic acids research*, 34(suppl 2):W293–W297.