

DIFFERENTIAL COUNT ANALYSIS WITH A TOPIC MODEL

PETER CARBONETTO*

1. Motivation, and overview of methods. The aim of this document is to derive, from first principles, a method for analysis of differential gene expression using a topic model (also known as “grade of membership model” [4]). This method may have other uses — say, to identify “key words” in a topic modeling analysis of text documents — but since our main motivation is analysis of gene expression data, we describe the methods with that application in mind.

For motivation, we begin with the “log-fold change” statistic commonly used in microarray and RNA sequencing experiments to quantify expression differences between two conditions (e.g., [3, 7]). The log-fold change for gene j and condition k is a ratio of two conditional expectations,

$$(1) \quad \text{lfc}(j, k) \equiv \log_2 \frac{E[x_j \mid \text{condition} = k]}{E[x_j \mid \text{condition} \neq k]},$$

where x_j is the measured expression level of gene j .¹ The way in which the expression level x_j is defined can lead to different log-fold change statistics. For example, several popular methods for analyzing differential expression compare changes in *relative* expression (say, relative to total expression in a cell) [2, 5, 6, 8]. As we will see, a topic modeling perspective provides a natural way to analyze differential gene expression when comparing either relative or absolute expression levels.

Topic modeling brings a new twist to analysis of differential gene expression: In conventional differential expression analysis, each sample is assigned to a single condition; in topic modeling, the assignments to each topic are *proportional*. This needs to be accounted for in developing the new methods for differential expression analysis.

1.1. The binomial model. We begin with a simple binomial model that predicts expression of a single gene given the proportional assignments to a topic:

$$(2) \quad x_i \sim \text{Binom}(s_i, \pi_i).$$

Here, x_i is the expression level of the target gene in sample i , s_i is the total expression in sample i , and $\text{Binom}(n, \theta)$ denotes the binomial distribution with n trials and success probability θ . In this simple model, the binomial probabilities are defined as

$$(3) \quad \pi_i = (1 - q_i)p_0 + q_i p_1,$$

where $q_i \in [0, 1]$ is the (known) proportion of sample i that is attributed to the topic, and $p_0, p_1 \in [0, 1]$ are two unknowns to be estimated.²

*Dept. of Human Genetics and the Research Computing Center, University of Chicago, Chicago, IL

¹Defining x_{jk} as the total gene expression for gene j among all samples (expression profiles) in condition k , x_j as the total gene expression for gene j in all samples, n_k as the number of samples in condition k , and n as the total sample size, the log-fold change can be computed as $\text{lfc}(j, k) = \log_2 \left\{ \frac{x_{jk}}{x_j - x_{jk}} \times \frac{n - n_k}{n_k} \right\}$.

²This is a special case of a well-studied model in econometrics called the *linear probability model* [9]. The linear probability model would be typically written as $\pi_i = \beta_0 + q_i \beta$, where $\beta_0 = p_0$ and $\beta = p_1 - p_0 \in [-1, +1]$. This is a regression model for the binomial probability π_i , in which π_i increases linearly with topic proportion q_i .

Statistical inference with this simple binomial model implements analysis of differential gene expression. In particular, $\log_2(p_1/p_0)$ is the *log-fold change in relative expression*. To show that this is so, consider the following statistical process for generating the counts x_1, \dots, x_n :

- for $i = 1, \dots, n$
 1. for $t = 1, \dots, s_i$
 - (a) Sample a topic, $z_{it} \sim \text{Binom}(1, q_i)$.
 - (b) Sample a gene, $w_{it} | z_{it} \sim \begin{cases} \text{Binom}(1, p_1) & \text{if } z_{it} = 1 \\ \text{Binom}(1, p_0) & \text{otherwise.} \end{cases}$
 2. Generate the final gene count, $x_i \leftarrow w_{i1} + \dots + w_{is_i}$.

In this statistical process, $q_i = p(z_{it} = 1)$ is the topic probability, $p_1 = p(w_{it} = 1 | z_{it} = 1)$ is the conditional probability that the gene is expressed given membership to the topic, and $p_0 = p(w_{it} = 1 | z_{it} = 0)$ is the probability that the gene is expressed when not belonging to the topic. Therefore, we have

$$(4) \quad \log_2 \frac{p_1}{p_0} = \log_2 \frac{p(w_{it} = 1 | z_{it} = 1)}{p(w_{it} = 1 | z_{it} = 0)}.$$

This is the log-fold change statistic for relative expression given proportional topic assignments q_1, \dots, q_n . The binomial model (2) can in fact be derived from this statistical process (proof not shown). Therefore, estimating p_0, p_1 for the binomial model (2) provides an estimate of the log-fold change (4).

Before continuing, we point out that the s_i 's, in practice, do not need to be the total expression for each sample i (*i.e.*, the total counts); for example, some researchers have suggested setting s_i to some pre-defined quantile of the sample's nonzero count distribution (e.g., [2]). This approach has been motivated in settings where a small number of genes are much more highly expressed than the others, and therefore these few genes have a large effect relative expression levels. In short, the binomial model is quite flexible, and can accommodate different relative differential expression analyses defines the s_i 's. The only constraint on s_i is that it cannot be smaller than x_i .

1.2. The Poisson model. A Poisson model similar to the binomial model (2) leads to a method for estimating log-fold change in absolute expression levels. We proceed in a similar way. The Poisson model predicts expression x_i given the topic proportions q_i :

$$(5) \quad x_i \sim \text{Poisson}(t_i \lambda_i),$$

in which the Poisson rates are defined as

$$(6) \quad \lambda_i = (1 - q_i)f_0 + q_i f_1,$$

and the two unknowns to be estimated are $f_0, f_1 \geq 0$. (The scalars $t_i > 0$ are additional “size factors” that give the model more flexibility; for now, assume $t_i = 1$ for all $i = 1, \dots, n$.) Observe that, unlike the binomial model, the unknowns for the Poisson model are not constrained to be in $[0, 1]$, and so they cannot represent probabilities. In the following, we show that $\log_2(f_1/f_0)$ is the log-fold change in absolute expression.

Consider the following process for generating the counts x_1, \dots, x_n :

- for $i = 1, \dots, n$

- | | |
|---|--------------------------------------|
| 1. $a_i \sim \text{Poisson}(f_1)$ | Sample the within-topic gene count. |
| 2. $b_i \sim \text{Poisson}(f_0)$ | Sample the outside-topic gene count. |
| 3. $a'_i \sim \text{Binom}(a_i, q_i)$ | Subsample the within-topic genes. |
| 4. $b'_i \sim \text{Binom}(b_i, 1 - q_i)$ | Subsample the outside-topic genes. |
| 5. $x_i \leftarrow a'_i + b'_i$ | Generate the final gene count. |

In this generative process, $f_1 = E[a_i]$ represents the within-topic gene rate, and $f_0 = E[b_i]$ is the outside-topic gene rate, and therefore

$$(7) \quad \log_2 \frac{f_1}{f_0} = \log_2 \frac{E[a_i]}{E[b_i]}$$

is the *log-fold change statistic in absolute expression given proportional topic assignments* q_1, \dots, q_n . For intuition, when all the topic proportions q_i are 0 or 1, this statistical process simplify to $x_i \sim \text{Poisson}(f_1)$ if $q_i = 1$, and $x_i \sim \text{Poisson}(f_0)$ if $q_i = 0$, and so (7) would reduce to the ratio of the mean expression levels inside and outside the topic. The Poisson model (5) can be derived from this statistical process, and so estimating f_0, f_1 for the Poisson model (5) provides an estimate of the log-fold change (7).

A practical question is the choice of size factors t_i . In the standard analysis, one would set $t_i = 1$ for all $i = 1, \dots, n$, but there may be situations in which other settings for t_i are warranted. For example, you may know in advance that one of the samples, say, the first sample, $i = 1$, due to some technical error, produced twice as much gene expression (say, two cells were accidentally combined into one during sample preparation, and so expression was measured in both). With $t_1 = 1$, this sample would bias the log-fold change statistics, whereas setting $t_1 = 2$ would “correct” this bias.

1.3. Binomial vs. Poisson: which to use. The binomial model (2) would be used to implement the log-fold change analysis for relative expression, and the Poisson model (5) would implement the log-fold change analysis for absolute expression. There is, as we have already established in related writeups, a close relationship between the Poisson and binomial models (the binomial being a special case of the multinomial), and in fact we have exploited that close relationship to use fast non-negative matrix algorithms to fit the multinomial topic model. In particular, it can be shown that the likelihood for the binomial model (2) is the same as the likelihood for this Poisson model:

$$\begin{aligned} x_i &\sim \text{Poisson}(s_i \pi_i) \\ y_i &\sim \text{Poisson}(s_i (1 - \pi_i)), \end{aligned}$$

where $y_i = s_i - x_i$. So the key difference between the relative and absolute expression analyses is actually not the use of the Poisson or binomial, *but the way in which the models are parameterized*. This difference in parameterization leads to analysis of absolute or relative expression. Therefore, we prefer to keep the development of the Poisson and binomial models separate.

The next sections derive the mathematical expressions needed to implement differential count analysis based on the binomial and Poisson models.

2. Binomial model derivations. *Add derivations here.*

3. Poisson model derivations. *Add derivations here.*

4. The multinomial topic model and Poisson non-negative matrix factorization. Here we briefly describe the multinomial topic model, and its connection to Poisson non-negative matrix factorization (Poisson NMF).

We begin with the “bag of words” description, which was used to describe the LDA model [1]. In this view, each document (or gene expression profile) i is represented as a vector of terms/genes, $w_i = (w_{i1}, \dots, w_{is_i})$, where s_i is the size of document i . (The order of the words or genes appearing in this vector doesn’t matter, hence the “bag of words.”) Each $w_{it} \in \{1, \dots, m\}$ is term/gene j with probability $p(w_{it} = j | z_{it} = k) = f_{jk}$, in which we have introduced z_{it} , a variable indicating which topic $k \in \{1, \dots, K\}$ the word/gene is drawn from. The topic indicator variables for document i are in turn generated according to $p(z_{it} = k) = l_{ik}$.

This process also defines a *multinomial* model for an $n \times m$ matrix of counts x_{ij} :

$$(8) \quad x_{i1}, \dots, x_{im} \sim \text{Multinom}(x_{i1}, \dots, x_{im}; s_i, \pi_i),$$

where $x_{ij} = \sum_{t=1}^{s_i} \delta_j(w_{it})$ is the number of times term/gene j appears in document/cell i , and the probabilities π_{ij} are weighted sums of the “factors” f_{jk} ,

$$(9) \quad \pi_{ij} = \sum_{k=1}^K l_{ik} f_{jk}.$$

The log-likelihood for the multinomial topic model, ignoring terms that do not depend on the model parameters, has a simple expression:

$$(10) \quad \log p(x) = \sum_{i=1}^n \sum_{j=1}^m x_{ij} \log(\sum_{k=1}^K l_{ik} f_{jk}).$$

As we have shown elsewhere, the multinomial topic model is closely related to a Poisson non-negative matrix factorization of the count data,

$$(11) \quad x_{ij} \sim \text{Poisson}(\lambda_{ij}),$$

where $\lambda_{ij} = \sum_{k=1}^K \hat{l}_{ik} \hat{f}_{jk}$. Given a Poisson NMF fit, an equivalent multinomial topic model can be easily recovered, as we have shown elsewhere.

5. Gene expression differences in topics. Returning to the question of assessing differential gene expression, there are two new twists when done in the context of topic modeling:

1. The cluster (topic) assignments are probabilistic.
2. The cluster assignments are made at the level of genes, not cells.

I propose a log-fold change statistic to address these two points. It compares the probability of gene j occurring ($w = j$) given topic k ($z = k$) versus the probability given assignment a topic other than k ($z \neq k$):

$$(12) \quad \text{lfc}^{\text{topics}}(j, k) \equiv \log_2 \frac{p(w = j | z = k)}{p(w = j | z \neq k)}.$$

For a given gene j and topic k , $\text{lfc}(j, k)$ can be calculated as

$$\begin{aligned}
\text{lfc}^{\text{topics}}(j, k) &= \log_2 \left\{ \frac{p(w = j, z = k)}{p(w = j, z \neq k)} \times \frac{p(z \neq k)}{p(z = k)} \right\} \\
&= \log_2 \left\{ \frac{\sum_{i=1}^n \sum_{t=1}^{s_i} \delta_j(w_{it}) \phi_{ijkt}}{\sum_{i=1}^n \sum_{t=1}^{s_i} \delta_j(w_{it}) (1 - \phi_{ijkt})} \right. \\
&\quad \times \left. \frac{\sum_{i=1}^n \sum_{j'=1}^m \sum_{t=1}^{s_i} \delta_{j'}(w_{it}) (1 - \phi_{ij'kt})}{\sum_{i=1}^n \sum_{j'=1}^m \sum_{t=1}^{s_i} \delta_{j'}(w_{it}) \phi_{ij'kt}} \right\},
\end{aligned} \tag{13}$$

where ϕ_{ijkt} denotes the posterior probability of $z_{it} = k$ given $w_{it} = j$,

$$\begin{aligned}
\phi_{ijkt} &\equiv p(z_{it} = k | w_{it} = j) \\
&= \frac{p(w_{it} = j | z_{it} = k) p(z_{it} = k)}{\sum_{k'=1}^K p(w_{it} = j | z_{it} = k') p(z_{it} = k')} \\
&= \frac{l_{ik} f_{jk}}{\sum_{k'=1}^K l_{ik'} f_{jk'}}.
\end{aligned} \tag{14}$$

Since the topic assignments z_{it} do not depend on t —that is, we can drop the “ t ” subscript from the ϕ_{ijkt} ’s—the expression for the lfc simplifies:

$$\text{lfc}^{\text{topics}}(j, k) = \log_2 \left\{ \frac{\sum_{i=1}^n x_{ij} \phi_{ijk}}{\sum_{i=1}^n x_{ij} (1 - \phi_{ijk})} \times \frac{\sum_{i=1}^n \sum_{j'=1}^m x_{ij'} (1 - \phi_{ij'k})}{\sum_{i=1}^n \sum_{j'=1}^m x_{ij'} \phi_{ij'k}} \right\}. \tag{15}$$

At the maximum-likelihood solution (MLE) of the l_{ik} ’s and f_{jk} ’s, the lfc statistic simplifies further:

$$\text{lfc}^{\text{topics}}(j, k) = \log_2 \left\{ \frac{\sum_{i=1}^n x_{ij} \phi_{ijk}}{\sum_{i=1}^n x_{ij} (1 - \phi_{ijk})} \times \frac{\sum_{i=1}^n s_i (1 - l_{ik})}{\sum_{i=1}^n s_i l_{ik}} \right\}. \tag{16}$$

This is because, at the MLE, the loadings l_{ik} , $k = 1, \dots, K$, for a given document/cell i should be equal to the average of the weighted counts $\frac{1}{s_i} \sum_{j=1}^m x_{ij} \phi_{ijk}$.

Finally, it is convenient that the lfc (13, 16) will be the same if we replace the multinomial topic model parameters l_{ik} and f_{jk} with the corresponding parameters of the Poisson NMF, \hat{l}_{ik} and \hat{f}_{jk} (proof not given). From the derivation of the EM algorithm for Poisson NMF, this identity holds at the MLE:

$$\hat{f}_{jk} = \frac{\sum_{i=1}^n \phi_{ijk}}{\sum_{i=1}^n \hat{l}_{ik}}.$$

Plugging this relationship into (16), we obtain the following simple expression for the log-fold change:

$$\text{lfc}^{\text{topics}}(j, k) = \log_2 \left\{ \frac{\hat{f}_{jk} \sum_{i=1}^n \hat{l}_{ik}}{\sum_{k' \neq k} \hat{f}_{jk'} \sum_{i=1}^n \hat{l}_{ik'}} \times \frac{\sum_{i=1}^n s_i (1 - \hat{l}_{ik})}{\sum_{i=1}^n s_i \hat{l}_{ik}} \right\}. \tag{17}$$

What is nice about this expression is that it can be computed without seeing the data. It is also plain to see from this expression that to arrive at a log-fold change, one must weight the factors f_{jk} by the sample-wide topic probabilities $\sum_i l_{ik}$ across i . This same expression also works with the parameters of multinomial topic model l_{ik} , f_{jk} , again, so long as they are MLEs (proof not shown).

REFERENCES

- [1] D. M. BLEI, A. Y. NG, AND M. I. JORDAN, *Latent Dirichlet allocation*, Journal of Machine Learning Research, 3 (2003), pp. 993–1022.
- [2] J. H. BULLARD, E. PURDOM, K. D. HANSEN, AND S. DUDOIT, *Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments*, BMC Bioinformatics, 11 (2010), p. 94.
- [3] X. CUI AND G. A. CHURCHILL, *Statistical tests for differential expression in cDNA microarray experiments*, Genome Biology, 4 (2003).
- [4] K. K. DEY, C. J. HSIAO, AND M. STEPHENS, *Visualizing the structure of RNA-seq expression data using grade of membership models*, PLoS Genetics, 13 (2017), p. e1006599.
- [5] C. W. LAW, Y. CHEN, W. SHI, AND G. K. SMYTH, *voom: precision weights unlock linear model analysis tools for RNA-seq read counts*, Genome Biology, 15 (2014), p. R29.
- [6] M. I. LOVE, W. HUBER, AND S. ANDERS, *Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2*, Genome Biology, 15 (2014), p. 550.
- [7] J. QUACKENBUSH, *Microarray data normalization and transformation*, Nature Genetics, 32 (2002), pp. 496–501.
- [8] M. D. ROBINSON, D. J. MCCARTHY, AND G. K. SMYTH, *edgeR: a Bioconductor package for differential expression analysis of digital gene expression data*, Bioinformatics, 26 (2010), pp. 139–140.
- [9] J. H. STOCK AND M. W. WATSON, *Introduction to Econometrics*, Pearson, Boston, MA, 3rd ed., 2015.