

Introduction

The discovery of DNA, RNA, and the development of genomic sequencing methods has proven to be a valuable tool of scientific research. However, the methods developed in the first-generation of technology relied on “bulk” sampling, which could only estimate population-average expression of RNA and DNA. In order to obtain a more complete understanding of how the cellular landscape functions, estimates of cell-to-cell variability would need to be estimated. Development of single-cell RNA sequencing (scRNAseq) technology has increased to satisfy this need, but a need for statistical analysis is still outstanding since the previous methods used to model bulk RNA sequencing data do not account for the correlated nature of scRNAseq data.

This paper will compare seven different modeling approaches on an observational scRNAseq data set obtained from a Lupus Nephritis Case/Control Study involving 33 patients across the United States. (Arazi et al. 2018) Two RNA genes were selected to be the predictor-response pair to simplify the modeling process. The main goal of this report is to investigate the ways in which parameter estimates vary as the modeling methodology is altered. It is hoped that the results of this investigation are useful for the development of future models for scRNAseq data sets.

Data

A single-cell RNA sequencing (scRNAseq) expression profile is a matrix of count-values representing a “snapshot” of the magnitude of activity of genomic features of a single cell. (“Gene Expression Profiling - Wikipedia,” n.d.) In its original form, the data being studied here had a data matrix that contained 9,560 single-cell observations clustered within 27 samples (5 control not included in data). Each observation contained the expression of 38,354 genetic features.

Single-Cell data is often unreliable, protocol dependent, and can often have batch effects. Data quality control (QC) filters out redundant measures, and dead cell observations. The Seurat guided tutorial (Satija and others 2018) was used to perform quality control, using the parameter values:

- Percent Mitochondrial DNA $> 60\%$
- Genetic Features Expressed $< 1,000$
- Genetic Features Expressed $> 5,000$
- B-cells only

These quality control measures reduced the original data by 88%, leaving only 1,110 observations clustered within 15 samples. Two genes (MALAT1 and CD19) were then selected from the set of genetic features in the initial data to be studied due to a higher correlation. MALAT1 has been linked with cancer metastasis, cell migration, and cell regulation. (“MALAT1 Gene - Genecards | Malat1 Rna Gene,” n.d.) CD19 encodes a cell surface molecule which regulates lymphocyte proliferation and differentiation. (“CD19 Gene - Genecards | Cd19 Protein | Cd19 Antibody,” n.d.).

Histograms and a joint distribution scatter plot were constructed to visualize the distributions of the selected variables (Appendix: Fig1-Fig3). The presence of zeros in the data indicated that the distribution might be well suited for a zero-inflated mixture model. Specifically, since the response was count-valued, the histograms indicated that a zero-inflated Poisson Generalized Linear Model or Generalized Linear Mixed Model would be appropriate. Additionally, while normality was not expected, log-transformations were also applied (Appendix: fig4-fig6), and resulted in approximate normality of the response MALAT1 and a bimodal distribution of the predictor CD19.

Methods

The seven different modeling methodologies explored in this investigation are:

1. Linear Models with Fixed Effects (LMwFE)
2. Linear Mixed Models with Random Effects (LMMwRE)
3. Poisson Generalized Linear Models without overdispersion (POI)
4. Poisson quasi-likelihood Generalized Linear Models with over dispersion (POIql)
5. Poisson Generalized Linear Mixed Models fit using Penalized Quasi-Likelihood (POIImm)
6. Zero-Inflated Poisson Generalized Linear Mixed Model – Fixed Effect Subject (ZIPwFE)
7. Zero-Inflated Poisson Generalized Linear Mixed Model – Random Effect Subject (ZIPwRE)

Models (1) and (2) are fit on log-transformed data (predictor and response), the remaining models are fit on untransformed data.

We assume that repeated measure residual errors are independent (when applicable), i.e for subject $i = 1, \dots, 15$ and repeated measure $j = 1, \dots, n_i$

$$\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$$

In an attempt to gain insight into the processes governing cellular relationships with their host subjects, a subject parameter was also fit using either a fixed or a random effect.

$$\begin{aligned} \text{Fixed Effect : } & \beta_{0i} \\ \text{Random Effect : } & b_{0i} \end{aligned}$$

where:

$$b_{0i} \sim N(0, \sigma_b^2)$$

A fixed effect, global intercept and a covariate (CD19) parameter are included in all models, these parameter estimates (along with their standard errors) will be the primary focus for comparing the modeling methods.

We let:

$$\gamma_{ij} = \beta_0 + \beta_1 \text{ CD19}$$

and unless otherwise stated, in situations where a model calls for a link function, the canonical link function for count data will be used:

$$\mu_{ij} = g^{-1}(\eta_{ij}) = \log(\eta_{ij})$$

Zero-Inflated models will use a fixed effect intensity model with only an intercept

$$R_{ij} \sim \text{bernoulli}(p_{ij}|\alpha_{0i}) \quad p(R_{ij} = 1) = \alpha_{0i}$$

the intensity model parameters will be varied using a fixed or random effect for subject. This terminology allows the models to be written as:

Model #	Model Name	Model Format	Assumed Data Distribution
1	LMwFE	$Y_{ij} = \beta_{0i} + \gamma_{ij} + \epsilon_{ij}$	$Y_{ij} \sim N(\beta_{0i} + E[\gamma_{ij}], \sigma_\epsilon^2)$
2	LMMwRE	$Y_{ij} = \gamma_{ij} + b_{0i} + \epsilon_{ij}$	$E[Y_{ij} b_{0i}] \sim N(E[\gamma_{ij} + b_{0i}], \sigma_\epsilon^2)$
3	POI	$\mu_{ij} = g^{-1}(\eta_{ij} = \beta_{0i} + \gamma_{ij})$	$Y_{ij} \sim \text{Poisson}(\beta_{0i} + E[\gamma_{ij}])$
4	POIql	$\mu_{ij} = g^{-1}(\eta_{ij} = \beta_{0i} + \gamma_{ij})$	$Y_{ij} \sim \text{Poisson}(\phi(\beta_{0i} + E[\gamma_{ij}]))$
5	POIImm	$\mu_{ij} = g^{-1}(\eta_{ij} = \gamma_{ij} + b_{0i})$	$E[Y_{ij} b_{0i}] \sim \text{Poisson}(E[\gamma_{ij}] + b_{0i})$
6	ZIPwFE	$R_{ij} \sim \text{bernoulli}(p_{ij} \alpha_{ij})$ $\mu_{ij} (r_{ij} = 1) = g(\eta_{ij} = \beta_{0i} + \gamma_{ij})$	$Y_{ij} \sim \text{ZerInfPoi}(\beta_{0i} + E[\gamma_{ij}], p_{ij})$
7	ZIPwRE	$R_{ij} \sim \text{bernoulli}(p_{ij} \alpha_{ij})$ $\mu_{ij} (r_{ij} = 1) = g(\eta_{ij} = \gamma_{ij} + b_{0i})$	$E[Y_{ij} b_{0i}] \sim \text{ZerInfPoi}(E[\gamma_{ij}] + b_{0i}, p_{ij})$

Results

Intercept Estimates

Model	Estimate	Std. Error	pvalue
LMwFE	8.3464	$4.981 * 10^{-2}$	$< 2 * 10^{-16}$
LMMwRE	8.3479	$1.3565 * 10^{-1}$	$< 2 * 10^{-16}$
POI	8.821	$4.856 * 10^{-4}$	$< 2 * 10^{-16}$
POIql	8.957	$3.007 * 10^{-2}$	$< 2 * 10^{-16}$
POIImm	8.8362	$1.0163 * 10^{-1}$	$< 1 * 10^{-5}$
ZIPfe	8.958	$3.689 * 10^{-4}$	$< 2 * 10^{-16}$
ZIPre	8.9402	$6.5229 * 10^{-4}$	$< 1 * 10^{-4}$

Slope Estimates

Model	Estimate	Std. Error	pvalue
LMwFE	$5.590 * 10^{-2}$	$1.534 * 10^{-2}$	$2.82 * 10^{-4}$
LMMwRE	$5.703 * 10^{-2}$	$1.528 * 10^{-1}$	$1.8935 * 10^{-4}$
POI	$3.246 * 10^{-4}$	$2.513 * 10^{-6}$	$< 2 * 10^{-16}$
POIql	$8.839 * 10^{-5}$	$1.913 * 10^{-4}$	0.644
POIImm	$3.16 * 10^{-4}$	$1.6525 * 10^{-4}$	$5.61 * 10^{-2}$
ZIPfe	$8.559 * 10^{-5}$	$2.176 * 10^{-6}$	$< 2 * 10^{-16}$
ZIPre	$2.9282 * 10^{-4}$	$2.1289 * 10^{-6}$	$1 * 10^{-4}$

Intercept Estimate Percent Change

	1	2	3	4	5	6	7
1	0	0.00	-0.06	-0.06	-0.06	-0.06	-0.07
2	0.00	0	-0.06	-0.06	-0.06	-0.06	-0.07
3	0.05	0.05	0	0.00	-0.00	0.00	-0.01
4	0.05	0.05	0.00	0	-0.00	0.00	-0.01
5	0.06	0.06	0.00	0.00	0	0.00	-0.01
6	0.05	0.05	0.00	0.00	-0.00	0	-0.01
7	0.07	0.07	0.01	0.01	0.01	0.01	0

Slope Estimate Percent Change

	1	2	3	4	5	6	7
1	0	-0.02	0.99	0.99	0.99	0.99	1.00
2	0.020	0	0.99	0.99	0.99	0.99	1.00
3	-171.23	-174.72	0	0.00	0.03	0.01	0.10
4	-171.23	-174.72	0.00	0	0.03	0.01	0.10
5	-175.85	-179.44	-0.03	-0.03	0	-0.02	0.07
6	-172.59	-176.10	-0.01	-0.01	0.02	0	0.09
7	-189.89	-193.76	-0.11	-0.11	-0.08	-0.10	0

Where the numerical model-mapping $model\ name \mapsto \{1, 2, \dots, 7\}$ is given in the model definition table.

Discussion

The results have shown that changes in modeling strategy have little impact on the estimate of the intercept parameter. This is supported by the fact that the maximum absolute difference between estimates for the intercept is a 7% change. This consistency also supports a stronger conclusion that parameter estimates for intercept agree in both sign and magnitude for all models.

These statements are not true for covariate parameter estimates. While these parameters are comparable within similar model frameworks, the estimates agree on only sign across all models. This is not an unexpected result, since parameter estimates for slopes are heavily dependent on how observations are correlated, and this concept is approached differently in all of the models employed.

It should be noted that, while the estimates for the slope parameters have been estimated to be “significant”, the magnitude of these parameter estimates deviates very little from 0. Given this information, it would not be unreasonable to conclude that the effect of CD19 can be almost completely ignored. Significant nested model comparisons were performed to test for the significance of the CD19 covariate. It was found that there was sufficient evidence to include the covariate when compared to the null model in most situations involving log-transformed variables, and when the subject term had already been included in as a Fixed Effect.

The standard errors of both parameter estimates (intercept and slope) indicate a trend of increased uncertainty in parameter estimates as methodologies move from incorporating subject as a fixed effect, to incorporating it as a random effect. While this result was not initially expected, it is plausible in light of the nested model comparison results, and would indicate that the subject effect is likely correlated with CD19. Which would imply that the random effect assumption: “unobserved... [individual/subject] heterogeneity is uncorrelated with independent variables” (Wooldridge 2010) is likely violated.

Limitations and Future Research

Initial quality control measures for the scRNAseq data found an extremely high presence of mitochondrial RNA (mRNA). Mitochondrial functionality is essential to cellular processes, upon death such functions cease, and mitochondria degrade. As a result, mRNA content is used as a quality control measure to indicate progression of cellular death. The Seurat tutorial for single-cell analysis recommended filter parameters of %mRNA < %5; however this threshold eliminated all data when used in conjunction with the other recommended quality control measures. A final determination of %mRNA < %60 was used to preserve as much of the clustering-structure as possible, with the trade-off of analyzing dead cells. Even after altering the QC thresholds to allow for more data, observational imbalance between subjects was observed, and remained mostly unaccounted-for throughout the analysis.

The results of this investigation are based upon mostly non-living biological samples. While interpretation of model parameters is certainly possible, the individual estimates are less meaningful considering the context than the estimate trends that we are seeking. The sign of the estimates for the slope parameter, which would be the most contextually illuminating (as there were no control-subjects in the data), demonstrate that increased values of CD19 are marginally (if at all) associated with higher values of MALAT1. However, as previously indicated, the magnitude of these estimates reinforce interpretational futility.

The analysis conducted in this investigation is inadequate and inaccurate. The analysis conducted here is inaccurate because of biologically poor quality data, and imbalance between subject information. Additionally, the model estimates being compared can not be used for inference on the source of estimate variances because of multi-factor changes across each estimate (model fitting technique, linear predictor changes, data transformation differences). The results of this investigation are inadequate because the original goal of the investigation was to investigate models over single-cell data and isolate the specific effects that the modeling strategies have over parameter estimates. Since only vague conclusions regarding slope and intercept were obtained, this investigation has fallen short of its objective.

Consequently, the future of research for this project is very broad in scope. Immediate considerations will be made to compare more simplistic modeling strategies over more similar models.

References

- Arazi, Arnon, Deepak A Rao, Celine C Berthier, Anne Davidson, Yanyan Liu, Paul J Hoover, Adam Chicoine, et al. 2018. "The Immune Cell Landscape in Kidneys of Lupus Nephritis Patients." *bioRxiv*. Cold Spring Harbor Laboratory, 363051.
- "CD19 Gene - Genecards | Cd19 Protein | Cd19 Antibody." n.d. <https://www.genecards.org/cgi-bin/carddisp.pl?gene=CD19&keywords=CD19>.
- "Gene Expression Profiling - Wikipedia." n.d. https://en.wikipedia.org/wiki/Gene_expression_profiling.
- "MALAT1 Gene - Genecards | Malat1 Rna Gene." n.d. <https://www.genecards.org/cgi-bin/carddisp.pl?gene=MALAT1>.
- Satija, R, and others. 2018. "Seurat: Guided Clustering Tutorial." *Satija Lab* [Http://Satijalab.Org/Seurat/Pbmc3k_tutorial.Html](http://Satijalab.Org/Seurat/Pbmc3k_tutorial.Html).
- Wooldridge, Jeffrey M. 2010. *Econometric Analysis of Cross Section and Panel Data*. MIT press.