Figure 5:

We see that the log-transformed response MALAT1 is approximately normal distribution. 168
Conversely, the log-transformed response FBLN1 is not inherently better than the un- 169
transformed response. We can clearly see the heavy influence of zero-inflation in these 170
variables as is apparent from the dominance of the "zero-bins" in Figure (5). 171

Regardless, we model each outcome under the assumption that: compensating for observa- 172
tional correlation will sufficiently account for non-normality of the responses. This may not 173
generally the case, and additional transformations or modeling methodologies may be needed 174
to improve model error distributions. However, for the purpose of comparing the previously 175
mentioned models on subject-correlated single-cell data, we will proceed with this assumption 176
and verify ridual homoscedasticity, normality and independence using fitted vs residual plots 177
and quantile-quantile plots. 178

11

# Model Descriptions

*[handwritten margin note top-left: Double?]*

*[handwritten note: What is meant by transformed? and is this necessary? don't we just need a variable?]*

We define our outcome(s) of interest to be one of the following transformed variables as 180

taken from Arazi, Rao, Berthier, et al. Let a single observation be designated as: $R_{hij}$. The 181

*[handwritten note: Lee, please describe the models in more general terms. No need to have the index h]*

index $h = \{1, 2\}$ represents the model pairing number ($CD19 \sim MALAT1$ is pairing #1 182

*[handwritten note: why not start at 1?]*

$CD34 \sim FBLN1$ is pairing #2) $i \in \{5, 6, \ldots, 26\}$ represents the subject (name of subject by 183

number) from which the observation originated, and the index $j = 1, \ldots, n_i$ represents the 184

single-cell observation within subject-i. We note that $n_i \in \{21, 22, 23, \ldots, 127\}$ in the context 185

of the Lupus Data. 186

*[handwritten note: please be more general, why?]*

*[handwritten note: as the # of single-cell obs. in each subject]*

We perform the transformations: 187

$$R_h = \log\left(Y_h^* + 1\right)$$

where 188

$$Y_1^* = \text{MALAT1} - 67 \quad \text{and} \quad Y_2^* = \text{FBLN1}$$

and 189

$$Y_h = Y_h^* - min(Y_h^*)$$

We aslo define the predictor attached to $R_h$ as: 190

$$P_h = \log\left(X_h + 1\right) \quad \text{for} \quad h = 1, 2$$

where 191

$$X_1 = \text{CD19} \quad \text{and} \quad Y_2 = \text{CD34}$$

We present the theoretical model frameworks here as "Less Than Full Rank" (LTFR) repre- 192

sentations. The Full-Rank model results presented in the *Results* section to follow are created 193

12

by droping the first level in all factors and using this as the referrence level. 194

## Linear Regression 195

We begin the model framework definitions by describing two Linear Regression models, with 196
Fixed Effect parameters estimated using maximum likelihood optimization. It should be 197
noted that these methods make the assumption that observations are independent, and 198
should therefore be used for comparison to modeling methods to come. However, the linear 199
regression models we present here can account for some observational correlation with the 200
use of a subject specific intercept term as we will see in the second model. 201

Ultimately, all the methods defined in this section assume an identical error structure across 202
all observations of the form: 203

$$\epsilon_{hij} \sim N\left(0, \sigma_\epsilon^2 * I_{1110}\right)$$

where we are assuming that $\sigma^2$ is a common variance parameter for all subjects and $I_{1110}$ is 204
the 1110 X 1110 identity matrix. 205

### Simple Linear Regression (Model 0) 206

Using the notation we defined above, we write the first model as: 207

$$R_{hij} = \beta_0 + \beta_1 P_{hij} + \epsilon_{hij}$$

which is equivalent to: 208

$$\log(Y_{hij}) = \beta_0 + \beta_1 \log(X_{hij}) + \epsilon_{hij}$$

We note that this model does not account for observational correlation, and instead provides 209
an estimation for population-averaged relationships, namely: 210

13

- What is the estimated average (across all observations, across all subjects) value of $R_{hij}$ when $P_{hij} = 0$ (intercept)

- On average (across all observations, across all subjects) what is the average rate of change in $R_{hij}$ per unit increase in $P_{hij}$ (slope)

## Fixed-Effect Subject-Specific Intercept (Model 1)

Adding a subject-specific intercept term allows us to account for within-subject correlation by uniformly shifting the fitted values specific to a subject. This model may be written as:

*the mean of the fitted values!*

$$R_{hij} = \beta_0 + \beta_{1i}(subject_i) + \beta_2 P_{hij} + \epsilon_{hij}$$

where we define the term:

*what do you use to define the total # of subjects?*

$$\beta_{1i}(subject_i) = \begin{cases} \beta_{1i} & \text{if } subject_i = i \\ 0 & \text{if } subject_i \neq i \end{cases}$$

*adds (N-1)*

This model provide the added estimated parameter $\hat{\beta}_{1i}$ which tells us a uniform estimated *is the* average deviation for each subjects fitted-response from the global estimated mean provided by Model 0 (Simple Linear Regression).

## Linear Mixed Effects Models

The next category of modeling approaches we describe is Linear Mixed Effect Models with Random Effects. Specifically, we describe two distinct Linear Mixed Effect Models that account for subject-correlation in a different manner than the previously discussed Linear Regression models. Linear Mixed Efffects Models do not necessarily assume-independence *require the assumption of* of observations. Correlation structures such as AR(1), spatial power, or unstructured can be

*change all m(e) to lower case*

*write out*

14

*we don't need the responses to be normally distributed? Just the random error, ε*

used ~~to estimate parameters~~ determining correlation amongst observations ~~within a subject~~ 228

and between observations ~~(across subjects)~~. Additionally, if we can ~~rationally~~ assume that the 229

**responses** shown in Figure 3 have a multivariate normal distribution, the model parameters 230

can be easily estimated using Maximum Likelihood Estimation techniques [9]. 231

*such as REML?*

## Linear Mixed Effects Model with Random Intercept (Model 2) 232

Model 1 (Linear Regression with Fixed Effect Intercept) accounts for subject correlation by 233

assuming that observations within a subject are uniformly ~~influenced~~ *correlated* by the nested nature 234

of the sampling ~~method~~ (i.e. observations are sampled so that they are identically correlate 235

within each subject). However, this assumption may not always be reasonable, as we could 236

imagine that responses within each subject also exhibit random variation that is related to 237

nested sampling methods. 238

*awkward & repetitive*

A Linear Mixed Effects Model that includes a Random Intercept accounts for subject-level 239

observational correlation by inducing individual-specific levels of random variation into all 240

observations specific to each subject. Such a model may be written as: 241

*But the random var is random so doesn't need to be modeled*

$$R_{hij} = \beta_0 + \beta_1 P_{hij} + b_{0i}\left(subject_i\right) + \epsilon_{hij}$$

where 242

$$b_{0i} \sim N\left(0, \sigma_b^2\right) \quad \text{for} \quad i \in \{5, 6, \ldots, 26\}$$

*make more general*

$$\epsilon_{hij} \sim N\left(0, \sigma_\epsilon^2 I_{n_i}\right)$$ 243

*does this need to be indexed by $n_i$?*

and we assume that $b_{0i}$ and $\epsilon_{hij}$ are independent. 244

We note that both random-components ~~can be~~ *are* assumed to have a mean of zero ~~as non-zero~~ 245

~~components are inherently deterministic and can be integrated into intercept terms.~~ 246

15

# Linear Mixed Effect Model with Random Slope (Model 3)

A further accounting for the effects of subject-level observational clustering may be made by
extrapolating on Model 2 (Linear Mixed Effects Model with Random Intercept) with the
addition of a random intercept.

The incoperation of a Fixed Effect subject-specific slope would account for subject-level
observational correlation by assuming that the relationship between predictor and response
are uniformly influenced across observations. Implying that, in addition to the average
response devation from the estimated average response, there is also an average uniform shift
in how each subjects' response changes with respect to a unit shift in the predictor. Again,
this assumption may not be reasonable, as we may expect variation in how responses within
a subject deviate from the estimated average change in response over the predictor space.

We will therefore incorperated a Random Slope into the format of the Random Intercept
model (Model 2) to attempt to reconcile these effects. This will allow for us to account for
observational correlation due to subject-level sampling as sourced from:

- subject-specific random variation associated with measurement instability
- predictor-dependent, subject-specific random variation associated with measurement instability

We write this model as:

$$R_{hij} = \beta_0 + \beta_1 P_{hij} + b_{0i} \left(subject_i\right) + \left[b_{1i}\left(subject_i\right) P_{hij}\right] + \epsilon_{hij}$$

where

$$\mathbf{b} = \begin{bmatrix} b_{0i} \\ b_{1i} \end{bmatrix} \sim N\left(\mathbf{0}, \mathbf{G}\right)$$

16

$$G = \begin{bmatrix} \sigma_b^2 & 0 \\ 0 & \sigma_b^2 \end{bmatrix}$$

*are these really the same?* (handwritten)

$$\epsilon_{hij} \sim N\left(0, \sigma_\epsilon^2 \mathbf{I}_{n_i}\right)$$

# Generalized Estimating Equations (Model 4)

~~Our final method for modeling scRNA-seq expression profiles is Generalized Estimating~~ 267

~~Equations (GEE). Dissimilar to each of the methods previously described,~~ GEE regression 268

esitimates are obtained using methodologies that allow for non-continuous responses. GEE 269

also extrapolates on the techniques used for modleing non-normal responses by incorperating 270

the effects of observational correlation. 271

*so do the others if we treat the use generalized* (handwritten)

*I would remove this section as LME & OLS regression can be easily extrapolated to generalized models* (handwritten)

GEE estimates are computed by solving the estimating equation(s): 272

*? Generalize please* (handwritten)
*define all terms please* (handwritten)

$$0 = U(\beta) = \sum_{i=1}^{15} \left\{ \mathbf{D}_{hi}^T \mathbf{V}_{hi}^{-1} \left( \mathbf{y}_{hi} - \mu_{hi} \right) \right\} \tag{1}$$

where: 273

*Instead, main difference is marginal (GEE) vs conditional (LME) controlling of the groups.* (handwritten)

$$\mu_{hi} = \mu_{hi}(\beta) = E\left[\mathbf{Y}_{hi}\right] = \eta_{hi}$$

represents the relationship between the expected value of the response $\mu_i$ (not necessarily 274

assumed to be a distribution) and the linear predictor $\eta_i$, 275

*☆ DISCUSS HOW TO INTERPRET LME (conditional) & GEE (marginal)* (handwritten)

$$\mathbf{D}_{hi} = \begin{bmatrix} \frac{\partial \mu_{hi1}}{\beta_1} & \frac{\partial \mu_{hi1}}{\beta_2} & \cdots & \frac{\partial \mu_{hi1}}{\beta_p} \\ \frac{\partial \mu_{hi2}}{\beta_1} & \frac{\partial \mu_{hi2}}{\beta_2} & \cdots & \frac{\partial \mu_{hi2}}{\beta_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \mu_{hin_i}}{\beta_1} & \frac{\partial \mu_{hin_i}}{\beta_2} & \cdots & \frac{\partial \mu_{hin_i}}{\beta_p} \end{bmatrix}$$
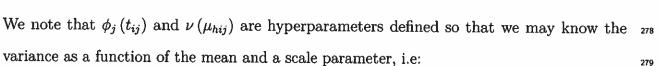
is the first derivative matrix, and

$$\mathbf{V}_{hi} = \mathbf{A}_{hi}^{\frac{1}{2}} Corr(\mathbf{Y}_{hi}) \mathbf{A}_{hi}^{\frac{1}{2}}$$

$$\mathbf{A}_{hi} = \underset{n_i}{diag} \left\{ \phi_j \left( t_{ij} \right) \nu \left( \mu_{hij} \right) \right\}$$

We note that $\phi_j \left( t_{ij} \right)$ and $\nu \left( \mu_{hij} \right)$ are hyperparameters defined so that we may know the variance as a function of the mean and a scale parameter, i.e:

$$Var \left( Y_{hij} \right) = \phi_j \left( t_{ij} \right) \nu \left( \mu_{hij} \right)$$

The GEE algorithm is iterative and used the following steps to converge at an estimate:

1. Generalized Linear Modeling methods employing Maximum Likelihood Estimation are used to obtain intial estimates for $\beta$

2. Estimates for $\beta$ used to compute hyper-parameters

3. New estimates for hyper-parameters and working covariance matrix $(\mathbf{V}_{hi})$ used to obtain new estimates for $\beta$ by solving (1)

4. Repeat Steps 2 & 3 until algorithm converges

The GEE algorithm has a quality which makes it very appealing for many applications with observational clustering. Specifically, the algorithm is robust to misspecification of the observational correlation structure. That is, the estimates $\hat{\beta}_{GEE}$ are consistent with $\beta$ irrespective of the estimates for within-subject correlation.

18

which implies we will be assuming the general modeling structure: 308

$$E\left[Y_{hij}\right] = \mu_{hij} = \eta_{hij} = \beta_0 + \beta_1 P_{hij}$$

we will assume a variance function of the form: 309

$$Var\left(Y_{hij}\right) = \phi$$

and we will be using a working covariance matrix structure for repeated measures that 310 corresponds to the assumption of independence of observations within a subject. 311

$$[Corr\left(Y_{hij}, Y_{hik}\right)]_{jk} = \begin{cases} 1 & \text{if} \quad j = k \\ 0 & \text{if} \quad j \neq k \end{cases}$$

$$for \quad j, k \in \{1, \ldots, n_i\}$$
312

# Results 313

Table 8 and table 9 display parameter value estimates, standard errors, test statistics, and 314 p-values for the main-effect slope term estimated by all five modeling approaches: 315

20

The GEE algorithm is also very stable, in-part due to the fact that the effect(s) that it [291] estimates are population-averaged. Each of the previous methods (Model 0 withstanding) had [292] subject-specific interpretations, but the GEE algorithm provides marginal parameter estimates. [293] These values do not represent any specific subject, but rather the population-average. [294] According to Fitzmaurice, Laird, and Ware [9] we also need to ensure that any responses [295] modeled in the GEE process are stationary, i.e: [296]

$$E\left[Y_{hij}|\mathbf{X}_{hi}\right] = E\left[Y_{hij}|X_{hi1}, \ldots, X_{hin_i}\right] = E\left[Y_{hij}|X_{hij}\right]$$

The scRNA-seq data has been assumed to be independent within-subject, therefore we have: [297]

$$E\left[Y_{hij}|X_{hij}\right] = E\left[Y_{hij}|X_{hij'}\right]$$

[298]

$$\forall j \in \{1, \ldots, n_i\} \quad j \neq j'$$

as needed. [299]

The three-part specification of the GEE framework includes: [300]

1. The link function and linear predictor [301]

2. Variance function [302]

3. A working covariance matrix [303]

The link function and linear predictor are chosen so that the resulting model estimates will [304] be comparable to preceeding estimates for intercept and slope. Therefore, we will use the [305] identity link function: [306]

$$g(x) = x$$

in conjunction with the linear predictor: [307]

$$g(\mu_{hij}) = \eta_{hij} = \beta_0 + \beta_1 P_{hij}$$

19