

Normalization and **variance** **stabilization** of single-cell RNA-seq data using **regularized** negative binomial regression

Christoph Hafemeister et al. (2020/04)

1-introduction

Characteristics of an effective normalization

- In general, the normalized expression level of a gene should not be correlated with the total sequencing depth of a cell.
- The variance of a normalized gene (across cells) should primarily reflect biological heterogeneity, independent of gene abundance or sequencing depth.

SC RNA-seq Counts are confounded by seq depth

- 10X 22k PBMC is used
- In B, divide the genes into 6 groups by gene mean cross cell.
- C, D show the trends of gene count vs. cell count.
- E divides cells into 5 groups based seq-depth.
- E calculate the contribution of each cell group to the variations within each gene group.

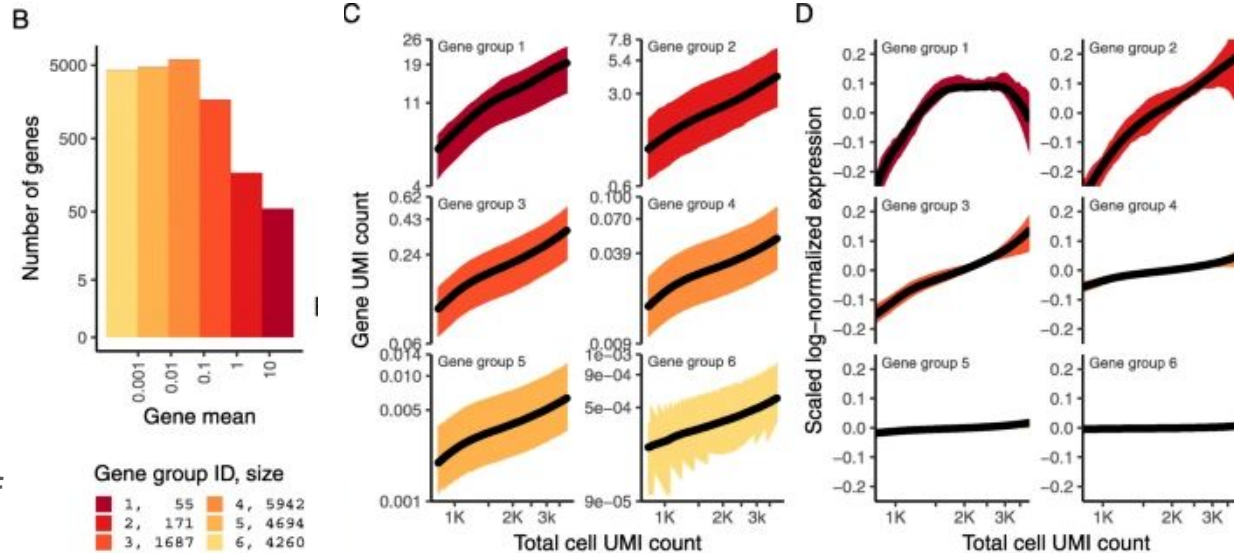
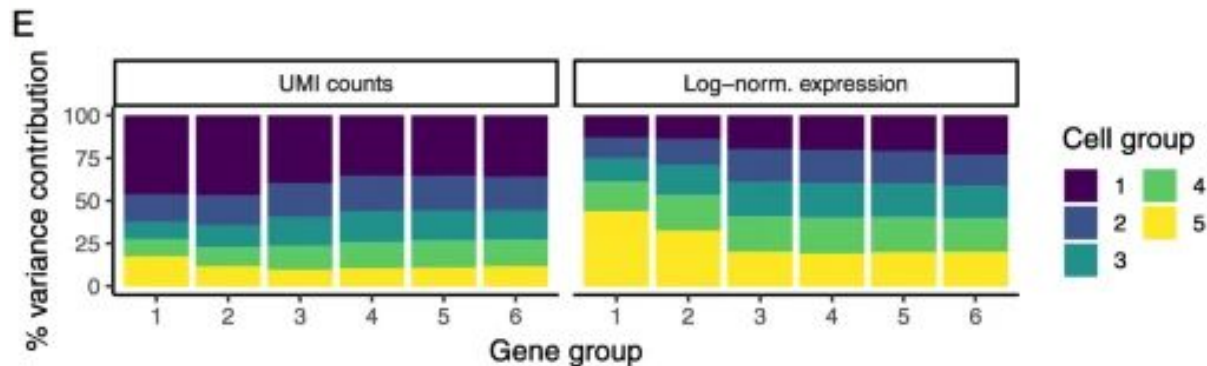


Fig1

SC RNA-seq variance are confounded by seq depth

- E divides cells into 5 groups based seq-depth.
- E calculate the contribution of each cell group to the variations within each gene group.



D and E shows that being normalized by single size factor will be confounded by sequence depth.

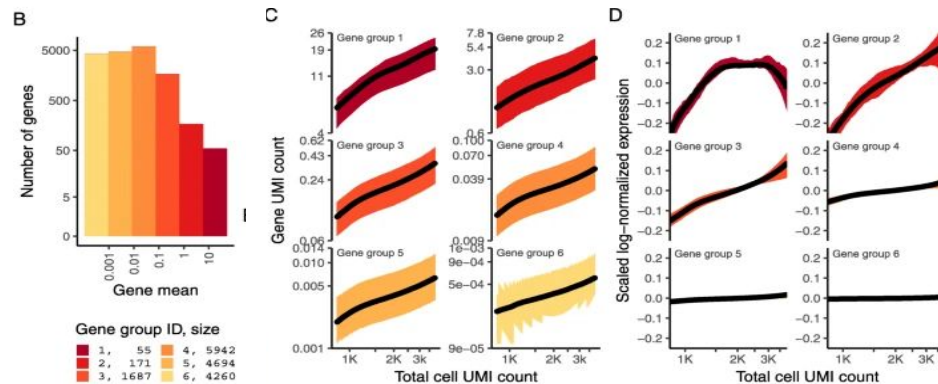


Fig1

Main idea of SCTransform

- That different groups of genes cannot be normalized by the same constant factor.
- A generalized linear model (GLM) for each gene with UMI counts as the response and sequencing depth as the explanatory variable.
- We can regularize parameter by pooling information across genes with similar abundances.

2-method

Review General Linear Model

Linear Model

$$Y|X \sim \mathcal{N}(\mu(X), \sigma^2 I)$$

$$\mathbb{E}(Y|X) = \mu(X) = X^\top \beta$$

General Linear Model

1. Random component:

$Y \sim$ some exponential family distribution

2. Link: between the random and covariates:

$$g(\mu(X)) = X^\top \beta$$

where g called link function and $\mu = \mathbb{E}(Y|X)$.

Review General Linear Model

Linear Model

$$Y|X \sim \mathcal{N}(\mu(X), \sigma^2 I)$$

$$\mathbb{E}(Y|X) = \mu(X) = X^\top \beta$$

MASS: glm.nb()

General Linear Model

1. Random component:

$Y \sim$ Negative Binomial

2. Link: between the random and covariates:

$$\log(\mathbb{E}(x_i)) = \beta_0 + \beta_1 \log_{10} m_i$$

Variance:

$$\mu + \frac{\mu^2}{\theta}$$

Regularize the GLM

Fig A is the dot plot of model parameters vs. gene mean.

Interestingly, they use bootstrap to calculate the uncertainty of each dot and color the dots by degree of uncertainty.

Fig B is the distribution of parameters' standard deviation during bootstrap.

So the estimated parameters were not reproducible across bootstraps.

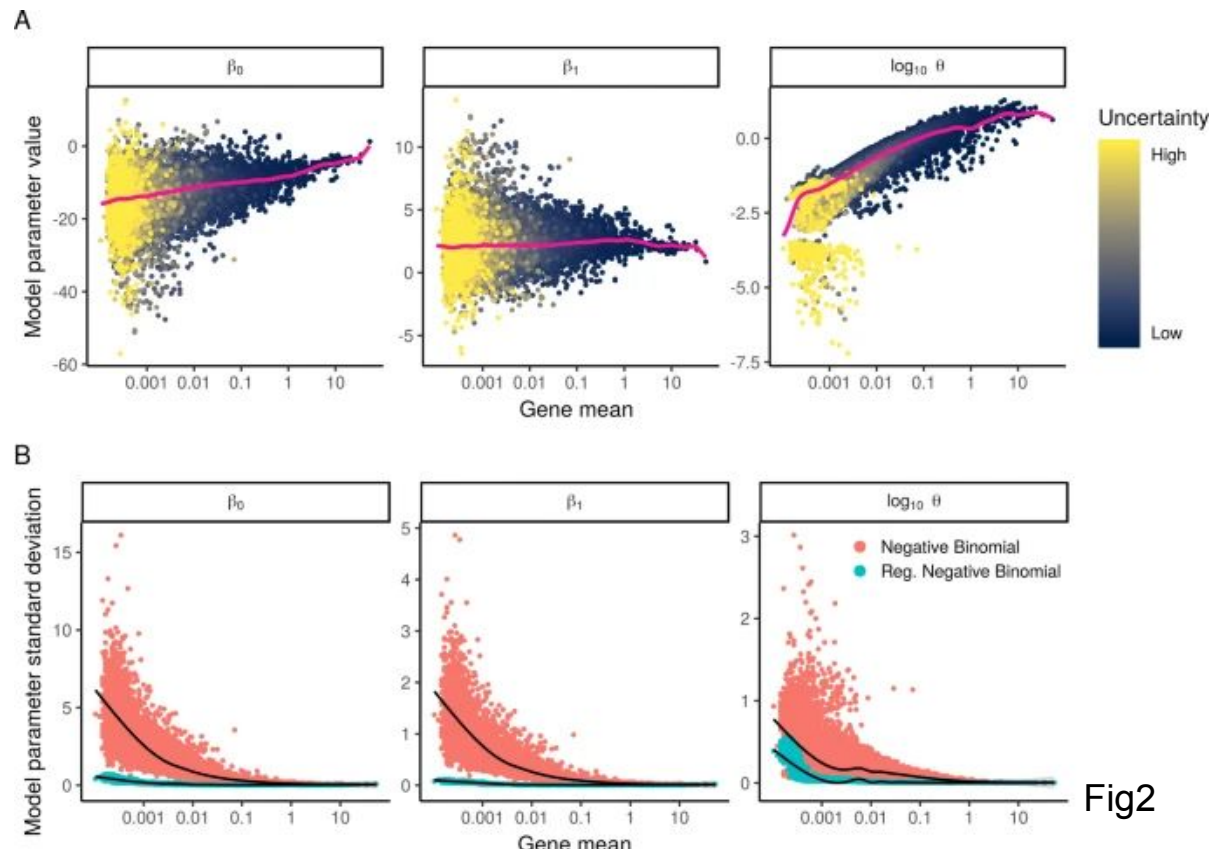


Fig2

Regularize the GLM

So the estimated parameters were not reproducible across bootstraps.

kernel regression is used to smoothing the distribution by considering the effect of neighbors within the kernel window.

It can be seen as a regularization to avoid overfitting to the technical noise.

The blue dots in fig B show that the parameters become stable during bootstrap after regularization.

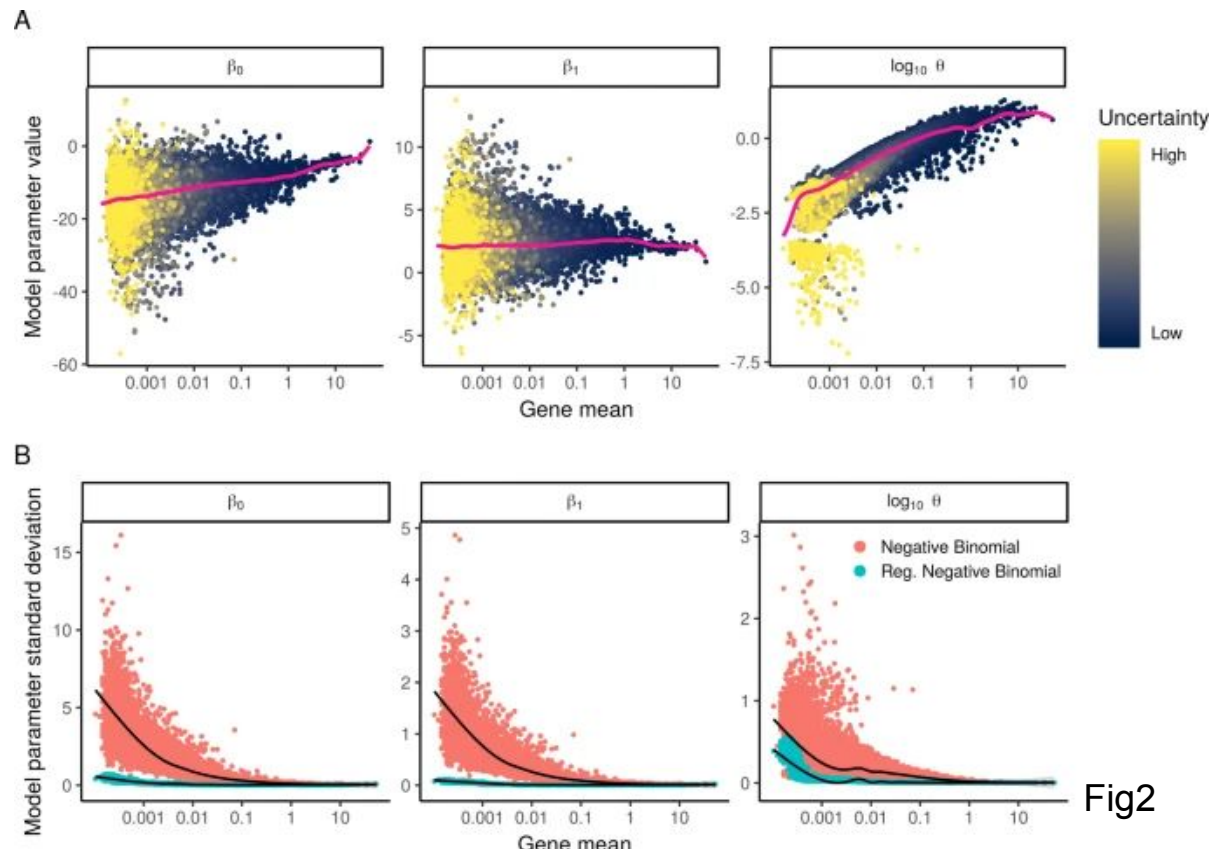


Fig2

Pearson residual effectively normalized scRNA-seq

Pearson residuals: response residuals divided by the expected standard deviation.

In fig A, it shows that Pearson residuals are independent with the total UMI counts in most of the cells.

In figure B, it shows that the contribution of variance become even if Pearson residuals are used as the corrected expression level.

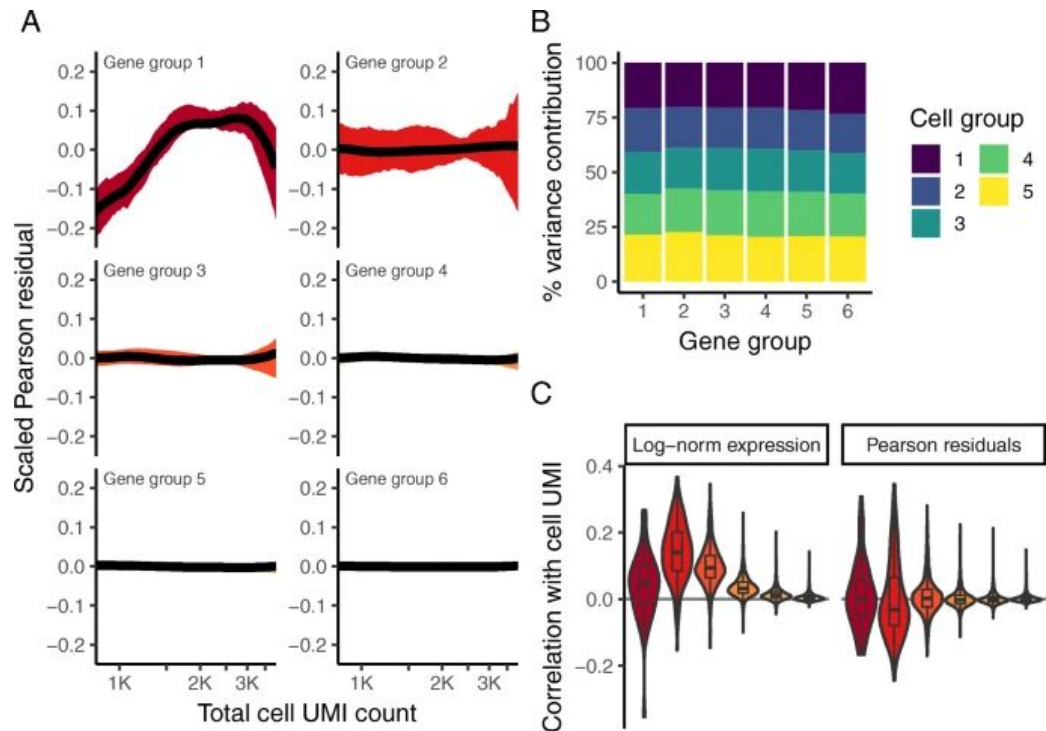


Fig3

3-result

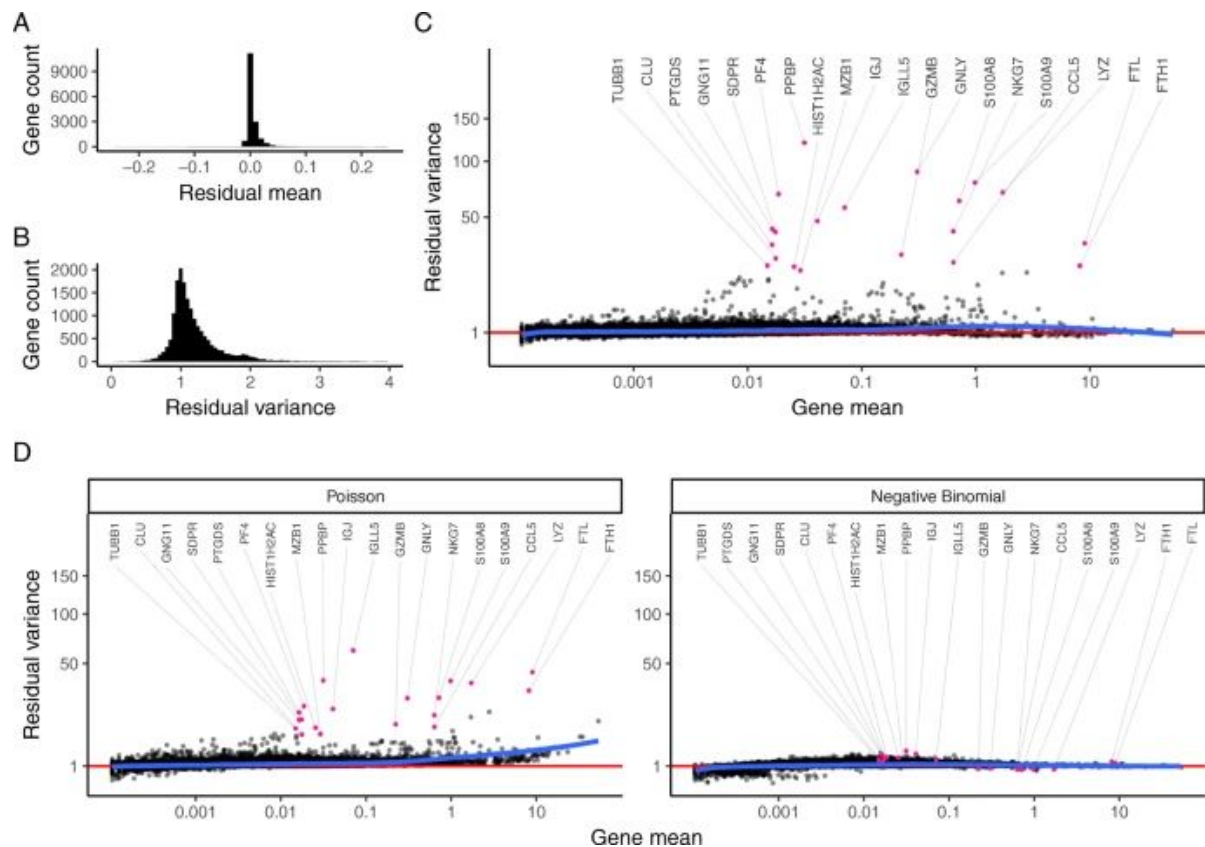
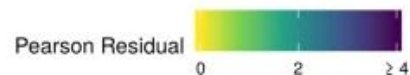
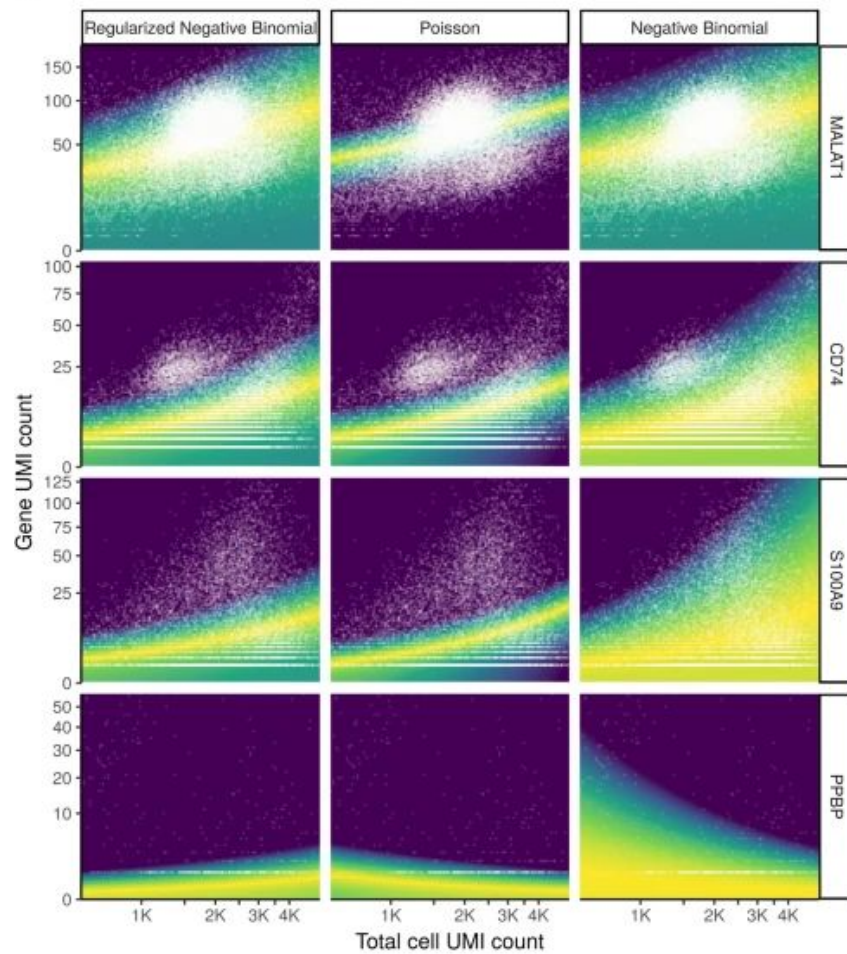


Fig4

A



B

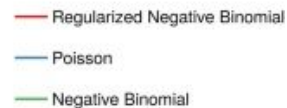
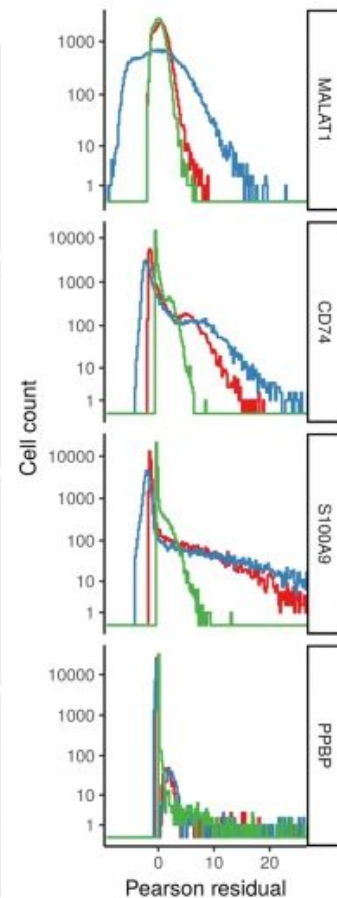
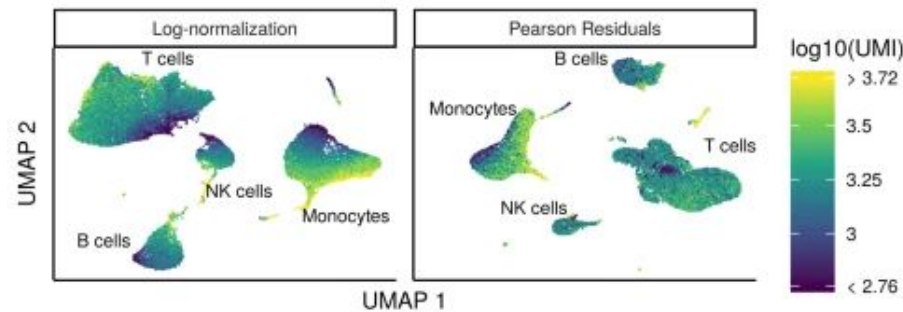
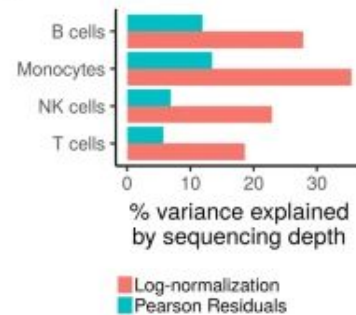


Fig5

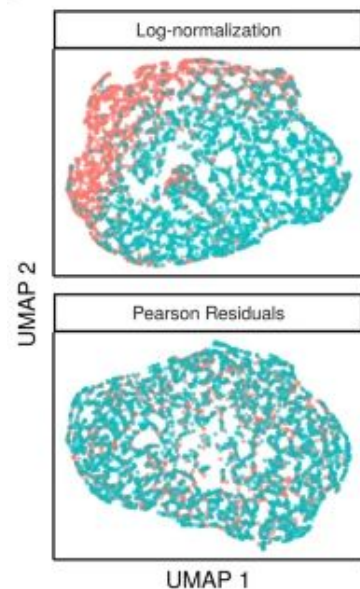
A



B



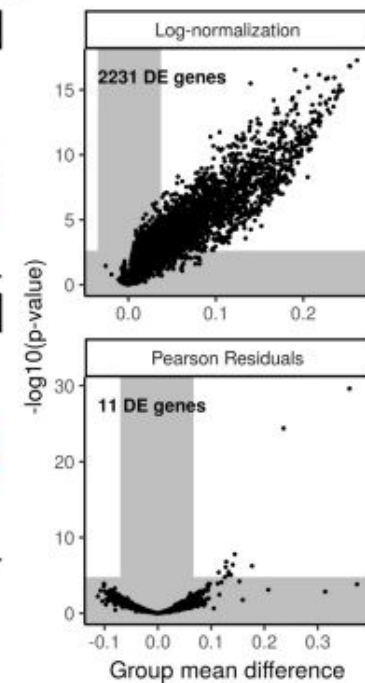
C



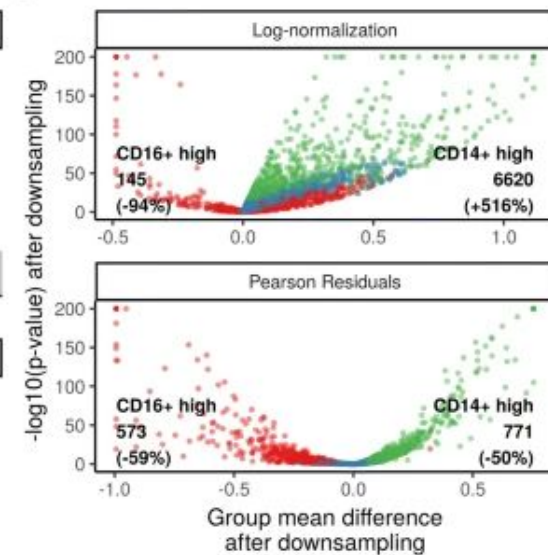
Original

Downsampled

D



E



DE classification of genes in monocytes before downsampling CD16+ group

High in CD16+

not DE

High in CD14+

Thank you!