

Deeper differential expression analysis with shrinkage correction

Jared Brown

Postdoctoral research fellow

Lab of Rafael Irizarry · DFCI Data Science

16 July 2025

Follow-along preliminaries

Software:

- R: <https://cran.r-project.org/>
- Rstudio: <https://posit.co/download/rstudio-desktop/>
- R packages:

```
install.packages("BiocManager")
```

```
BiocManager::install(c("DESeq2", "pasilla",...))
```

```
DESeq2, pasilla, ggplot2, matrixStats, apeglm, ashR, TENxPBMCDData, Matrix,  
    irlba, scan, scater, sparseMatrixStats
```

Analysis scripts:

- https://github.com/JBrownBiostat/DifferentialExpressionTraining_July2025

From the original developers, Michael Love, Simon Anders, Wolfgang Huber:
<https://bioconductor.org/packages/release/bioc/vignettes/DESeq2/inst/doc/DESeq2.html>

General (single-cell RNA) pipelines from some greats in the field:
<https://bioconductor.org/books/3.17/OSCA/>

Carpentry workshop from Bioconductor:
<https://carpentries-incubator.github.io/bioc-rnaseq/>

DE definitions: data structure

A $G \times J$ matrix of abundance measures across:

- G many features: **genes**, transcripts, exons, protein binding peaks, methylation sites, etc.
- J many samples: **experimental libraries**, single cells, spatial spots, binned pixels, etc.
- First few lines from the *pasilla* dataset of pasilla gene knock-down in drosophila melanogaster

	Samp1	Samp2	Samp3	Samp4	Samp5	Samp6	Samp7
FBgn0000003	0	0	0	0	0	0	1
FBgn0000008	92	161	76	70	140	88	70
FBgn0000014	5	1	0	0	4	0	0
FBgn0000015	0	2	1	2	1	0	0
FBgn0000017	4664	8714	3564	3150	6205	3072	3334
FBgn0000018	583	761	245	310	722	299	308
FBgn0000022	0	1	0	0	0	0	0
FBgn0000024	10	11	3	3	10	7	5
FBgn0000028	0	1	0	0	0	1	1
FBgn0000032	1446	1713	615	672	1698	696	767

Switch over to `bulkRNA_AnalysisScript.R`

DE definitions: what are we testing

Most generally:

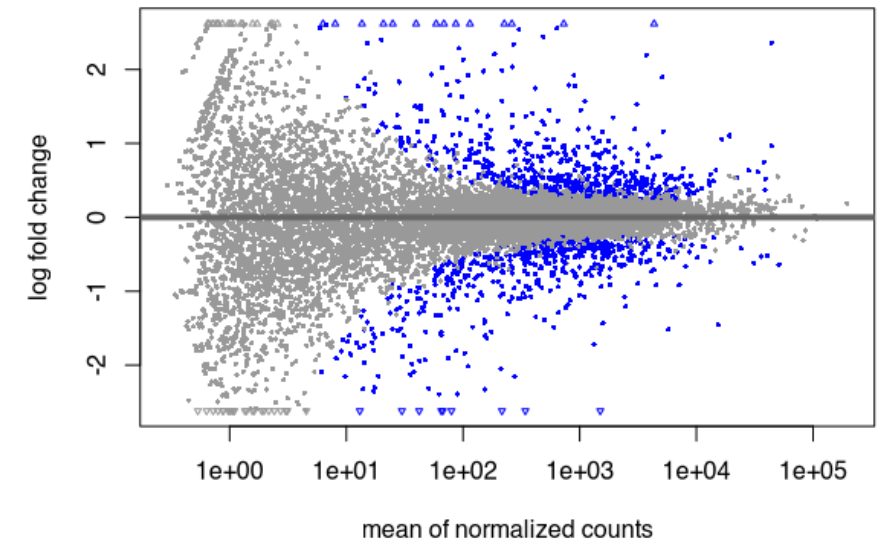
- Is the effect size of a given combination of covariates significantly non-zero

In practice:

- Is the abundance of gene g in condition A significantly different than in condition B controlling for appropriate nuisance variation

Implied comparison of interest:

- Ratios measured as log (2) fold-changes



DE definitions: typical model spec.

DESeq2 (and other models) assume observed expression is effectively modeled as a *Negative Binomial GLM*, i.e.:

$$y_{gj} \sim NB(\mu_{gj}, \phi_g)$$
$$\log(\mu_{gj}) := \vec{\beta}_g^T [X]_j + \log(s_j)$$

For observed counts y_{gj} , design matrix X , coefficient vector $\vec{\beta}$, normalization factor/offset s_j , and dispersion parameter ϕ_g .

DE definitions: typical model spec.

DESeq2 (and other models) assume observed expression is effectively modeled as a *Negative Binomial GLM*, i.e.:

$$y_{gj} \sim NB(\mu_{gj}, \phi_g)$$
$$\log(\mu_{gj}) := \vec{\beta}_g^T [X]_j + \log(s_j)$$

For observed counts y_{gj} , design matrix X , coefficient vector $\vec{\beta}$, normalization factor/offset s_j , and dispersion parameter ϕ_g .

NB GLM: Comp. to linear regression

Traditional “simple linear regression” can be written as a Normal GLM:

$$y_{gj} = \vec{\beta}_g^T [X]_j + \epsilon_{gj} \text{ s.t. } \epsilon_{gj} \sim N(0, \sigma_g^2)$$

is equivalent to:

$$\begin{aligned} y_{gj} &\sim N(\mu_{gj}, \sigma_g^2) \\ \mu_{gj} &:= \vec{\beta}_g^T [X]_j \end{aligned}$$

Compared to the NB case (same as previous slide):

$$\begin{aligned} y_{gj} &\sim NB(\mu_{gj}, \phi_g) \\ \log(\mu_{gj}) &:= \vec{\beta}_g^T [X]_j + \log(s_j) \end{aligned}$$

Identity link (Normal):

$$\mu_{gj} = \vec{\beta}_g^T [X]_j = \beta_{g0}x_j + \beta_{g1}x_j \cdots + \beta_{gK}x_j$$

Log link (Negative Binomial):

$$\begin{aligned} \log(\mu_{gj}) &= \vec{\beta}_g^T [X]_j + \log(s_j) \\ \Rightarrow \mu_{gj} &= (e^{\beta_{g0}x_j} e^{\beta_{g1}x_j} \cdots e^{\beta_{gK}x_j}) s_j \end{aligned}$$

Action of log link indicates **data Y should be raw expression** and not “normalized expression”

Why NB GLM (summary)

- We usually want to test **fold-change (ratios)**
 1. Statistical tests of *ratios* of random variables are challenging
 2. $\log\left(\frac{\mu_A}{\mu_B}\right) = \log(\mu_A) - \log(\mu_B) \approx \log(\bar{x}_A) - \log(\bar{x}_B) \neq \overline{\log(x_A)} - \overline{\log(x_B)}$
 3. NB GLM (with log-link) *directly* tests ratios
- Our data is often (small) **integers**, not continuous
 1. Central limit theorem is limited for small sample-sizes / high variance / highly discrete data
 2. NB GLM *assumes* data are non-negative integers
- Empirically, sequencing data is **overdispersed**
 1. Poisson testing yields *false positives*
 2. NB GLM accepts flexible dispersion parameters

(Pre) Calculation of normalization offset

In DESeq2 (and other methods) normalization factors s_j are offsets:

- s_j are calculated prior to model estimation (fixed)
- s_j have an implicit fixed coefficient $\beta_s = 1$

Popular norm. factor calculation methods include:

1. Reads-per-kilobase-million (RPKM)
2. Transcripts/Counts-per-million (TPM/CPM)
3. Library-size normalization
4. Trimmed-mean-of-m-values (TMM, default in edgeR)

Poor normalization (example)

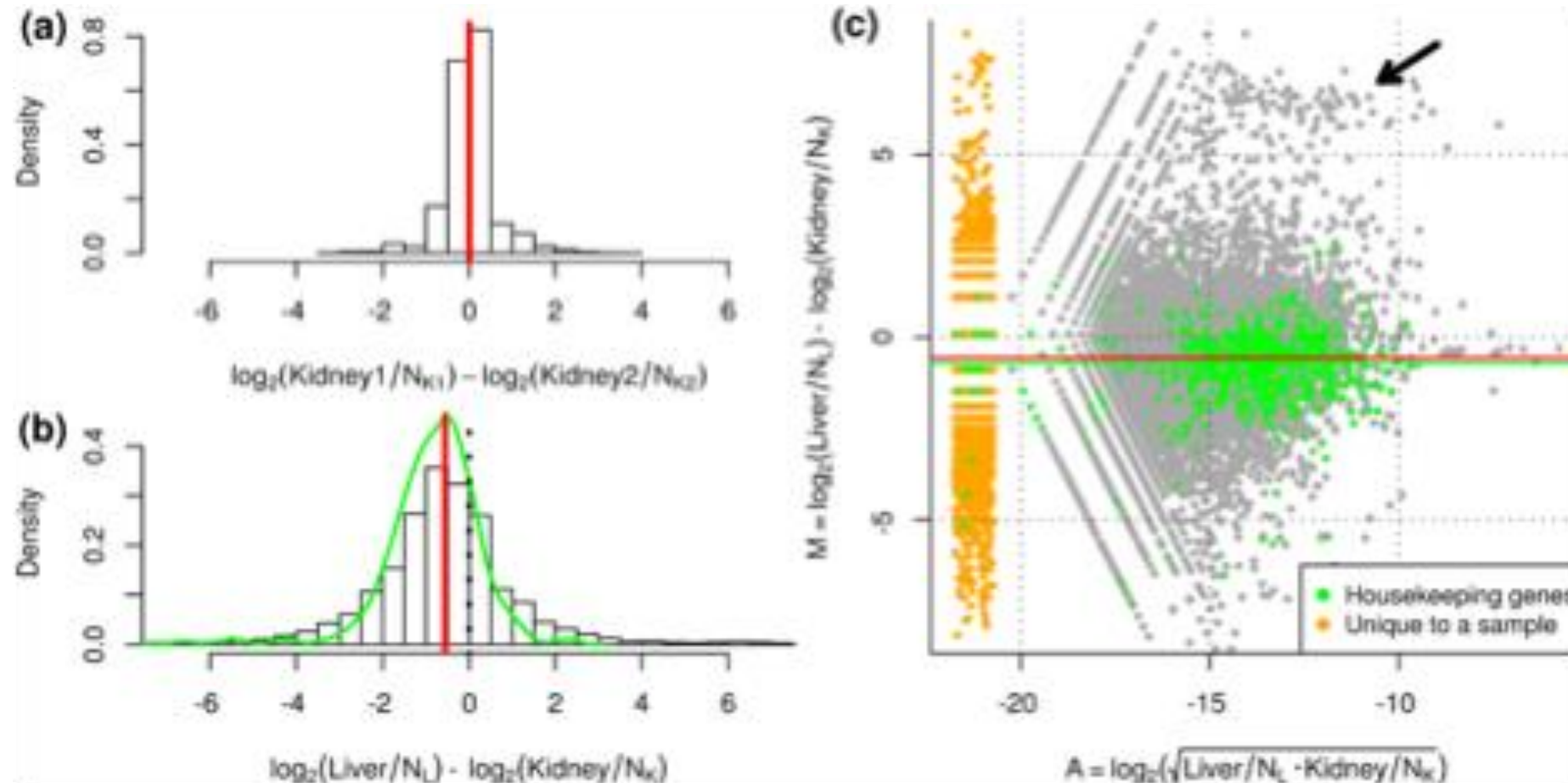


Fig. Ref: Robinson and Oshlack,
Genome Biology, 2010

DESeq2 normalization default

Method of *Median of Ratios* and *Scran*:

- Median of Ratios (MR) from Anders and Huber 2010
- Scran based on Lun, Bach, and Marioni 2016

Both methods based on the median of ratios relationship:

$$\hat{s}_j := \text{median}_g \left\{ \frac{\hat{y}_{gj}}{(\prod_j^J \hat{y}_{gj})^{\frac{1}{J}}} \right\} = \text{median}_g \left\{ \frac{\hat{y}_{gj}}{\exp \left[\frac{1}{J} \sum_j^J \log(\hat{y}_{gj}) \right]} \right\}$$

For MR, $\hat{y}_{gj} = y_{gj}$ and $s_j = \hat{s}_j$

For Scran, \hat{y}_{gj} is based on a pooling of cells and s_j is derived from a deconvolution of (multiple) estimates \hat{s}_j that are based on sample j

Switch over to `bulkRNA_AnalysisScript.R`

NB GLM: Dispersion estimates

Dispersion is analogous to variance, specifically:

$$y \sim NB(\mu, \phi) \Rightarrow \mathbb{E}[y] = \mu; \mathbb{V}[y] = \mu + \phi\mu^2$$

In typical bulk experiments, replicate counts are low:

- High standard error on fitted ϕ
- Reduced power to identify significant β_{gk}

So pool information across genes with similar expression levels...

Note: MLE of ϕ only exists when sample variance is greater than sample mean...

Switch over to `bulkRNA_AnalysisScript.R`

Standard multiple-regression is encoded in the design matrix X as:

$$\log(\mu_{gj}) = \vec{\beta}_g^T [X]_j + \log(s_j)$$

- Factors encoded in usual “dummy variable” format
- Default model specification chooses first factor as intercept
- For three group model, encoded as:

$$\log(\mu_{gj}) = \beta_{g1}x_j + \beta_{g2}x_j + \beta_{g3}x_j = \beta_{g1} \mid j \in \mathcal{A}$$

$$\log(\mu_{gj}) = \beta_{g1}x_j + \beta_{g2}x_j + \beta_{g3}x_j = \beta_{g1} + \beta_{g2} \mid j \in \mathcal{B}$$

$$\log(\mu_{gj}) = \beta_{g1}x_j + \beta_{g2}x_j + \beta_{g3}x_j = \beta_{g1} + \beta_{g3} \mid j \in \mathcal{C}$$

Contrasts and complex comparisons

We can use “contrasts” or *linear combinations of coefficients* to test non-default or complex hypotheses

- Typical DESeq2 reduces general *contrasts* to the difference between the sums of two groups of coefficients (general linear combination possible):

$$c = \sum_{k \in \mathcal{C}_1} \beta_k - \sum_{k \in \mathcal{C}_2} \beta_k$$
$$H_0: c = 0 \quad H_1: c \neq 0$$

- Equivalent to testing whether fold-change (ratio) between groups is different from 1; using groups from before:

$$\frac{\mu_{gB}}{\mu_{gC}} = \frac{e^{\beta_1 + \beta_2}}{e^{\beta_1 + \beta_3}} = e^{(\beta_1 + \beta_2) - (\beta_1 + \beta_3)} = e^{\beta_2 - \beta_3} \Rightarrow c = \beta_2 - \beta_3$$

Example 2 – ratio of ratios

More complex relationships possible:

- Desired test:

$$r = \frac{\left(\frac{\mu_{A1}}{\mu_{A2}}\right)}{\left(\frac{\mu_{B1}}{\mu_{B2}}\right)}$$
$$H_0: \log(r) = 0 \quad H_1: \log(r) \neq 0$$

- Key is to reformulate as single exponential:

$$\frac{\left(\frac{\mu_{A1}}{\mu_{A2}}\right)}{\left(\frac{\mu_{B1}}{\mu_{B2}}\right)} = \frac{\left(\frac{e^{\beta_{A1}}}{e^{\beta_{A2}}}\right)}{\left(\frac{e^{\beta_{B1}}}{e^{\beta_{B2}}}\right)} = \frac{e^{\beta_{A1}-\beta_{A2}}}{e^{\beta_{B1}-\beta_{B2}}} = e^{(\beta_{A1}+\beta_{B2})-(\beta_{A2}+\beta_{B1})} \Rightarrow \log(r) = (\beta_{A1} + \beta_{A2}) - (\beta_{A2} + \beta_{B1})$$

Switch over to `bulkRNA_AnalysisScript.R`

Shrinkage and calling sig. features

Goals:

- Maximize power
- Maintain control on FDR (or something like it...)
- Improve interpretability/reliability

Note 1: traditional methods for controlling FDR tend to be under-powered

Note 2: leveraging extra information can boost power

q-value and local false discovery rate

q-value (Corollary 2, J. Storey 2003) for statistic t and rejection region Γ_α :

$$q(t) := \inf_{\{\Gamma_\alpha: t \in \Gamma_\alpha\}} \mathbb{P}(H = 0 | T \in \Gamma_\alpha)$$

lfdr for a typical null hypothesis ($\beta_j = 0$):

$$lfdr_j := \mathbb{P}(\beta_j = 0 | \hat{\beta}_j, \hat{s}_j, \hat{\pi})$$

For a particular set of observed effect sizes, ordered by lfdr,

$$q(\hat{\beta}_{(j)}) = \frac{1}{j} \sum_{i=1}^j lfdr_{(i)}$$

Local false sign rate and s-value

lfsr (Eqn. 2.7 from M. Stephens 2017):

$$lfsr_j := \min[\mathbb{P}(\beta_j \geq 0 | \hat{\pi}, \hat{\beta}, s), \mathbb{P}(\beta_j \leq 0 | \hat{\pi}, \hat{\beta}, s)]$$

Tukey (1991): “All we know about the world teaches us that the effects of A and B are always different – in some decimal place – for any A and B”

For a particular set of observed effect sizes, ordered by lfsr,

$$s(\hat{\beta}_{(j)}) = \frac{1}{j} \sum_{i=1}^j lfsr_{(i)}$$

In practice lfsr can be more powerful than lfdi in the sense that calculated lfsr is closer to the true lfsr while still being conservative

Adaptive SHrinkage (ASH) model

Model from *M. Stephens 2017*

Consider *true* effects (model coefficients) $\beta = (\beta_1, \dots, \beta_J)$ and observed effect sizes $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_J)$ with corresponding estimated standard errors $\hat{s} = (\hat{s}_1, \dots, \hat{s}_J)$

We can conduct hypothesis testing based on the model of the true effects:

$$\mathbb{P}(\beta | \hat{\beta}, \hat{s}) \propto \mathbb{P}(\beta | \hat{s}) \mathbb{P}(\hat{\beta} | \beta, \hat{s})$$

Model the conditional distribution on true effects and observed effects as:

$$\begin{aligned}\mathbb{P}(\beta | \hat{s}, \pi) &= \prod_j g(\beta_j | \pi) \\ g(\cdot | \pi) &= \pi_0 \delta_0(\cdot) + \sum_{k=1}^K \pi_k N(\cdot | 0, \sigma_k^2) \\ \mathbb{P}(\hat{\beta} | \beta, \hat{s}) &= \prod_j N(\hat{\beta}_j | \beta_j, \hat{s}_j^2)\end{aligned}$$

Switch over to `bulkRNA_AnalysisScript.R`

Some considerations for single cell

Data sparsity causes normalization problems

- Use `scrn`

Data sparsity also affects lower bound on fitted dispersion parameters

- Set `minmu = 1e-6`

High sample (cell) counts trigger outlier correction

- Set `minReplicatesForReplace = Inf`

Higher variance (sparsity, heterogeneity, etc.) slow convergence to standard normal

- Set `useT = TRUE`

Switch over to `scRNA_AnalysisScript.R`