

Analysing single cell RNAseq

Laurent Modolo

March 16 2018

Table of Contents

- 1 single-cell RNA sequencing
- 2 Counting molecules
- 3 Cell quality control
- 4 Normalization
- 5 Dimension reduction
- 6 Differential expression analysis
- 7 Clustering

Table of Contents

- 1 single-cell RNA sequencing
- 2 Counting molecules
- 3 Cell quality control
- 4 Normalization
- 5 Dimension reduction
- 6 Differential expression analysis
- 7 Clustering

classical RNASeq (bulk RNASeq)

For each gene or transcript :

- We have a number of reads
- We can estimate the amount of RNA transcripts

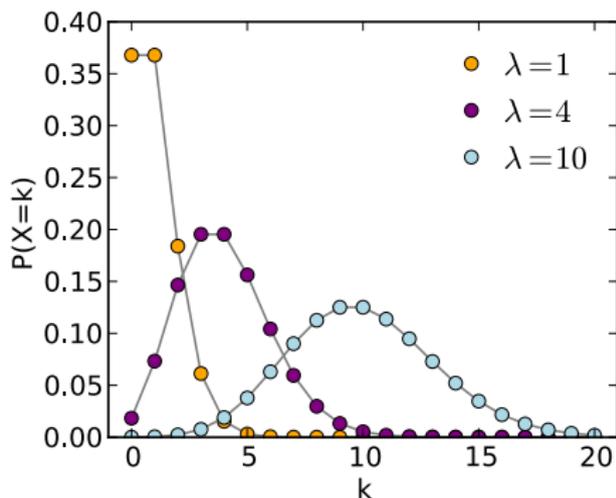
Compared to other RNA measurements we handle counts (Y positive integers)

The natural count distribution is the Poisson distribution.

$$Y_{ig} \sim \mathcal{P}(\lambda_{ig})$$

for the gene g in condition i .

- λ_{ig} RNA transcription rate



¹Gierlinski2015

classical RNASeq (bulk RNASeq)

For each gene or transcript :

- We have a number of reads
- We can estimate the amount of RNA transcripts

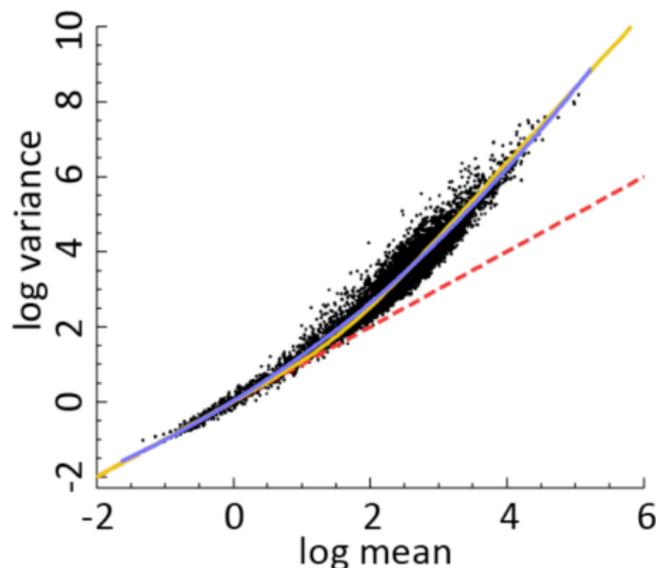
Compared to other RNA measurements we handle counts (Y positive integers)

The natural count distribution is the Poisson distribution.

$$Y_{ig} \sim \mathcal{P}(\lambda_{ig})$$

for the gene g in condition i .

- λ_{ig} RNA transcription rate



¹Gierlinski2015

classical RNASeq (bulk RNASeq)

For each gene or transcript :

- We have a number of reads
- We can estimate the amount of RNA transcripts

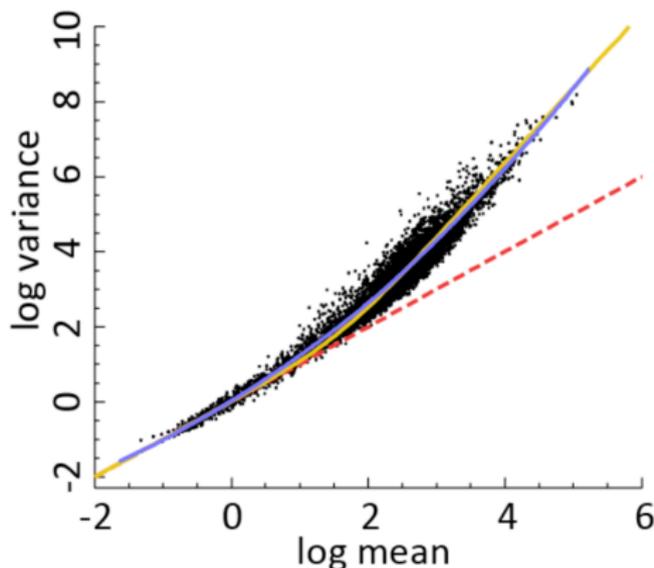
Compared to other RNA measurements we handle counts (Y positive integers)

The NegativeBinomial distribution allows for more or less variability.

$$Y_{ig} \sim NB(\mu_{ig}, \alpha_{ig})$$

for the gene g in condition i .

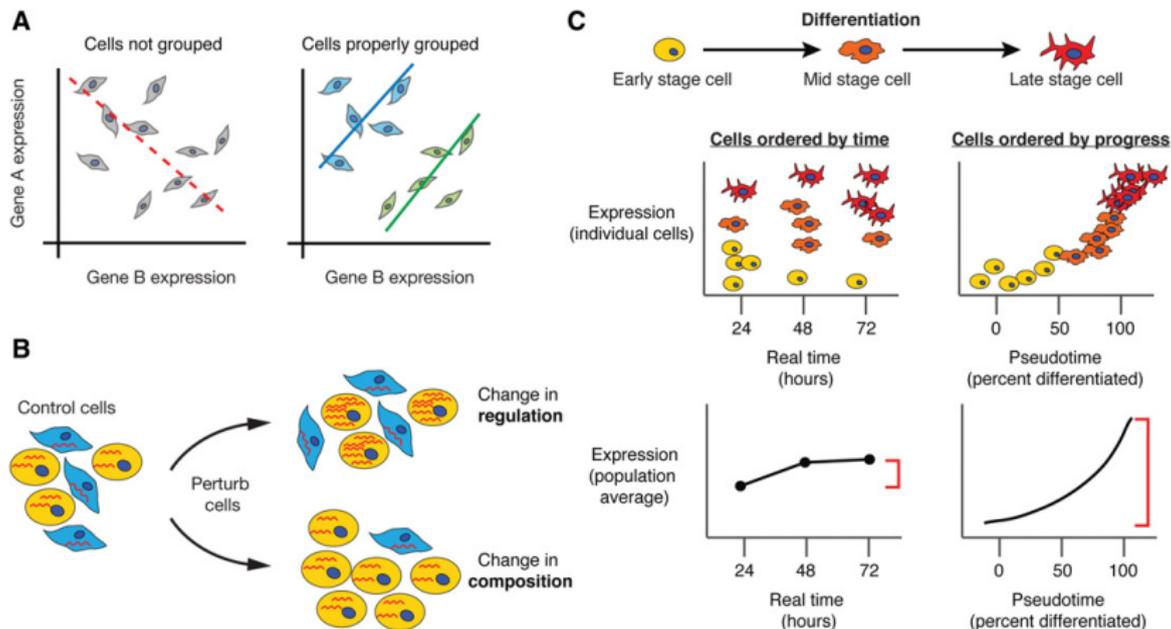
- μ_{ig} RNA transcription rate
- α_{ig} technical or biological variability



¹Gierlinski2015

single-cell RNA sequencing (scRNASeq)

Examples where biological variability stay hidden in bulk RNASeq experiments



¹Trapnell2015

Drawbacks

Compared to bulk RNASeq, we are not working with measurement on a population of cells:

- low starting amount of RNA (**zeros**)
- need to amplify (**more errors**)
- transcription status of each gene in each cell (**more zeros**)

Drawbacks

Compared to bulk RNASeq, we are not working with measurement on a population of cells:

- low starting amount of RNA (**zeros**)
- need to amplify (**more errors**)
- transcription status of each gene in each cell (**more zeros**)

$$\text{scRNASeq} = \text{bulk RNASeq} + \text{zeros} \times 2 + \text{noise}$$

single-cell RNA sequencing (scRNASeq)

The zero-inflated Negative Binomial distribution accounts for the excess of zeros.

$$Y_{ig} \sim \pi_{ig} \delta_0 + (1 - \pi_{ig}) NB(\mu_{ig}, \alpha_{ig})$$

for the gene g in condition i .

- μ_{ig} RNA transcription rate
- α_{ig} technical or biological variability
- π the proportion of additional zeros

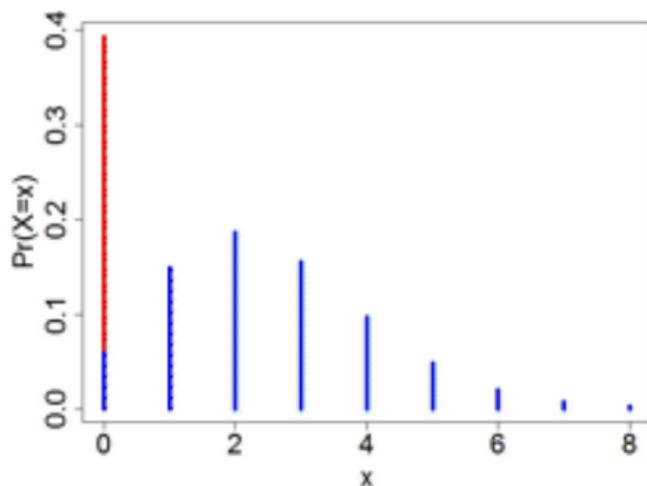


Table of Contents

- 1 single-cell RNA sequencing
- 2 Counting molecules**
- 3 Cell quality control
- 4 Normalization
- 5 Dimension reduction
- 6 Differential expression analysis
- 7 Clustering

Bulk or scRNASeq reads are the same thing

You can use the same quality control and mapping tools

¹**Smith2017**

²**Petukhov2017**

³**Bray2016**

Bulk or scRNASeq reads are the same thing

You can use the same quality control and mapping tools

Unique Molecular Identifiers UMI

- control of cDNA amplification
- we have access to the RNA transcripts counts (dropEst¹)
- risk of UMI collisions²
- alternative splicing?

¹Smith2017

²Petukhov2017

³Bray2016

Bulk or scRNASeq reads are the same thing

You can use the same quality control and mapping tools

Unique Molecular Identifiers UMI

- control of cDNA amplification
- we have access to the RNA transcripts counts (dropEst¹)
- risk of UMI collisions²
- alternative splicing?

Classical RNASeq

- alternative splicing (Kallisto³ with the pseudo mode)
- possible cDNA amplification bias
- we work with read counts

No tools to infer the transcript counts (yet)

¹Smith2017

²Petukhov2017

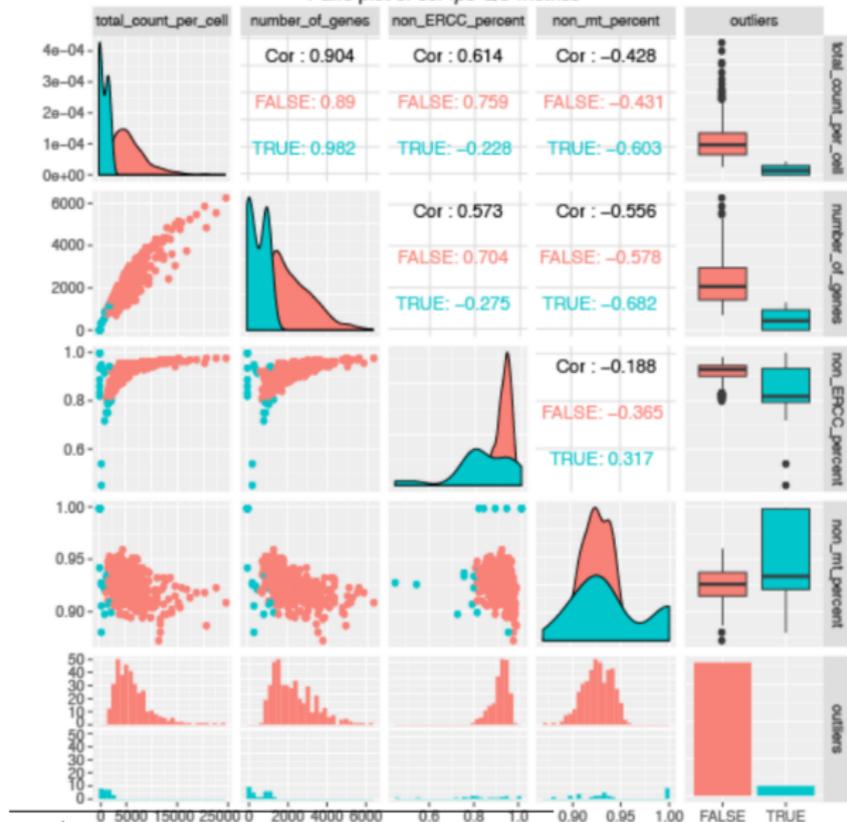
³Bray2016

Table of Contents

- 1 single-cell RNA sequencing
- 2 Counting molecules
- 3 Cell quality control**
- 4 Normalization
- 5 Dimension reduction
- 6 Differential expression analysis
- 7 Clustering

Low quality cells must be removed

Pairs plot of *scPipe* QC metrics



Identifying low quality cells:

- detection of outliers
- sequencing of blanks
- ERCC counts

Gaussian mixture model (*mclust*)

Support vector machines (*e1071*)

Table of Contents

- 1 single-cell RNA sequencing
- 2 Counting molecules
- 3 Cell quality control
- 4 Normalization**
- 5 Dimension reduction
- 6 Differential expression analysis
- 7 Clustering

Normalization

bulk RNASeq

- scaling between replicates
- as many replicate as we want
- (batch effects)

scRNASeq

- scaling between cells
- cells can be measured only once (for now)
- 10-60% mRNA capture efficiency
- large number of batches

Normalization

bulk RNASeq

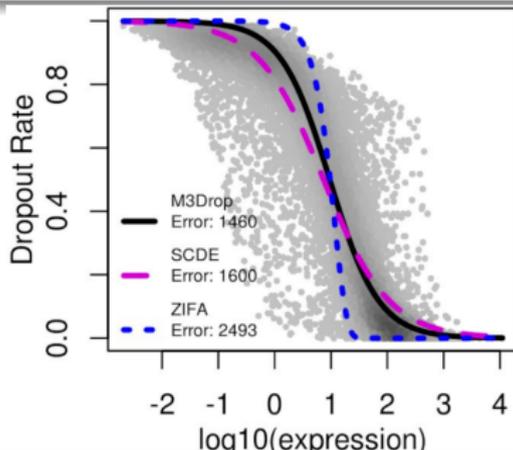
- scaling between replicates
- as many replicate as we want
- (batch effects)

scRNASeq

- scaling between cells
- cells can be measured only once (for now)
- 10-60% mRNA capture efficiency
- large number of batches

Need to normalize for:

- differences between cells
- differences between batches
- differences between genes



¹Andrews2018

Normalization

- single-cell scaling factor with SCnorm¹
- batch effect DASC²

¹**Batcher2017**

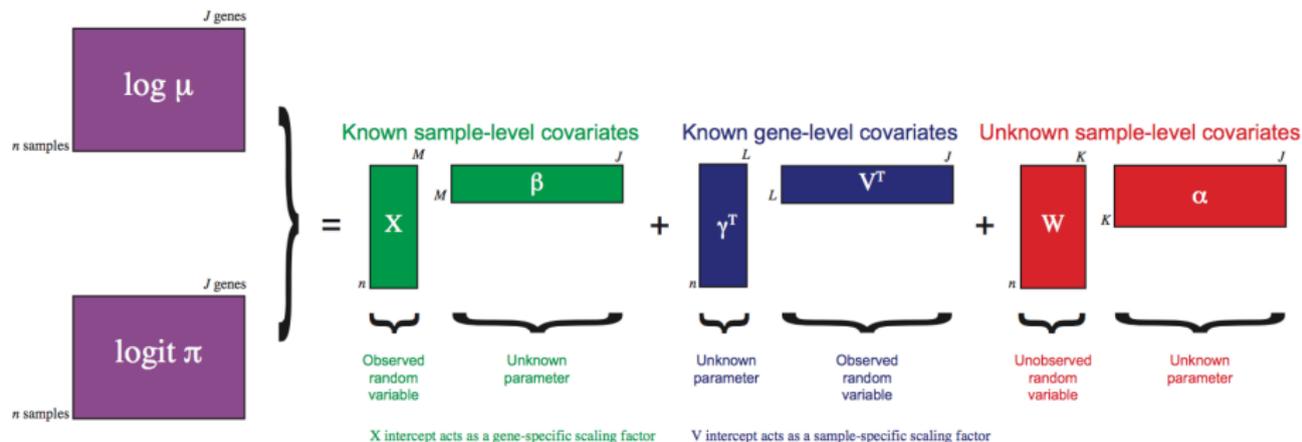
¹**Yi2017**

³**Risso2017**

⁴**Li2018**

Normalization

- single-cell scaling factor with SCnorm¹
- batch effect DASC²
- both ZINB-WaVe³



¹Batcher2017

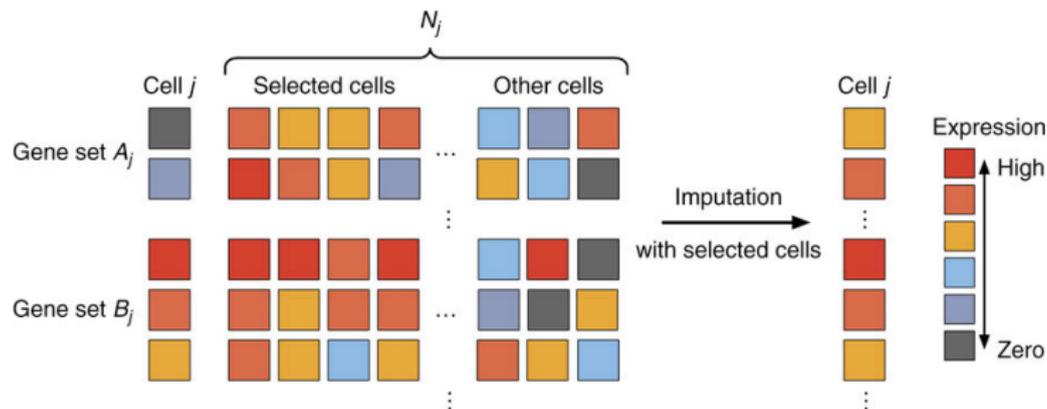
¹Yi2017

³Risso2017

⁴Li2018

Normalization

- single-cell scaling factor with SCnorm¹
- batch effect DASC²
- both ZINB-WaVe³
- correct dropout scImpute⁴



¹Batcher2017

¹Yi2017

³Risso2017

⁴Li2018

Scaling

scaled counts

$$\frac{Y_g}{\exp(\widehat{\alpha}_g)} \times (1 - \widehat{\pi}_g)$$

- zeros stay zeros
- the more zeros the less the gene will contribute
- we use the empirical dispersion

spread counts

- log-transform: $\log(Y + 1)$
- Anscomb transform: $\sqrt{Y + 3/8}$

Table of Contents

- 1 single-cell RNA sequencing
- 2 Counting molecules
- 3 Cell quality control
- 4 Normalization
- 5 Dimension reduction**
- 6 Differential expression analysis
- 7 Clustering

Dimension reduction

We have a large number of cells and a large number of genes

PCA

Classical PCA on scaled and spread counts

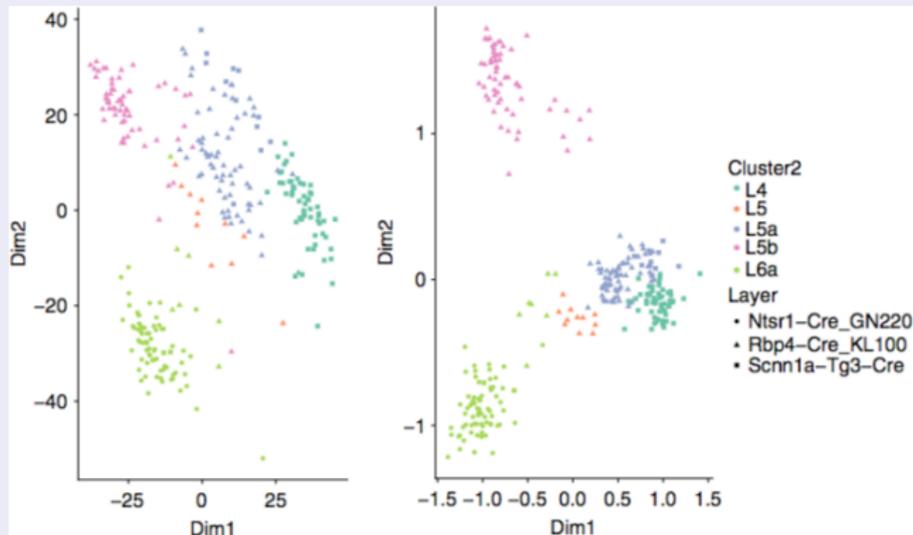
Dimension reduction

We have a large number of cells and a large number of genes

PCA

Classical PCA on scaled and spread counts

ziNB-WaVe



Dimension reduction

We have a large number of cells and a large number of genes

PCA

Classical PCA on scaled and spread counts

ziNB-WaVe

Using W

t-SNE on the above

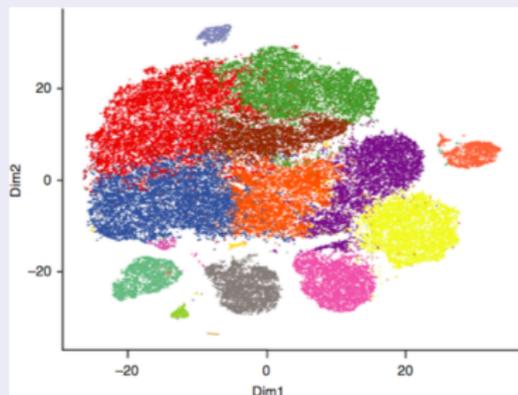
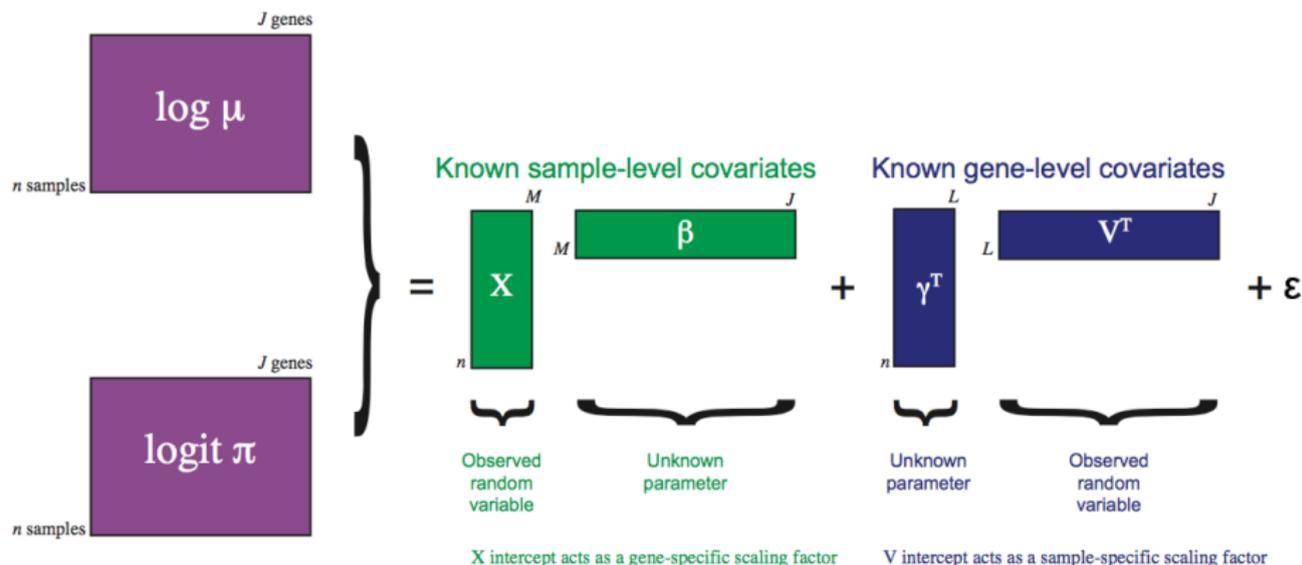


Table of Contents

- 1 single-cell RNA sequencing
- 2 Counting molecules
- 3 Cell quality control
- 4 Normalization
- 5 Dimension reduction
- 6 Differential expression analysis**
- 7 Clustering

Differential expression analysis



- with classical RNASeq tools `zinbwaveZinger`¹
- with zero-inflated NegativeBinomial GLM `pscl`, `glmmADMB`
- with dropout modelization `M3Drop`²

¹Risso2018

²Tallulah2018

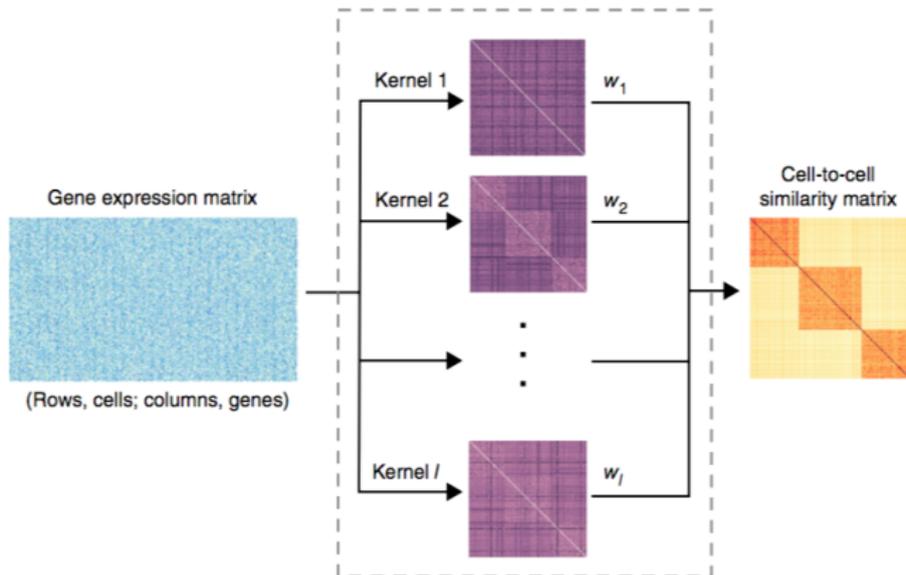
Table of Contents

- 1 single-cell RNA sequencing
- 2 Counting molecules
- 3 Cell quality control
- 4 Normalization
- 5 Dimension reduction
- 6 Differential expression analysis
- 7 Clustering**

Clustering

Euclidian distance don't work with more than 80% of zeros

- multiple kernel method SIMLR¹
- dropout imputation scImpute²



¹Wang2017

²Li2018

Thank you

$$Y_{ig} \sim \pi_{ig} \delta_0 + (1 - \pi_{ig}) NB(\mu_{ig}, \alpha_{ig})$$

for the gene g in condition i .

More than 200 tools at:

<http://www.scrna-tools.org/>

Tutorials and tools at:

<https://github.com/seandavi/awesome-single-cell>