Evolutionary Biology Department Seminar

Geteiltes Leid ist halbes Leid

Common difficulties in analysis and interpretation of single-cell RNA-seq data

2025-03-25 • Niko Papadopoulos

Other good resources

Current best practices in single-cell RNA-seq analysis: a tutorial Luecken & Theis, Molecular Systems Biology, 2019

Single-cell best practices book

(maintained & expanded by the Theis lab and the single-cell community) https://www.sc-best-practices.org/preamble.html

	Cell1	Cell2		CellN
Gene1	3	2		13
Gene2	2	3		1
Gene3	1	14		18
	•			
	•		·	•
GeneM	25	0		0

	Cell1	Cell2		CellN
Gene1	3	2		13
Gene2	2	3		1
Gene3	1	14		18
	•	•	•	•
	•			
GeneM	25	0	·	0























Why is scRNA-seq analysis so hard?

- To get the most out of the data you (probably) need non-standard analysis
- Not clear what standard analysis is anyway
- Too many knobs to turn \rightarrow analysis paralysis
- Often: biological questions unclear \rightarrow hard to formulate testable hypotheses
- "Best" scRNA-seq skillset:
 - Scripting
 - Programming
 - Statistics
 - (cellular/molecular) biology
 - (organismal/tissue/disease) biology

 \rightarrow does not exist

Why is scRNA-seq analysis so hard?

- To get the most out of the data you (probably) need non-standard analysis
- Not clear what standard analysis is anyway
- Too many knobs to turn \rightarrow analysis paralysis
- Often: biological questions unclear \rightarrow hard to formulate testable hypotheses
- "Best" scRNA-seq skillset:
 - Scripting
 - Programming
 - Statistics
 - (cellular/molecular) biology
 - (organismal/tissue/disease) biology
 - \rightarrow does not exist

- Recognize gaps
- Prioritize
- Develop "working understanding"
- Focus on biology/questions
- Ask for help









it's all about cell-cell distances

cell filtering

normalisation

variance stabilisation

it's all about cell-cell distances

dimensionality reduction

informative genes

distance measures



























Genes						
	0	0	0	0	0	2
	0	0	5	4	6	6
	0	0	6	7	6	0
	7	6	6	0	0	0



All about distances (2) - mitochondrial content

Tangent - mitochondrial content

Mitochondrial genes are expressed in most cells, and their expression level is cell type-specific.

High expression levels of mitochondrial genes could be an indicator of:

- 1. Poor sample quality, leading to a high fraction of apoptotic or lysing cells.
- 2. Biology of the particular sample, for example, tumor biopsies, may have increased mitochondrial gene expression due to metabolic activity and/or necrosis.

- 10X Genomics, 2017
• dead cells% ~ MT%?

• dead cells% ~ MT%?

 \rightarrow try different combinations of happy/unhappy/dying cells* & sequence

• dead cells% ~ MT%?

 \rightarrow try different combinations of happy/unhappy/dying cells* & sequence

* human PBMCs from a healthy donor

• dead cells% ~ MT%?

 \rightarrow try different combinations of happy/unhappy/dying cells* & sequence

Sample	Name	What happened
1	control	sample prep into seq
2	Room temperature (24h)	Clearly unhappy cells
3	Digitonin Low	1:1 mix of dying (digitonin) and happy cells
4	Digitonin High	5:1 mix of dying/happy

* human PBMCs from a healthy donor



https://support.10xgenomics.com/single-cell-gene-expression/index/doc/technical-note-removal-of-dead-cells-from-single-cell-suspensions-improves-performance-for-10x-genomics-single-cell-applications.

















• Doublets & low-quality cells (usually) increase noise in data*

- Doublets & low-quality cells (usually) increase noise in data*
- Mammalian cut-offs probably don't apply

- Doublets & low-quality cells (usually) increase noise in data*
- Mammalian cut-offs probably don't apply
- No silver bullet :(

- Doublets & low-quality cells (usually) increase noise in data*
- Mammalian cut-offs probably don't apply
- No silver bullet :(
- Compare to other indicators of cell quality in context

- Doublets & low-quality cells (usually) increase noise in data*
- Mammalian cut-offs probably don't apply
- No silver bullet :(
- Compare to other indicators of cell quality in context
 - #genes detected
 - #total reads
 - %ribosomal content
 - %highest expressed genes
 - Compare to housekeeping genes**
 - Ratio of useful reads/mito+housekeeping

* technically possible to redeem, under very specific circumstances. It very probably doesn't apply to your project. ** if you can get a list of them for your species

- Doublets & low-quality cells (usually) increase noise in data*
- Mammalian cut-offs probably don't apply
- No silver bullet :(
- Compare to other indicators of cell quality in context
 - #genes detected
 - #total reads
 - %ribosomal content
 - %highest expressed genes
 - Compare to housekeeping genes**
 - Ratio of useful reads/mito+housekeeping

* technically possible to redeem, under very specific circumstances. It very probably doesn't apply to your project. ** if you can get a list of them for your species







- Doublets & low-quality cells (usually) increase
- Mammalian cut-offs probably don't a plan
- No silver bullet :(
- Compare to other inclusions of a quality in context
 - #gen de tou
 - ALOTA reads
 - %ribo mal content
 - %higher expressed genes
 - Compare to housekeeping gent
 - Ratio of useful reads/mito puse

* technically possible to reducin, under very specific circumstances. It very probably doesn't apply to your project. ** if you can get a list of them for your species

Overall summary

- scRNA-seq analysis has many moving parts from different areas of expertise
 → easy to have blind spots
- Cell-to-cell distances at the core of (most) analysis
- How is <insert choice> affecting the (pattern of) distances between cells?
- QC: try to combine "orthogonal" indicators
- QC: use common sense + knowledge of system
- QC: When in doubt...
 - Is this crucial for answering my question?
 - Can I validate this with other data?

Ideas graveyard











Cell barcode UMI	RNA	Demultiplexing	Cell-specific reads
ACAGTATAAAGACT.			TGACAATAAAGACTTCTAGCTG TGACAAGTTACGTCACAATGCT TGACAATGATGCCGGTCACATC
CGTTAGGTTACGTC. TGACAAGTTACGTC. GTTAGCTGATGCCG.	GATTATAG		ACAGTATAAAGACTGGGCCCCG ACAGTAGTTACGTCGTCACATC ACAGTATGATGCCGTCGACGAT
GTTAGCTGATGCCG. CGTTAGTGATGCCG. ACAGTAGTTACGTC. TGACAATGATGCCG.			GTTAGCTGATGCCGCTTTGCAT GTTAGCTGATGCCGTCTCGACT GTTAGCTAAAGACTACATGCTG
ACAGTATGATGCCG. GTTAGCTAAAGACT. CGTTAGGTTACGTC.	TCGACGAT ACATGCTG TAGCCAGT		CGTTAGGTTACGTC GATTATAG CGTTAGTGATGCCG CCTCGAGC CGTTAGGTTACGTC TAGCCAGT





Cell barcode UMI	RNA	Demultiplexing	Cell-specific reads
ACAGTATAAAGACT. TGACAATAAAGACT.			TGACAATAAAGACTTCTAGCTG TGACAAGTTACGTCACAATGCT TGACAATGATGCCGGTCACATC
CGTTAGGTTACGTC. TGACAAGTTACGTC. GTTAGCTGATGCCG.	GATTATAG		ACAGTATAAAGACTGGGCCCCG ACAGTAGTTACGTCGTCACATC ACAGTATGATGCCGTCGACGAT
CGTTAGCTGATGCCG. CGTTAGTGATGCCG. ACAGTAGTTACGTC. TGACAATGATGCCG.	CCTCGAGC GTCACATC GTCACATC		GTTAGCTGATGCCGCTTTGCAT GTTAGCTGATGCCGTCTCGACT GTTAGCTAAAGACTACATGCTG
ACAGTATGATGCCG. GTTAGCTAAAGACT. CGTTAGGTTACGTC.	TCGACGAT ACATGCTG TAGCCAGT		CGTTAGGTTACGTC GATTATAG CGTTAGTGATGCCG CCTCGAGC CGTTAGGTTACGTC TAGCCAGT

Reference genome







Cell barcode UMI	RNA	Demultiplexing	Cell-specific reads
ACAGTATAAAGACT.	GGGCCCCG		TGACAATAAAGACTTCTAGCTG TGACAAGTTACGTCACAATGCT TGACAATGATGCCGGTCACATC
CGTTAGGTTACGTC. TGACAAGTTACGTC. GTTAGCTGATGCCG.	GATTATAG		ACAGTATAAAGACTGGGCCCCG ACAGTAGTTACGTCGTCACATC ACAGTATGATGCCGTCGACGAT
GTTAGCTGATGCCG. CGTTAGTGATGCCG. ACAGTAGTTACGTC. TGACAATGATGCCG.	TCTCGACT CCTCGAGC GTCACATC GTCACATC		GTTAGCTGATGCCGCTTTGCAT GTTAGCTGATGCCGTCTCGACT GTTAGCTAAAGACTACATGCTG
ACAGTATGATGCCG. GTTAGCTAAAGACT. CGTTAGGTTACGTC.	TCGACGAT ACATGCTG TAGCCAGT		CGTTAGGTTACGTCGATTATAG CGTTAGTGATGCCGCCTCGAGC CGTTAGGTTACGTCTAGCCAGT





	Cell1	Cell2	 CellN
Gene1	3	2	13
Gene2	2	3	1
Gene3	1	14	18
	· ·		
	.		•
	·		
GeneM	25	0	0





Cell barcode UMI	RNA	Demultiplexing	Cell-specific reads
ACAGTATAAAGACT			TGACAATAAAGACTTCTAGCTG TGACAAGTTACGTCACAATGCT TGACAATGATGCCGGTCACATC
CGTTAGGTTACGTC TGACAAGTTACGTC GTTAGCTGATGCCG	GATTATAG		ACAGTATAAAGACTGGGCCCCG ACAGTAGTTACGTCGTCACATC ACAGTATGATGCCGTCGACGAT
GTTAGCTGATGCCG CGTTAGTGATGCCG ACAGTAGTTACGTC TGACAATGATGCCG	TCTCGACT CCTCGAGC GTCACATC GTCACATC		GTTAGCTGATGCCGCTTTGCAT GTTAGCTGATGCCGTCTCGACT GTTAGCTAAAGACTACATGCTG
ACAGTATGATGCCG GTTAGCTAAAGACT CGTTAGGTTACGTC	TCGACGAT ACATGCTG TAGCCAGT		CGTTAGGTTACGTCGATTATAG CGTTAGTGATGCCGCCTCGAGC CGTTAGGTTACGTCTAGCCAGT

Reference genome



	Cell1	Cell2		CellN
Gene1	3	2		13
Gene2	2	3		1
Gene3	1	14		18
	· ·			
	.		•	
	.			
GeneM	25	0		0














Why methods matter/hidden assumptions

- scRNA-seq == distances between cells; the perpetual struggle between noise and signal
- Anything that distorts the distances can change our analysis!
 - Normalisation: changes the values
 → changes distances!
 - Variance stabilisation: changes values
 - \rightarrow changes distances!
 - Mitochondrial/ribosomal counts: keep/remove? before/after normalisation?
 → changes distances!
 - Doublets: keep/remove?
 - \rightarrow changes distance *patterns*

Confusing methods for biology

- Clustering == cell types
- Highly variable genes == informative genes

Data hygiene

- Keep track of (and upload!) mapping reference (transcriptome/genome/proteome)
 - What choices were made (e.g. only protein-coding genes?)
 - State of annotation/reference: e.g. UTRs often extremely important for mapping rate
 - Include gene IDs in sc file/plots (ScanPy: use as .var index)
 - Keep track of annotation (BLAST/emapper/whatever tables)
- Bare minimum:
 - Raw count matrix
 - Filtering results (genes/cells)
 - Clustering results (corresponds to paper)