

AARHUS UNIVERSITY

BIOINFORMATICS MASTER THESIS

BIOINFORMATICS RESEARCH CENTER & DEPARTMENT OF MOLECULAR MEDICINE

Sequence-based neural networks applied to denoising of single-cell RNAseq data

Author:

Jens Brogård Stenbye

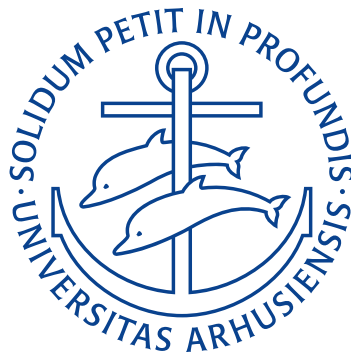
Stdnr: 201505646

Supervisors:

Lasse Maretty Sørensen

Palle Villesen

April 30, 2022



Contents

1	Acknowledgements	3
2	Abstract	3
3	Introduction	4
3.1	Single cell RNA-sequencing, application and problems	4
3.2	Denoising scRNA data	6
3.2.a	Non model-based approaches	6
3.2.b	Model-based approaches	6
3.3	Neural-network model approaches	8
3.3.a	Neural network introduction	8
3.4	Autoencoders in the field of scRNAseq	11
3.4.a	scRNAseq data denoising Autoencoders	11
3.5	Sequence-information in expression	12
3.5.a	Neural networks in genomics	12
3.6	Aim of this project	13
4	Results	14
4.1	Data and preprocessing	14
4.1.a	Cleaning and preprocessing scRNA count data	14
4.1.b	Transcription start site sequence extraction	16
4.2	Implementation of a benchmark model, using DCA architecture	16
4.2.a	Implementing the DCA architecture	17
4.2.b	Training the DCA model	18
4.3	Implementing a sequence- and count-based autoencoder	19
4.3.a	Implementing the sequence dependent module	19
4.3.b	Combining the modules	19
4.3.c	Training the SCELD model	21
4.4	Evaluating model imputation performance	21
4.5	Investigating clustering performance	24
5	Discussion	25
6	Conclusion	27

7	Methods	28
7.1	Code availability	28
7.2	Preprocessing of data	28
7.3	Neural network implementation	28
7.3.a	DCA	28
7.3.b	SCELD	29
7.3.c	Negative binomial loss function	30
7.4	Data-corruption/downsampling	30
7.5	Clustering analysis	31
A	Appendix	36
A.1	supplementary figures	36

1 Acknowledgements

I would like to thank the following people who have been integral to me completing this project. First and foremost a massive thanks to my supervisor Lasse Maretty Sørensen not only for his insight and supervision regarding the project but also his kindness and understanding of the troubles i encountered during the project.

Furthermore i would like to extend a special thanks to Christian Storm and Palle Villesen from BiRC for their help, patience and understanding in the administrative issues i encountered.

Finally i would like to thank my family and also my good friend Erik Gunnersen whose support have been invaluable during the making of this thesis.

2 Abstract

Single-cell RNA-sequencing(scRNAseq) offers many exciting analysis opportunities. Unlike bulk RNA-sequencing, scRNAseq gives an expression profile for each individual cell in a sample. This allows for identification of rare sample cell-subtypes, analysis of cell differentiation among other interesting analysis. The sequencing depth of each cell in scRNAseq is however very small, leading to many dropout events and general overdispersed data.

Several methods have been developed in order to denoise scRNAseq data, imputing the dropout events and correcting counts. In recent years, several neural network models have been developed in the field of scRNAseq denoising. Common for all current modelling approaches is that they only utilize the scRNAseq gene-counts themselves to denoise the counts.

In this thesis we argue that the inherent sequence context of each gene holds expression information, and that this could be utilized in inferring the denoised gene-counts. To that end we develop the Sequence- and Count-based Encoder and Linear Decoder(SCELD) model which utilizes sequence information in addition to the count-matrix. We also implement the Deep-CountAutoencoder(DCA) an existing scRNA-denoising model as a benchmark model.

In general the SCELD was slightly outperformed by DCA across the several denoising metrics. Especially when investigating the ability to retain the clustering of cells when trained on a corrupted dataset did SCELD perform poorly while DCA was able to reproduce the cell clusters. The choice of benchmarking data and the model design may not have been optimal however, and we argue that further research is needed to properly evaluate the use of sequence information in scRNAseq denoising.