

SCNet: Automatic Multi-Channel Genome Network Inference from Single-Cell RNA Sequences

Shiyin Wang (wangshiy16@mails.tsinghua.edu.cn)

December 23, 2019

Outline

1 Introduction

- Single-Cell RNA Sequencing
- Mathematics Preliminaries

2 Methods

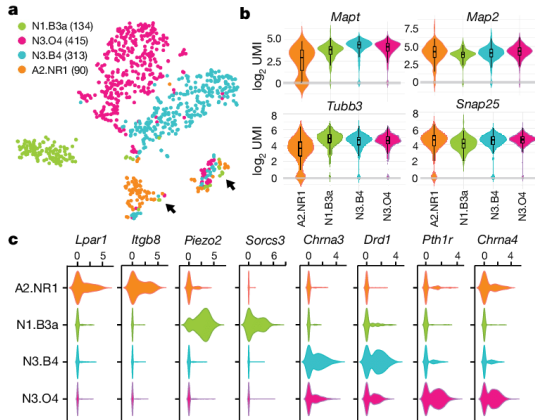
- Exploratory Data Analysis
- Probabilistic models
- Optimization

3 Results

4 Future Work

5 Q & A

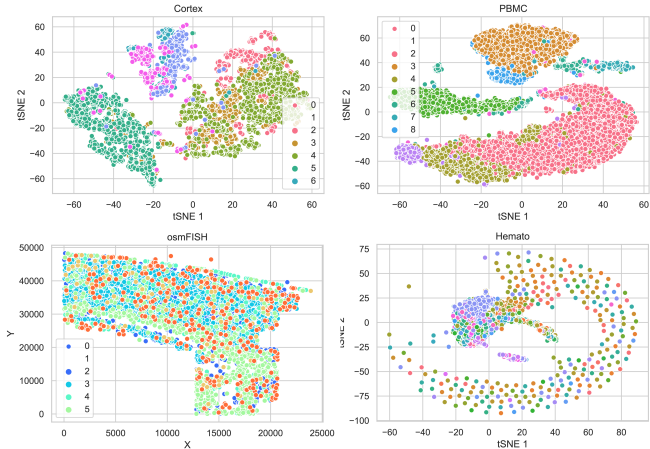
What is single-cell RNA sequencing?



Tsunemoto, Rachel, et al. "Diverse reprogramming codes for neuronal identity." *Nature* 557.7705 (2018): 375.



Datasets



Bayesian Methods

Likelihood

How probable is the evidence
given that our hypothesis is true?

Prior

How probable was our hypothesis
before observing the evidence?

$$P(H | e) = \frac{P(e | H) P(H)}{P(e)}$$

Posterior

How probable is our hypothesis
given the observed evidence?
(Not directly computable)

Marginal

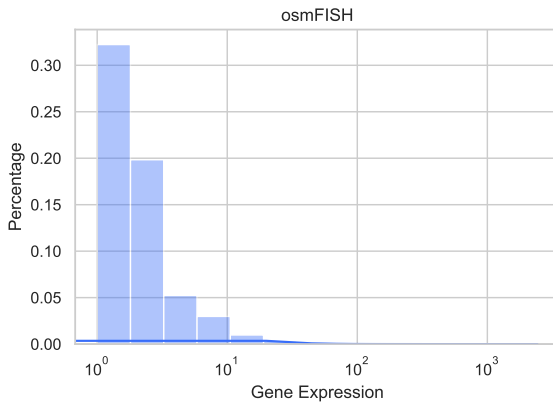
How probable is the new evidence
under all possible hypotheses?
 $P(e) = \sum P(e | H_i) P(H_i)$

Markov Chain Monte Carlo

- Initialise $x^{(0)}$.
- For $i = 0$ to $N - 1$
 - Sample $u \sim U_{[0,1]}$.
 - Sample $x^* \sim q(x^*|x^{(i)})$.
 - If $u < A(x^{(i)}, x^*) = \min \left\{ 1, \frac{p(x^*)q(x^{(i)}|x^*)}{p(x^{(i)})q(x^*|x^{(i)})} \right\}$
$$x^{(i+1)} = x^*$$

else
$$x^{(i+1)} = x^{(i)}$$

Choose Log-Scale to Model Gene Expression



Bayesian Hierarchical Linear Model

X and Y are the expression profiles of two genes. The edge weights in the desired network correspond to the regression coefficient k in the equation.

$$\begin{aligned} \log Y &= k \log X + \epsilon \\ k &\sim N(\beta, \sigma) \\ \epsilon &\sim N(0, \gamma^2) \end{aligned} \tag{1}$$

This model estimates k from single-cell RNA sequencing records.

Define Transition Probability through Edge Weights

Once we have a weighted network, we can define the association probability between two nodes X and Y by a very trivial model.

$$Pr(X \rightarrow Y) = 1 - (1 - w_{X,Y}) \prod_{a \in V} (1 - w_{X,a} w_{a,Y}) \prod_{a \in V, b \in V} (1 - w_{X,a} w_{a,b} w_{b,Y}) \quad (2)$$

For simplicity, I expanded the search for three steps. Walk length is flexible to choose.

Maximal Likelihood Optimization

Bayesian hierarchical linear model and transition probability explain the network dynamics from two different angles. Now we can combine them together to infer the edge weights of networks (V, E, W) .

$$\text{maximize } L(W;_k, \Sigma_k) = \sum_{v_1 \in V} \sum_{v_2 \in V} Pr_{N(k, \Sigma_k)}(w = Pr(v_1 \rightarrow v_2)) \quad (3)$$

Add regulation term to restrain the number of edges in the network.

$$\text{maximize } L(W;_k, \Sigma_k) = \sum_{v_1 \in V} \sum_{v_2 \in V} Pr_{N(k, \Sigma_k)}(w = Pr(v_1 \rightarrow v_2)) - \lambda |V| \quad (4)$$

Markov Chain Monte Carlo

Possible Directions

- Integrate existing knowledge from protein-protein interaction networks (STRING, OmniPATH, ConsensusPath, etc) as priors
- Design better probability models
- Make interactive transition videos

Resources

- 10x Genomics: Datasets providing single cell and spatial views of biological systems (<https://www.10xgenomics.com>)

Questions & Answers

