

RNA-Seq and Single-Cell RNA-Seq Tertiary Analysis

Yue Zhang^{1,2} and Ramesh Nair²

¹Genetics Bioinformatics Service Center (GBSC) and ²Stanford Center for Genomics and Personalized Medicine (SCGPM)

RNA-Seq Tertiary Analysis		Single-Cell RNA-Seq Tertiary Analysis																																					
Weighted Gene Co-expression Network Analysis (WGCNA) R Package		Single-Cell Differential Expression (SCDE) R Package																																					
<h3>1. WGCNA Publication</h3> <p>Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network analysis. <i>BMC Bioinformatics</i> 9, 559 (2008).</p> <p>Citations: 1354</p>	<h3>4. Module Detection</h3> <p>Genes correspond to rows and columns</p> <p>Hierarchical clustering dendrogram</p> <p>Connectivity matrix</p> <p>Module: Correspond to branches</p>	<h3>1. SCDE Publications</h3> <p>Kharchenko, P.V., Silberstein, L. & Scadden, D.T. Bayesian approach to single-cell differential expression analysis. <i>Nat. Methods</i> 11, 740–742 (2014).</p> <p>Fan, J. <i>et al.</i> Characterizing transcriptional heterogeneity through pathway and gene set overdispersion analysis. <i>Nat. Methods</i> 13, 241–244 (2016).</p>	<h3>4. DE Gene Detection</h3> <p>IK1</p> <p>IK2</p> <p>MLE: 1.19 95% CI: 0.93 - 1.46 Z = 7.16 pZ = 17.16</p>																																				
<h3>2. Input Data – Gene Expression Matrix</h3> <p>Gene expression</p> <p>sample1 sample2 sample3 sample4 sample5</p> <p>gene1 gene2 gene3 gene4 gene5</p>	<h3>5. Relate Module With Sample Trait</h3> <p>Network visualization (Cytoscape)</p> <p>Module and sample trait relationship</p>	<h3>2. Input Data – Raw Count Matrix</h3> <table border="1"> <thead> <tr> <th></th> <th>cell1</th> <th>cell2</th> <th>cell3</th> <th>cell4</th> <th>cell5</th> </tr> </thead> <tbody> <tr> <td>gene1</td> <td>0</td> <td>2</td> <td>0</td> <td>0</td> <td>0</td> </tr> <tr> <td>gene2</td> <td>12</td> <td>34</td> <td>52</td> <td>24</td> <td>0</td> </tr> <tr> <td>gene3</td> <td>0</td> <td>456</td> <td>756</td> <td>365</td> <td>0</td> </tr> <tr> <td>gene4</td> <td>0</td> <td>68</td> <td>76</td> <td>1</td> <td>1</td> </tr> <tr> <td>gene5</td> <td>0</td> <td>18</td> <td>0</td> <td>14</td> <td>0</td> </tr> </tbody> </table>		cell1	cell2	cell3	cell4	cell5	gene1	0	2	0	0	0	gene2	12	34	52	24	0	gene3	0	456	756	365	0	gene4	0	68	76	1	1	gene5	0	18	0	14	0	<h3>5. Subpopulation Detection</h3> <p>Weighted PCA on noisy datasets with missing values make the result more sensitive to the true underlying signal variations.</p> <p>Cell hierarchical clustering is performed on the eigengene for each geneset.</p>
	cell1	cell2	cell3	cell4	cell5																																		
gene1	0	2	0	0	0																																		
gene2	12	34	52	24	0																																		
gene3	0	456	756	365	0																																		
gene4	0	68	76	1	1																																		
gene5	0	18	0	14	0																																		
<h3>3. Scale-free Network Fitness</h3> <p>$P(k) \sim k^{-\gamma}$</p> <p>Many nodes have few connections</p> <p>Few nodes have many connections</p> <p>$k_i = \sum_j a_{ij}$</p> <p>Scale-free network fitness</p>	<h3>6. BaaS Use Case</h3> <p>BaaS project: Drug-induced pathway disruption for cardiac differentiation.</p> <p>Input data and clean data: FPKM matrix in log scale was generated from study samples. ~10K genes remained after filtering out low expression genes.</p> <p>Day 6, Treated Day 6, Treated Day 2, Treated Day 2, Control Day 0, Control</p> <p>GO:0072358 cardiovascular system development</p> <p>GO:0007507 heart development</p>	<h3>3. Statistical Methods</h3> <ol style="list-style-type: none"> 1. Find a group of nearest neighbor cells for each cell using Pearson correlation. 2. Estimate expected expression for each gene in each cell. 3. Build probabilistic model with the observed expression value of each gene in each cell modeled as a mixture of a dropout and an amplified component. 4. Generate a corresponding weight matrix measuring dropout event probability which is used for weighted PCA $p(x r_c, \Omega_c) = p_d(x) p_{\text{Poisson}}(x) + (1 - p_d(x)) p_{NB}(x r_c)$ <table border="1"> <thead> <tr> <th></th> <th>cell1</th> <th>cell2</th> <th>cell3</th> <th>cell4</th> <th>cell5</th> </tr> </thead> <tbody> <tr> <td>gene1</td> <td>0</td> <td>0.002</td> <td>0</td> <td>0</td> <td>0</td> </tr> <tr> <td>gene2</td> <td>0.003</td> <td>0.034</td> <td>0.052</td> <td>0.024</td> <td>0</td> </tr> <tr> <td>gene3</td> <td>0</td> <td>0.156</td> <td>0.589</td> <td>0.265</td> <td>0</td> </tr> <tr> <td>gene4</td> <td>0</td> <td>0.109</td> <td>0.276</td> <td>0.009</td> <td>0</td> </tr> <tr> <td>gene5</td> <td>0</td> <td>0.028</td> <td>0</td> <td>0.021</td> <td>0</td> </tr> </tbody> </table> <p>Prior distribution: Learned from expression data Likelihood function: From probabilistic model Posterior distribution: Prior * Likelihood Joint posterior distribution of a gene being expressed at an average level x in a subpopulation of cells S is determined as an expected value from Bootstrap sampling</p> $p_S(x) = E \left[\prod_{c \in S} p(x r_c, \Omega_c) \right]$		cell1	cell2	cell3	cell4	cell5	gene1	0	0.002	0	0	0	gene2	0.003	0.034	0.052	0.024	0	gene3	0	0.156	0.589	0.265	0	gene4	0	0.109	0.276	0.009	0	gene5	0	0.028	0	0.021	0	
	cell1	cell2	cell3	cell4	cell5																																		
gene1	0	0.002	0	0	0																																		
gene2	0.003	0.034	0.052	0.024	0																																		
gene3	0	0.156	0.589	0.265	0																																		
gene4	0	0.109	0.276	0.009	0																																		
gene5	0	0.028	0	0.021	0																																		

Bioinformatics-as-a-Service (BaaS)

- Available to any Stanford faculty or affiliate
- Hourly rate affordable by small labs
- FY 2016: Subsidized by Dean Ann Arvin and Dean Harry Greenberg
- FY 2017: Subsidized by Prof. Mike Snyder
- Contacts for BaaS
Ramesh Nair, Bioinformatics Service Supervisor (rnair@stanford.edu)
Yue Zhang, Bioinformatics Data Scientist (yuezhang@stanford.edu)