

1     **Polar Gini Curve: a Technique to Discover Single-cell Biomarker**  
2                                    **Using 2D Visual Information**

3     Thanh Minh Nguyen<sup>1</sup>, Jacob John Jeevan<sup>1</sup>, Nuo Xu<sup>2</sup>, Jake Chen<sup>1\*</sup>

4     <sup>1</sup>*Informatics Institute, the University of Alabama at Birmingham, AL, United States*

5     <sup>2</sup>*Collat School of Business, the University of Alabama at Birmingham, AL, United States*

6     \*Corresponding author: Jake Chen

7     Email: jakechen@uab.edu

8

9     **Running title:** *Nguyen et al / Polar Gini Curve single cell*

10

11     **Authors' ORCID No**

12     Thanh Nguyen: 0000-0002-8440-1594

13     Jake Chen: 0000-0001-8829-7504

14

15

16     Total word counts: 3446

17     Total figures: 10

18     Total tables: 0

19     Total supplementary figures: 0

20     Total supplementary tables: 0

21     Total supplementary files: 3

22

## 23 **Abstract**

24 In this work, we design the Polar Gini Curve (PGC) technique, which combines the gene  
25 expression and the 2D embedded visual information to detect biomarkers from single-cell data.  
26 Theoretically, a Polar Gini Curve characterizes the shape and ‘evenness’ of cell-point  
27 distribution of cell-point set. To quantify whether a gene could be a marker in a cell cluster, we  
28 can combine two Polar Gini Curves: one drawn upon the cell-points expressing the gene, and the  
29 other drawn upon all cell-points in the cluster. We hypothesize that the closers these two curves  
30 are, the more likely the gene would be cluster markers. We demonstrate the framework in several  
31 simulation case-studies. Applying our framework in analyzing neonatal mouse heart single-cell  
32 data, the detected biomarkers may characterize novel subtypes of cardiac muscle cells. The  
33 source code and data for PGC could be found at  
34 [https://figshare.com/projects/Polar\\_Gini\\_Curve/76749](https://figshare.com/projects/Polar_Gini_Curve/76749).

35  
36 **KEYWORDS:** Single-cell gene expression; Gini coefficient; Polar Gini Curve; Biomarker

## 38 **Introduction**

39 Discovering biomarkers from the single-cell gene expression data is an interesting yet  
40 challenging problem [1]. Compared to the well-established bulk gene expression data, the  
41 expression distribution in single-cell is significantly more heterogeneous [2-4]. Therefore, as  
42 shown in [5, 6], the bulk-analysis strategies [7, 8] achieve low sensitivity in detecting markers. In  
43 addition, as embedding [9-11] and clustering [12-14] are the essential components in many  
44 single-cell expression analytical pipelines [15, 16], the biomarker detection techniques would  
45 need to tackle the challenges and errors from embedding and clustering [17, 18].

46  
47 From the statistical point of view, there are two different directions among the current state-  
48 of-the-art methods in solving the single-cell biomarker discovery problem. The first direction is  
49 using non-parametric approaches [19]. Non-parametric approaches do not attempt to construct  
50 the model characterizing the gene expression distribution [20]. They do not require too many  
51 prior assumptions about the expression data. Therefore, in theory, they could be applied in most  
52 of the heterogeneous scenarios in single-cell expression. For example, Seurat [16] and the  
53 SINCERA [21] pipelines use the Mann–Whitney test [22]. The disadvantages of non-parametric

54 approaches include lacking the point-estimator (for example, we could not tell how much of  
55 fold-change when comparing the expressions of the same gene in two populations) and the lower  
56 true positive rate [5, 6]. On the other hand, the parametric approaches model the underlying  
57 expression distribution. For example, [23] applies Bayesian statistics, Monocle2 [11, 24] and  
58 MAST [2] apply different linear models, and [25] applies the Poisson models to single-cell  
59 differential expression analysis. The parametric approaches, compared to the non-parametric  
60 ones, are significantly more sensitive [5, 6], especially in detecting markers in small cell-cluster  
61 since they may require less number of cell-samples. However, these approaches assume that the  
62 gene expression distribution has specific shapes; therefore, these approaches tend to have higher  
63 false-positive rates.

64  
65 In this work, we developed a new framework based on the novel idea of integrating  
66 expression and the embedded visual information of single-cell data into one metric to identify  
67 biomarkers. This idea has been successfully implemented in spatial single-cell data, in which the  
68 visualization space reflects the relative position of the cells in a tissue image [26, 27]. In this  
69 framework, we decided to take advantage of cluster shape and cell-point distribution from the 2D  
70 visual space. Our strategy was to project the single-cell 2D cluster onto multiple angle-axes to  
71 explore all viewing angles of the cluster. On each ‘viewing angle’, we captured the visual  
72 distribution using the Gini coefficient [28]. Together, for each set of points in 2D, we constructed  
73 a Polar Gini Curve (PGC) from the correspondent between viewing angle and Gini coefficient.  
74 We hypothesized that for the marker gene, its expressing cell set should have its PGC close to  
75 the PGC computed from the whole cluster cell-set. We demonstrated the framework in several  
76 simulation case-studies. Applying our framework in analyzing neonatal mouse heart single-cell  
77 data [29], the detected biomarkers may characterize novel subtypes of cardiac muscle cells. We  
78 named the framework PGC-RSMD (Polar Gini Curve – Root Mean Square Deviation). The  
79 source code and dataset, including supplemental data, used in this manuscript could be found in  
80 [https://figshare.com/projects/Polar\\_Gini\\_Curve/76749](https://figshare.com/projects/Polar_Gini_Curve/76749).

81

82

## 83 Material and Method

### 84 Computing PGC-RSMD for one gene in one cluster

85 **Figure 1** demonstrates the workflow to compute PGC-RSMD for one gene in a cell cluster  
86 from the single-cell expression data. Our approach used the 2D embedding [9] and clustering  
87 results from single-cell expression data as the input. Starting from the 2D  $x$ - $y$  embedding space,  
88 for an arbitrary angle  $\theta$ , the pipeline projects the  $x$ - $y$  coordinate [30] *for every cell-point* onto the  
89  $\theta$ -axis ( $z$  score)

$$z = x\cos(\theta) + y\sin(\theta) \quad (1)$$

90 We subtracted the scores from (1) with the smallest  $z$  to ensure that all  $z$  scores are non-negative,  
91 which is the requirement for computing the Gini coefficient. Then, it computed two Gini  
92 coefficients  $g_{\text{sub}}$  and  $g_{\text{whole}}$  to measure the inequality among the  $z$  scores. The  $g_{\text{sub}}$  coefficient only  
93 used the distribution of  $z$  scores obtained from cells expressing the gene. The  $g_{\text{whole}}$  coefficient  
94 would use the distribution of all  $z$  scores. The Gini coefficient formula is as in [28]

$$g = \frac{1}{2n^2\bar{z}} \sum_{i=1}^n \sum_{j=1}^n |z_i - z_j| \quad (2)$$

95 Here,  $i$  and  $j$  are arbitrary indices in the list of  $z$  scores being used in the computation,  $\bar{z}$  is the  
96 average of these  $z$  scores, and  $n$  is the size of the  $z$  score list. Repeating (1) and (2) for multiple  
97 angles  $\theta$  spanning from 0 to  $2\pi$  would yield the corresponding lists between  $g$  and  $\theta$ , as shown in  
98 the bottom-right table in **Figure 1**. This would lead to two polar curves for Gini coefficients, one  
99 for the cell-points expressing the gene in the cluster, and one for all cell-points in the cluster. We  
100 hypothesize that *the two curves would be closer in the marker-gene scenario than in the non-*  
101 *marker gene scenario*. Therefore, we used the root-mean-square deviation (RSMD) metric,  
102 which is popular in computing fitness in Bioinformatics [31], to determine whether a gene is a  
103 marker in the cluster.

$$RSMD = \frac{\sum_{\forall\theta} (g_{\text{sub}}(\theta) - g_{\text{whole}}(\theta))^2}{n_{\theta}} \quad (3)$$

104 Here,  $n_{\theta}$ , also called resolution, is the number of angles  $\theta$  for which we repeat (1) and (2). In this  
105 work, we chose  $n_{\theta} = 1000$ , which makes the angle list  $\theta = 0, \pi/500, 2\pi/500, \dots, 999\pi/500, 2\pi$ .

106

107 To compute the RSMD statistical p-value for each gene in each cluster, first, we linearly  
108 normalized (scaled) the RSMD computed in (3) such that the normalized RSMD is between 0  
109 and 1. This could be done by dividing (3) by the largest RSMD among all genes in each cluster.  
110 Then, we applied the estimated p-value calculation in [32] to assign a p-value for each gene in  
111 each cluster. Briefly, from the RSMD scores in (3), we verified that the RSMD scores followed a  
112 bell-shaped distribution. Then, we computed the mean  $\mu$  and standard deviation  $\sigma$  of the  
113 normalized RSMD. Then, the p-value for each gene in the cluster is

$$p - value(i) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^U e^{-\frac{(u-\mu)^2}{2\sigma^2}} du \quad (4)$$

114 In (4),  $U$  stands for the normalized RSMD.

115

116

### 117 **Setting up simulation**

118 In this work, to demonstrate how the PGC-RSMD functions, we setup two simulations. In the  
119 first simulation, the cell cluster in the  $x$ - $y$  embedding space had 5000 points, which were  
120 uniformly generated in the unit circle  $x^2 + y^2 \leq 1$ . In the second simulation, the 2D visualization  
121 of cell clusters had the shape identical to the real-world cluster obtained from visualizing the  
122 mouse fetal lung single-cell data [29] using tSNE [33]. We applied the sampling-by-rejection  
123 technique [34] to generate these cluster points as follows. In the first simulation, we randomly  
124 generated a point whose coordinates are between -1 and 1 using uniform sampling, then accepted  
125 the point if it had  $x^2 + y^2 \leq 1$ . In the second simulation, the random point coordinates were within  
126 the cluster coordinate range. We extracted the cluster boundary points, compute the polygon  
127 from these boundary points, which allowed deciding whether a point was inside the polygon  
128 using Matlab [35, 36]. In each simulation, we randomly chose  $m$  percentage of points ( $m = 5, 10,$   
129  $15, \dots, 95$ ) and assumed that they represent the cells expressing gene. For each  $m$  percentage, we  
130 repeat the simulation 1000 times.

131 In addition, to evaluate how the performance of PGC-RSMD would change in drop-out  
132 scenario, we modified the single-cell data simulator in [37] as follow. First, we use [37] default  
133 parameters to synthesize 2 clusters such that each cluster has 6000 cells, 250 markers (total 500  
134 cluster markers) and other 4500 genes. For each cluster marker, the average expression fol-  
135 low change when comparing two clusters was between 4 and 1000. We assigned the drop-out

136 probability for each gene from 0, 5, 10, ..., to 45% such that there were 25 markers for each  
137 drop-out probability. Then, in each cell, we randomly change the marker expression to 0  
138 according to the markers' drop-out probability. For each of the 4500 non-cluster markers, the  
139 expression in each cell was randomly between 0 and 500. We assigned the sparsity – defined as  
140 the percentage of non-expressing cells (0 expression) – for each non-cluster marker from 0, 5,  
141 10, .... to 95%. In each cell, we randomly changed the non-cluster marker to 0 according to its  
142 sparsity. We used the AUC metric to evaluate whether the PGC-RSMD score could differentiate  
143 the 500 cluster-markers: whether each marker is specific for the first or the second clusters.

144

145

#### 146 **Identifying cardiac muscle cell clusters and marker genes from the neonatal mouse** 147 **heart single-cell data**

148 We obtained the neonatal mouse heart single-cell case-study from the Mouse Cell Atlas [29].  
149 We processed the data as specified in [29]. After preprocessing, the dataset covered 19,494 genes  
150 expression in 5075 cells. We use tSNE [33] (without dimensional reduction) to embed the  
151 dataset into the 2D space. We used the density-based clustering algorithm [38] implemented in  
152 Matlab [39] to identify 9 cell clusters. In the implementation [39], we chose the clustering  
153 parameters  $\epsilon = 4$ ,  $\text{minpts} = 40$ . There were 788, 397, 2966, 156, 288, 123, 76, 125, 87 cell-  
154 points in cluster 1, 2, ..., 9, correspondingly. There were 69 cell-points for which the algorithm  
155 is unable to assign to any clusters (Supplemental Data 3).

156

157 We computed the percentage of expressing cells (the naïve approach) and PGC-RSMD for  
158 all genes in all clusters. We removed genes expressing in less than 10% of the cluster cells. For  
159 comparison, in the naïve approach, in each cluster, we selected the top genes sorted by the  
160 highest percentage of expressing cells as the cluster markers. In the PGC-RSMD approach, we  
161 selected the smallest-RSMD genes with its p-value  $< 0.05$  as the cluster marker. In this work, we  
162 focused on identifying the heart muscle cell clusters and their markers. We manually examine  
163 the distribution of cells expressing the well-known heart muscle cell markers: *Myh7*, *Actc1*, and  
164 *Tnnt2* [40-47].

165

166

## 167        **Setting up the re-identifying cluster ID problem**

168        To compare the robustness of our PGC-RSMD markers with other approaches, we setup the  
169 re-identifying cluster ID as follow. From the visual coordinates and 9 clusters of 5075 cells in  
170 [29], we randomly divided the dataset into the training set (4060 cells – 80%) and the test set  
171 (1015 cells – 20%) such that set has samples of all 9 clusters. Using the training set and markers’  
172 expression found by PGC-RSMD, in comparison with other approaches, we applied the neural  
173 network algorithm [48] to train models that identify cluster ID. We evaluated these models in the  
174 test set and recorded the classification accuracy and area-under-receiving characteristic curve  
175 (AUC). Here, we hypothesized that the ‘better’ markers would yield higher classification  
176 accuracy and AUC. The other approaches being compared with PGC-RSMD are:

177        - The baseline approach: in this approach, we would train the classification models using all  
178 genes expression.

179        - The differential expression approach: in this approach, we use Fisher’s exact test [49],  
180 which computes the likelihood of a gene being expressed (raw expression  $> 0$ ) in a cluster,  
181 compare to the likelihood of the gene being expressed outside the cluster. In this work, we select  
182 the DEG markers in each cluster according to the following criteria: odd ratio  $> 5$  and the  
183 percentage of expressing cell ( $m$ )  $> 50\%$ .

184        - The SpatialDE [26] approach: SpatialDE finds the gene with high variance regarding the  
185 distribution of ‘point’ on the spatial 2D space. The ‘null’ hypothesis in this approach is the gene  
186 distribution in the ‘spatial space’ follows a multivariate normal distribution. The marker is  
187 selected if the gene expression distribution is significantly different from the null distribution,  
188 recorded in the p-value. In this work, we select the SpatialDE marker according to the following  
189 criteria: q-value (adjusted p-value)  $< 0.05$  and percentage of expressing cell ( $m$ )  $> 50\%$ .

190        In both the DEG and the SpatialDE approach, we sort the markers according to the decreasing  
191 order of  $m$ . To make a fair comparison, we use the same number of markers, ranging from 5 to  
192 100, found by PGC-RSMD, DEG and SpatialDE to train the classification models.

193

194



## 195 Results

### 196 PGC-RSMD strongly correlates to the percentage of expressing cell in a cluster

197 In **Figure 2**, we show that the fitness between the cluster PGC and the sub-cluster PGC  
198 strongly correlates to the percentage of expressing cell-points as the ‘sub-cluster’  $m$  in the circle-  
199 shaped simulation. In addition, as  $m$  increases, the RSMD variance decreases. We represented  
200 the fitness by the root-mean-square deviation (RSMD) as showed in the method section. In this  
201 figure, for each  $m$  (from 5 to 95), we randomly generate 1000 sub-clusters and their PGCs. The  
202 detailed result of this simulation could be found at the Supplemental Data 1.

203 In addition, we observed a similar correlation when experimenting with the mouse fetal lung  
204 single-cell data [29]. **Figure 3a** shows the dataset clusters visualization using tSNE [33] and the  
205 chosen cluster. To synthesize a 3000-point cluster with the same shape to the chosen cluster, we  
206 still applied the random-by-rejection [34] as presented in the Material and Method section.  
207 **Figure 3b** still shows a strong correlation between  $m$  and RSMD. The detailed result of this  
208 simulation could be found at the Supplemental Data 2.

209  
210 On the other hand, the PGC approach has the potential to answer whether the marker could  
211 identify subpopulations of cells in a cluster. **Figure 4a** demonstrates the 30000-point cluster with  
212 ring-shape  $0.25 \leq x^2 + y^2 \leq 1$ , which appears to be a sub-cluster marker. In this case,  $m = 0.75$ . In  
213 this example, RSMD = 0.033 (**Figure 4b**), which is greater than the RSMD distribution  
214 computed from the random and uniformly-distributed cluster with the same  $m$  (**Figure 4c**).

215  
216 **Figure 5** shows a decrease of PGC-RSMD performance in the drop-out scenario. Briefly, the  
217 synthetic data has 2 clusters, 250 distinct markers for each cluster. Each gene has a specific drop-  
218 out rate as presented in the Material and Method section. Using the PGC-RSMD scores in each  
219 cluster to differentiate these 500 the cluster-specific markers, we observed that PGC-RSMD  
220 achieves very high area-under-receiver-characteristic curve (AUC) ( $>0.95$ ) when the drop-out  
221 probability is small ( $\leq 5\%$ ). However, AUC decreases significantly with the probability of drop-  
222 out (**Figure 5a**). This phenomenon further demonstrates the strong association between RSMD  
223 and the percentage of expressing-cell. When the drop-out rate increases, the percentage of  
224 expressing-cell decreases; therefore, RSMD may mischaracterize a high-dropout marker as non-  
225 marker.



226

227

## 228 **Case-study: PGC identifies heart muscle cell in neonatal mouse heart single-cell**

229 *PGC-RSMD detects markers to support cell-type identification in single-cell mouse*  
230 *neonatal heart data*

231 **Figure 6** summarizes the neonatal mouse heart single-cell data [29] and its 9-cluster markers.  
232 **Figure 6a** visualizes these 9 clusters with tSNE. The PGC-RSMD finds 258 genes, which are  
233 the union of the smallest 100-PCG-RSMD genes found in each cluster, marking these clusters  
234 (Supplemental Data 3). **Figures 6b** and **6c** showed that the gene-cluster marker-association  
235 reflects the underlying gene expression in the single-cell data. In these heatmap figures, each row  
236 corresponds to one gene.

237

238 We identified the muscle-cell clusters 1, 4 and 9 by the expression of *Myh7*, *Actc1*, and  
239 *Tnnt2*, which strongly express in muscle cell type [40-47] (**Figure 7**). Compared to the naïve  
240 method using the percentage of expressing cell, our PGC-RSMD is significantly better by  
241 detecting *Actc1*, which are missed by the naïve approach (**Figure 8**). Furthermore, our approach  
242 identified *Mgrn1* [50, 51], *Ifitm3* [52], *Myl6b* [53] marking cluster 1, which could play important  
243 roles in cardiac muscle functionality, heart failure, and heart development. These genes are not  
244 identified using the naïve approach (**Figure 8**). On the other hand, among genes having a high  
245 percentage of expressing cell, our PGC-RSMD suggests that *Ndufa4l2*, *Mdh2*, and *Atp5g1* may  
246 not be heart muscle cell markers. However, they could suggest a subtype of heart muscle cells  
247 (**Figure 9**). The percentage of expressing cells, PGC-RSMD, statistical p-value and ranks for all  
248 genes could be found in Supplemental Data 3.

249

## 250 *Re-identifying the cells' cluster ID from markers*

251 We observe that the markers found by the PGC-RSMD approach achieve better performance  
252 than the similar SpatialDE [26] markers, and similar performance to the differentially-expressed-  
253 gene (DEG) when being used to re-identify cell's cluster ID. Briefly, after computing the visual  
254 coordinate and cluster ID of all cells, we randomly split the dataset [29] into the training (80%)  
255 and test (20%) sets. We only applied the baseline PGC-RSMD, SpatialDE and DEG approaches  
256 to find the markers and built machine learning models to predict the cells' cluster ID from these

257 markers in the training set. In this experiment, we used all genes to train the predictor in the  
258 baseline approach. The detailed description of this experiment could be found in the method  
259 section. Evaluating the prediction models in the test set, the PGC-RSMD approach performs  
260 closely to the DEG; both have cluster ID prediction accuracy above 0.9 and AUC above 0.95 on  
261 average (**Figure 10**). These two approaches significantly outperform SpatialDE, whose accuracy  
262 is just above the baseline.

263

264

## 265 **Discussion**

266 In this work, we show that integrating the embedded information, which does not often have  
267 a deterministic relationship with gene expression and is primarily for clustering a visualization,  
268 could lead to new insights to biomarkers in single-cell data. In the mouse neonatal heart case-  
269 study, our PGC-RSMD approach could recall *Actc1* as the marker characterizing heart muscle  
270 cell. Meanwhile, the approach using the ratio of expressing cell may fail to recall because a large  
271 percentage of cells does not capture *Actc1* transcript. Therefore, our proposed technique has the  
272 potential to handle analytical issues due to single-cell data quality, such as short-read and low  
273 sequencing depth [54-56]. On the other hand, for genes having high percentage of expressing  
274 cell, the PGC approach could further show that these genes may characterize novel cardiac  
275 muscle cell sub-types for future studies, such as in *Mdh2* and *Myl6b*. Therefore, we suggest that  
276 the biomarker discovery problem could be divided into two sub-problems: the ‘global markers’  
277 specify cell types and the ‘local markers’ specify subtypes. We could solve these two sub-  
278 problems by the right integration of gene expression and visual information.

279

280 In addition to our proposed PGC approach, we could apply several alternative strategies to  
281 integrate the gene expression and visual information to solve the single-cell biomarker discovery  
282 problem. For example, the fractal dimension analysis strategies [57, 58], which focus on  
283 evaluating the uniformity of cell-point distribution, could be applied to identify markers in which  
284 the expressing cells distribute more densely than they are in the overall cluster. In addition, we  
285 could also customize the statistical texture analysis in image processing, such as homogeneity  
286 and integrity [59, 60], to analyze the difference between the overall cluster cell-point and cell-  
287 expressing gene point as the metric to determine markers. On the other hand, choosing the

288 appropriate visual approach depends on the nature of the data and the problem. Our experiment  
289 with the re-identifying cluster ID shows that the well-established SpatialDE [26] does not  
290 outperform our approach and the DEG approach. One explanation is that in our problem, a good  
291 marker for identifying cell type usually follows a good ‘default’ distribution over the visual  
292 space; meanwhile, the SpatialDE aims to find markers that express significantly different from a  
293 default distribution.

294

295 The major limitation of our proposed PGC-RSMD approach is the long computational time,  
296 especially when comparing to the DEG approaches. This is similar to SpatialDE, which also used  
297 visual information to detect marker genes. The DEG approaches may only need to compute one  
298 statistical test to determine whether a gene is a marker for all clusters. Meanwhile, to draw the  
299 curves, PGC-RSMD would need to compute hundreds to thousands, which depends on the  
300 desired curve resolution, to characterize one gene in one cluster. Due to the long computational  
301 time, we were not able to create multiple simulations, which is the ideal approach, run to  
302 compute the statistical [32] p-value for the RSMD score. Therefore, we decided to reapply the  
303 estimation presented to compute the p-value. This approach is computationally more efficient but  
304 may not well-reflect the statistical characteristic of the single-cell data. In addition, we have not  
305 fully tackled the problem of choosing the right threshold to determine whether a gene expresses  
306 in a cell. Because of the strong association between PGC-RSMD and the percentage of  
307 expressing-cell, we expect that the result would significantly different when choosing a different  
308 threshold to determine whether a gene expresses in a cell. In this work, choosing 0 as the  
309 threshold still yields good performance because of the high sparsity in the real dataset.

310

## 311 **Conclusions**

312 In this work, we have presented Polar Gini Curve, a novel technique to detect markers from  
313 the single-cell RNA expression data using visual information. In principle, our technique could  
314 complement the state-of-the-art approach: the PGC technique finds markers such that the  
315 expressing cells are evenly distributed throughout the cluster space; meanwhile, the state-of-the-  
316 art approach finds markers assuming a multivariate normal distribution of gene expression in the  
317 visual space. We have demonstrated that the PGC technique performs better in some tasks in  
318 single-cell analysis.

319

320

321

322

323

## 324 **Authors' contribution**

325 TN designed and implemented the core polar Gini algorithm, designed the sampling  
326 strategies in the simulation, performed the neonatal heart single-cell case-study, and primarily  
327 prepared the manuscript. JJ prepared the software package, preprocessed the single-cell data, and  
328 executed the simulation designs. NX designed different simulation scenarios, interpreted the  
329 statistical outcomes, and prepared the literature review. JC originated the idea of using Gini  
330 coefficient curves to integrate gene expression, cell-point distribution, and cluster shape to solve  
331 the biomarker discovery problem, designed the performance evaluation, and supervised the  
332 overall technical development. All authors reviewed, revised and approved the manuscript.

333

334

## 335 **Competing interests**

336 The authors have declared that no competing interests exist.

337

338

## 339 **Acknowledgments**

340 This work was supported by the 'startup budget' granted to Jake Chen from the University of  
341 Alabama at Birmingham.

342

343

## 344 **References**

- 345 1. Zhu, Z., et al., Single-cell transcriptome in the identification of disease biomarkers:  
346 opportunities and challenges. *J Transl Med*, 2014. 12: p. 212.

- 347 2. Finak, G., et al., MAST: a flexible statistical framework for assessing transcriptional  
348 changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome*  
349 *Biol*, 2015. 16: p. 278.
- 350 3. McCarthy, D.J., et al., Scater: pre-processing, quality control, normalization and  
351 visualization of single-cell RNA-seq data in R. *Bioinformatics*, 2017. 33(8): p. 1179-  
352 1186.
- 353 4. Poirion, O.B., et al., Single-Cell Transcriptomics Bioinformatics and Computational  
354 Challenges. *Front Genet*, 2016. 7: p. 163.
- 355 5. Jaakkola, M.K., et al., Comparison of methods to detect differentially expressed genes  
356 between single-cell populations. *Brief Bioinform*, 2017. 18(5): p. 735-743.
- 357 6. Wang, T., et al., Comparative analysis of differential gene expression analysis tools for  
358 single-cell RNA sequencing data. *BMC Bioinformatics*, 2019. 20(1): p. 40.
- 359 7. Love, M.I., W. Huber, and S. Anders, Moderated estimation of fold change and  
360 dispersion for RNA-seq data with DESeq2. *Genome Biol*, 2014. 15(12): p. 550.
- 361 8. Robinson, M.D., D.J. McCarthy, and G.K. Smyth, edgeR: a Bioconductor package for  
362 differential expression analysis of digital gene expression data. *Bioinformatics*, 2010.  
363 26(1): p. 139-40.
- 364 9. Pezzotti, N., et al., Approximated and User Steerable tSNE for Progressive Visual  
365 Analytics. *IEEE Trans Vis Comput Graph*, 2017. 23(7): p. 1739-1752.
- 366 10. Becht, E., et al., Dimensionality reduction for visualizing single-cell data using UMAP.  
367 *Nat Biotechnol*, 2018.
- 368 11. Qiu, X., et al., Reversed graph embedding resolves complex single-cell trajectories. *Nat*  
369 *Methods*, 2017. 14(10): p. 979-982.
- 370 12. Yang, Y., et al., SAFE-clustering: Single-cell Aggregated (from Ensemble) clustering for  
371 single-cell RNA-seq data. *Bioinformatics*, 2019. 35(8): p. 1269-1277.
- 372 13. Kiselev, V.Y., et al., SC3: consensus clustering of single-cell RNA-seq data. *Nat*  
373 *Methods*, 2017. 14(5): p. 483-486.
- 374 14. Aibar, S., et al., SCENIC: single-cell regulatory network inference and clustering. *Nat*  
375 *Methods*, 2017. 14(11): p. 1083-1086.
- 376 15. Zheng, G.X., et al., Massively parallel digital transcriptional profiling of single cells. *Nat*  
377 *Commun*, 2017. 8: p. 14049.

- 378 16. Satija, R., et al., Spatial reconstruction of single-cell gene expression data. Nat  
379 Biotechnol, 2015. 33(5): p. 495-502.
- 380 17. Kiselev, V.Y., T.S. Andrews, and M. Hemberg, Challenges in unsupervised clustering of  
381 single-cell RNA-seq data. Nat Rev Genet, 2019. 20(5): p. 273-282.
- 382 18. Yuan, G.C., et al., Challenges and emerging directions in single-cell analysis. Genome  
383 Biol, 2017. 18(1): p. 84.
- 384 19. Conover, W.J. and W.J. Conover, Practical nonparametric statistics. 1980.
- 385 20. Hollander, M., D.A. Wolfe, and E. Chicken, Nonparametric statistical methods. Vol. 751.  
386 2013: John Wiley & Sons.
- 387 21. Guo, M., et al., SINCERA: A Pipeline for Single-Cell RNA-Seq Profiling Analysis.  
388 PLoS Comput Biol, 2015. 11(11): p. e1004575.
- 389 22. Birnbaum, Z. On a use of the Mann-Whitney statistic. in Proceedings of the Third  
390 Berkeley Symposium on Mathematical Statistics and Probability, Volume 1:  
391 Contributions to the Theory of Statistics. 1956. The Regents of the University of  
392 California.
- 393 23. Korthauer, K.D., et al., A statistical approach for identifying differential distributions in  
394 single-cell RNA-seq experiments. Genome Biol, 2016. 17(1): p. 222.
- 395 24. Trapnell, C., et al., The dynamics and regulators of cell fate decisions are revealed by  
396 pseudotemporal ordering of single cells. Nat Biotechnol, 2014. 32(4): p. 381-386.
- 397 25. Kharchenko, P.V., L. Silberstein, and D.T. Scadden, Bayesian approach to single-cell  
398 differential expression analysis. Nat Methods, 2014. 11(7): p. 740-2.
- 399 26. Svensson, V., S.A. Teichmann, and O. Stegle, SpatialDE: identification of spatially  
400 variable genes. Nat Methods, 2018. 15(5): p. 343-346.
- 401 27. Edsgard, D., P. Johnsson, and R. Sandberg, Identification of spatial expression trends in  
402 single-cell gene expression data. Nat Methods, 2018. 15(5): p. 339-342.
- 403 28. Gini, C., Concentration and dependency ratios. Rivista di politica economica, 1997. 87:  
404 p. 769-792.
- 405 29. Han, X., et al., Mapping the Mouse Cell Atlas by Microwell-Seq. Cell, 2018. 173(5): p.  
406 1307.
- 407 30. Strang, G., et al., Introduction to linear algebra. Vol. 3. 1993: Wellesley-Cambridge Press  
408 Wellesley, MA.

- 409 31. Hyndman, R.J. and A.B. Koehler, Another look at measures of forecast accuracy.  
410 International journal of forecasting, 2006. 22(4): p. 679-688.
- 411 32. Yue, Z., et al., WIPER: Weighted in-Path Edge Ranking for biomolecular association  
412 networks. Quantitative Biology, 2019. 7(4): p. 313-326.
- 413 33. Maaten, L.v.d. and G. Hinton, Visualizing data using t-SNE. Journal of machine learning  
414 research, 2008. 9(Nov): p. 2579-2605.
- 415 34. Bishop, C.M., Pattern recognition and machine learning. 2006: springer.
- 416 35. MathWorks. Matlab - inpolygon. 2019 2019/06/05]; Available from:  
417 <https://www.mathworks.com/help/matlab/ref/inpolygon.html>.
- 418 36. MathWorks. Matlab - boundary. 2019 2019/06/05]; Available from:  
419 <https://www.mathworks.com/help/matlab/ref/boundary.html>.
- 420 37. Baruzzo, G., I. Patuzzi, and B. Di Camillo, SPARSim Single Cell: a count data simulator  
421 for scRNA-seq data. Bioinformatics, 2019.
- 422 38. Ester, M., et al. A density-based algorithm for discovering clusters in large spatial  
423 databases with noise. in Kdd. 1996.
- 424 39. MathWorks. Matlab - dbscan. 2019; Available from:  
425 <https://www.mathworks.com/help/stats/dbscan.html>.
- 426 40. Bashyam, M.D., et al., Molecular genetics of familial hypertrophic cardiomyopathy  
427 (FHC). J Hum Genet, 2003. 48(2): p. 55-64.
- 428 41. Finsterer, J., C. Stollberger, and J.A. Towbin, Left ventricular noncompaction  
429 cardiomyopathy: cardiac, neuromuscular, and genetic factors. Nat Rev Cardiol, 2017.  
430 14(4): p. 224-237.
- 431 42. Keren, A., P. Syrris, and W.J. McKenna, Hypertrophic cardiomyopathy: the genetic  
432 determinants of clinical disease expression. Nat Clin Pract Cardiovasc Med, 2008. 5(3):  
433 p. 158-68.
- 434 43. Morita, H., et al., Shared genetic causes of cardiac hypertrophy in children and adults. N  
435 Engl J Med, 2008. 358(18): p. 1899-908.
- 436 44. Jiang, H.K., et al., Reduced ACTC1 expression might play a role in the onset of  
437 congenital heart disease by inducing cardiomyocyte apoptosis. Circ J, 2010. 74(11): p.  
438 2410-8.



- 439 45. Kwon, C., et al., A regulatory pathway involving Notch1/beta-catenin/Isl1 determines  
440 cardiac progenitor cell fate. *Nat Cell Biol*, 2009. 11(8): p. 951-7.
- 441 46. Wei, B. and J.P. Jin, TNNT1, TNNT2, and TNNT3: Isoform genes, regulation, and  
442 structure-function relationships. *Gene*, 2016. 582(1): p. 1-13.
- 443 47. Ju, Y., et al., Troponin T3 expression in skeletal and smooth muscle is required for  
444 growth and postnatal survival: characterization of Tnnt3(tm2a(KOMP)Wtsi) mice.  
445 *Genesis*, 2013. 51(9): p. 667-75.
- 446 48. Russell, S. and P. Norvig, *Artificial Intelligence: A Modern Approach*. 2003. Prentice  
447 Hall, Upper Saddle River, New Jersey.
- 448 49. Mehta, C.R. and N.R. Patel, A network algorithm for performing Fisher's exact test in  $r \times$   
449  $c$  contingency tables. *Journal of the American Statistical Association*, 1983. 78(382): p.  
450 427-434.
- 451 50. Mukherjee, R. and O. Chakrabarti, Regulation of Mitofusin1 by Mahogunin Ring Finger-  
452 1 and the proteasome modulates mitochondrial fusion. *Biochim Biophys Acta*, 2016.  
453 1863(12): p. 3065-3083.
- 454 51. Liu, X., et al., Differential microRNA Expression and Regulation in the Rat Model of  
455 Post-Infarction Heart Failure. *PLoS One*, 2016. 11(8): p. e0160920.
- 456 52. Lau, S.L., et al., Interferons induce the expression of IFITM1 and IFITM3 and suppress  
457 the proliferation of rat neonatal cardiomyocytes. *J Cell Biochem*, 2012. 113(3): p. 841-7.
- 458 53. Wang, L., et al., Mutations in myosin light chain kinase cause familial aortic dissections.  
459 *Am J Hum Genet*, 2010. 87(5): p. 701-7.
- 460 54. Shalek, A.K., et al., Single-cell RNA-seq reveals dynamic paracrine control of cellular  
461 variation. *Nature*, 2014. 510(7505): p. 363-9.
- 462 55. Rizzetto, S., et al., Impact of sequencing depth and read length on single cell RNA  
463 sequencing data of T cells. *Sci Rep*, 2017. 7(1): p. 12781.
- 464 56. McDavid, A., et al., Data exploration, quality control and testing in single-cell qPCR-  
465 based gene expression experiments. *Bioinformatics*, 2013. 29(4): p. 461-7.
- 466 57. Fortin, C., et al., Fractal dimension in the analysis of medical images. *IEEE Engineering*  
467 *in Medicine and Biology Magazine*, 1992. 11(2): p. 65-71.
- 468 58. Davies, S. and P. Hall, Fractal analysis of surface roughness by using spatial data. *Journal*  
469 *of the Royal Statistical Society: Series B (Statistical Methodology)*, 1999. 61(1): p. 3-37.

- 470 59. Bharati, M.H., J.J. Liu, and J.F. MacGregor, Image texture analysis: methods and  
471 comparisons. *Chemometrics and intelligent laboratory systems*, 2004. 72(1): p. 57-71.
- 472 60. Kunitatsu, A., et al., Comparison between glioblastoma and primary central nervous  
473 system lymphoma using MR image-based texture analysis. *Magnetic Resonance in*  
474 *Medical Sciences*, 2017: p. mp. 2017-0044.

475  
476

## 477 **Figure legends**

478 Figure 1. Overall workflow to compute the PGC-RSMD metric for one gene in one cluster of  
479 cells. Here, the data points, histogram, and PGC for cells expressing the gene are cyan. The ones  
480 for the whole cells in the cluster are red.

481

482 Figure 2. Boxplot showing a strong correlation between ‘subcluster’ percentage ( $m$ ) and cluster-  
483 subcluster PGC fitness (RSMD) in uniformly-distributed and a circular cluster.

484

485 Figure 3. a) The selected cluster for the experiment in [29]. b) Correlation between ‘subcluster’  
486 percentage ( $m$ ) and cluster-subcluster PGC fitness (RSMD) in the selected cluster.

487

488 Figure 4. The ring-shape simulation study: a) Visualization of the cluster and ring-shape sub-  
489 cluster ( $m = 0.75$ ); b) PGC yield RSMD = 0.033; c) Distribution of RSMD, extracted from  
490 Figure 2 with  $m = 75\%$ , when the sub-cluster uniformly distributed on the cluster area.

491

492 Figure 5. PGC-RSMD performance in recalling cluster marker in drop-out simulation. a)  
493 heatmap showing the simulation design of 500 markers and 4500 neutral genes, with drop out /  
494 percentage of cell expressing between 5 and 100%; b) The simulation data 2D visualization; c)  
495 the AUC drops when drop-out increases.

496

497 Figure 6. The result from mouse neonatal heart single-cell [29] analysis. a) the tSNE plot shows  
498 9 clusters. b) gene-cluster marker relationship (from 258 genes) found by PGC-RSMD; ■ gene  
499 is found as marker, ■ gene is found as non-marker. c) expression heatmap for these genes.

500

501 Figure 7: Heart muscle cell clusters, identified by *Myh7*, *Actc1*, and *Tnnt2*

502

503 Figure 8. PGC-RSMD highlight makers that do not have high percentage of expressing cells:

504 PGCs of *Actc1*, *Mgrn1*, *Ifitm3*, *Myl6b* in cluster 1. The numbers in the parenthesis are ranks of

505 these genes in each metric

506

507 Figure 9. PGC-RSMD shows that gene has high percentage of expressing cells: *Ndufa4l2*,

508 *Mdh2*, and *Atp5g1*, may not be markers in cluster 1. These genes appear to highlight a local

509 subcluster. The numbers in the parenthesis are ranks of these genes in each metric.

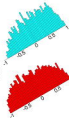
510

511 Figure 10. Performance of the PGC-RSMD, SpatialDE, and DEG in re-identifying the cell's

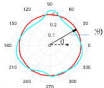
512 cluster ID problem using dataset [26]. The x-axis shows the number of top-significant markers

513 being selected to train the prediction models. a) accuracy; b) AUC over 9 clusters.

2D visualization



$$Gini(g) = \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|$$

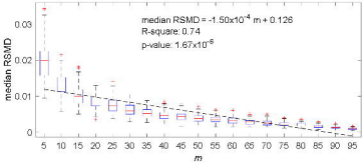


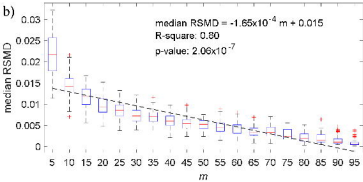
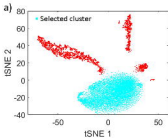
| Angle $\theta$ | $S_{sub}$ | $S_{whole}$ | $f_{\theta} =  S_{sub} - S_{whole} ^2$ |
|----------------|-----------|-------------|--|
| ...            | ...       | ...         | ...                                    |
| 30°            | 0.26      | 0.29        | 0.03                                   |
| ...            | ...       | ...         | ...                                    |

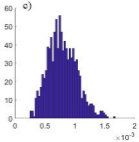
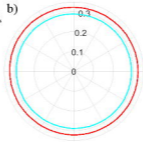
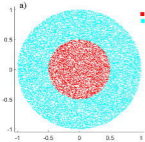
$$RSMD = \sum f_{\theta} / n_{\theta}$$

■ Cell expressing genes (sub cluster)

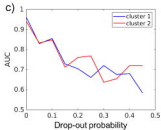
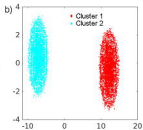
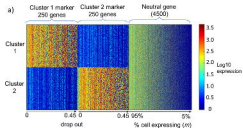
■ All cells (whole cluster)

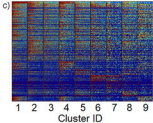
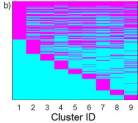
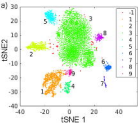


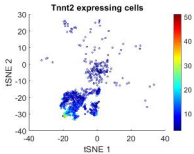
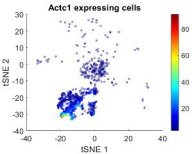
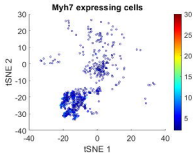
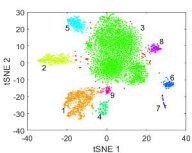


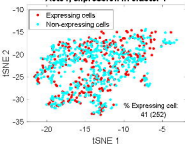
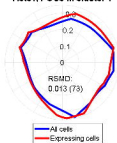
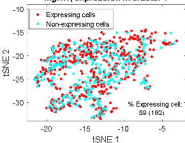
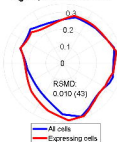
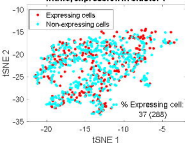
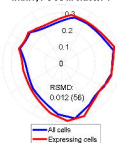
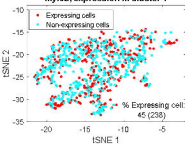
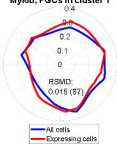


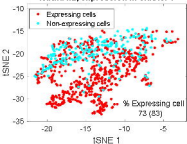
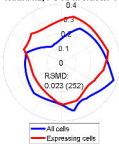
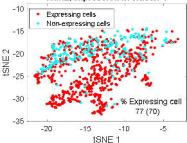
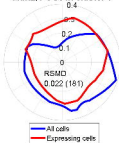
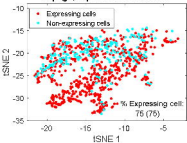








**Actc1, expression in cluster 1****Actc1, PGCs in cluster 1****Mgrn1, expression in cluster 1****Mgrn1, PGCs in cluster 1****Ifitm3, expression in cluster 1****Ifitm3, PGCs in cluster 1****Myli6b, expression in cluster 1****Myli6b, PGCs in cluster 1**

**Ndufa4l2, expression in cluster 1****Ndufa4l2, PGCs in cluster 1****Mdh2, expression in cluster 1****Mdh2, PGCs in cluster 1****Atp5g1, expression in cluster 1****Atp5g1, PGCs in cluster 1**