1    # Polar Gini Curve: a Technique to Discover Single-cell Biomarker

2    # Using 2D Visual Information

3    Thanh Minh Nguyen[1], Jacob John Jeevan[1], Nuo Xu[2], Jake Chen[1*]

4    *[1]Informatics Institute, the University of Alabama at Birmingham, AL, United States*

5    *[2]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    Jacob John Jeevan: 0000-0003-0910-5610

14    Jake Chen: 0000-0001-8829-7504

15

16

17    Total word counts: 3446

18    Total figures: 10

19    Total tables: 0

20    Total supplementary figures: 0

21    Total supplementary tables: 0

22    Total supplementary files: 3

23

## Abstract

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

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

## Introduction

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

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

55  comparing the expressions of the same gene in two populations) and the lower true positive rate

56  [5, 6]. On the other hand, the parametric approaches model the underlying expression distribution.

57  For example, [23] applies Bayesian statistics, Monocle2 [11, 24] and MAST [2] apply different

58  linear models, and [25] applies the Poisson models to single-cell differential expression analysis.

59  The parametric approaches, compared to the non-parametric ones, are significantly more sensitive

60  [5, 6], especially in detecting markers in small cell-cluster since they may require less number of

61  cell-samples. However, these approaches assume that the gene expression distribution has specific

62  shapes; therefore, these approaches tend to have higher false-positive rates.

63

64      In this work, we developed a new framework based on the novel idea of integrating expression

65  and the embedded visual information of single-cell data into one metric to identify biomarkers.

66  This idea has been successfully implemented in spatial single-cell data, in which the visualization

67  space reflects the relative position of the cells in a tissue image [26, 27]. In this framework, we

68  decided to take advantage of cluster shape and cell-point distribution from the 2D visual space.

69  Our strategy was to project the single-cell 2D cluster onto multiple angle-axes to explore all

70  viewing angles of the cluster. On each 'viewing angle', we captured the visual distribution using

71  the Gini coefficient [28]. Together, for each set of points in 2D, we constructed a Polar Gini Curve

72  (PGC) from the correspondent between viewing angle and Gini coefficient. We hypothesized that

73  for the marker gene, its expressing cell set should have its PGC close to the PGC computed from

74  the whole cluster cell-set. We demonstrated the framework in several simulation case-studies.

75  Applying our framework in analyzing neonatal mouse heart single-cell data [29], the detected

76  biomarkers may characterize novel subtypes of cardiac muscle cells. We named the framework

77  PGC-RSMD (Polar Gini Curve – Root Mean Square Deviation). The source code and dataset,

78  including    supplemental    data,    used    in    this    manuscript    could    be    found    in

79  https://figshare.com/projects/Polar_Gini_Curve/76749.

80

81

82  **Material and Method**

83      **Computing PGC-RSMD for one gene in one cluster**

84      **Figure 1** demonstrates the workflow to compute PGC-RSMD for one gene in a cell cluster

85  from the single-cell expression data. Our approach used the 2D embedding [9] and clustering

86    results from single-cell expression data as the input. Starting from the 2D *x-y* embedding space,

87    for an arbitrary angle θ, the pipeline projects the x-y coordinate [30] *for every cell-point* onto the

88    θ-axis (*z* score)

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

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 $g_{sub}$ and $g_{whole}$ to measure the inequality among the *z* scores. The $g_{sub}$ coefficient only

93    used the distribution of *z* scores obtained from cells expressing the gene. The $g_{whole}$ coefficient

94    would use the distribution of all *z* scores. The Gini coefficient formula is as in [28]

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

96    Here, *i* and *j* are arbitrary indices in the list of *z* scores being used in the computation, $\bar{z}$ is the

97    average of these *z* scores, and *n* is the size of the *z* score list. Repeating (1) and (2) for multiple

98    angles θ spanning from 0 to 2π would yield the corresponding lists between *g* and θ, as shown in

99    the bottom-right table in **Figure 1**. This would lead to two polar curves for Gini coefficients, one

100   for the cell-points expressing the gene in the cluster, and one for all cell-points in the cluster. We

101   hypothesize that *the two curves would be closer in the marker-gene scenario than in the non-*

102   *marker gene scenario*. Therefore, we used the root-mean-square deviation (RSMD) metric, which

103   is popular in computing fitness in Bioinformatics [31], to determine whether a gene is a marker in

104   the cluster.

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

106   Here, $n_\theta$, also called resolution, is the number of angles θ for which we repeat (1) and (2). In this

107   work, we chose $n_\theta = 1000$, which makes the angle list θ = 0, π/500, 2π/500, …, 999π/500, 2π.

108

109       To compute the RSMD statistical p-value for each gene in each cluster, first, we linearly

110   normalized (scaled) the RSMD computed in (3) such that the normalized RSMD is between 0 and

111   1. This could be done by diving (3) by the largest RSMD among all genes in each cluster. Then,

112   we applied the estimated p-value calculation in [32] to assign a p-value for each gene in each

113   cluster. Briefly, from the RSMD scores in (3), we verified that the RSMD scores followed a bell-

114     shaped distribution. Then, we computed the mean $\mu$ and standard deviation $\sigma$ of the normalized

115     RSMD. Then, the p-value for each gene in the cluster is

116    
$$p - value(i) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{U} e^{-\frac{(u-\mu)^2}{2\sigma^2}} du \quad (4)$$

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

118

119

120     **Setting up simulation**

121     In this work, to demonstrate how the PGC-RSMD functions, we setup two simulations. In the

122     first simulation, the cell cluster in the $x$-$y$ embedding space had 5000 points, which were uniformly

123     generated in the unit circle $x^2 + y^2 \leq 1$. In the second simulation, the 2D visualization of cell clusters

124     had the shape identical to the real-world cluster obtained from visualizing the mouse fetal lung

125     single-cell data [29] using tSNE [33]. We applied the sampling-by-rejection technique [34] to

126     generate these cluster points as follows. In the first simulation, we randomly generated a point

127     whose coordinates are between -1 and 1 using uniform sampling, then accepted the point if it had

128     $x^2 + y^2 \leq 1$. In the second simulation, the random point coordinates were within the cluster

129     coordinate range. We extracted the cluster boundary points, compute the polygon from these

130     boundary points, which allowed deciding whether a point was inside the polygon using Matlab

131     [35, 36]. In each simulation, we randomly chose $m$ percentage of points ($m$ = 5, 10, 15, …, 95) and

132     assumed that they represent the cells expressing gene. For each $m$ percentage, we repeat the

133     simulation 1000 times.

134     In addition, to evaluate how the performance of PGC-RSMD would change in drop-out

135     scenario, we modified the single-cell data simulator in [37] as follow. First, we use [37] default

136     parameters to synthesize 2 clusters such that each cluster has 6000 cells, 250 markers (total 500

137     cluster markers) and other 4500 genes. For each cluster marker, the average expression fol-change

138     when comparing two clusters was between 4 and 1000. We assigned the drop-out probability for

139     each gene from 0, 5, 10, …, to 45% such that there were 25 markers for each drop-out probability.

140     Then, in each cell, we randomly change the marker expression to 0 according to the markers' drop-

141     out probability. For each of the 4500 non-cluster markers, the expression in each cell was randomly

142     between 0 and 500. We assigned the sparisity – defined as the percentage of non-expressing cells

143     (0 expression) – for each non-cluster marker from 0, 5, 10, …. to 95%. In each cell, we randomly

144  changed the non-cluster marker to 0 according to its sparsity. We used the AUC metric to evaluate

145  whether the PGC-RSMD score could differentiate the 500 cluster-markers: whether each marker

146  is specific for the first or the second clusters.

147

148

149  **Identifying cardiac muscle cell clusters and marker genes from the neonatal mouse heart**

150  **single-cell data**

151  We obtained the neonatal mouse heart single-cell case-study from the Mouse Cell Atlas [29].

152  We processed the data as specified in [29]. After preprocessing, the dataset covered 19,494 genes

153  expression in 5075 cells. We use tSNE [33] (without dimensional reduction) to embed the dataset

154  into the 2D space. We used the density-based clustering algorithm [38] implemented in Matlab

155  [39] to identify 9 cell clusters. In the implementation [39], we chose the clustering parameters

156  epsilon = 4, minpts = 40. There were 788, 397, 2966, 156, 288, 123, 76, 125, 87 cell-points in

157  cluster 1, 2, …, 9, correspondingly. There were 69 cell-points for which the algorithm is unable to

158  assign to any clusters (Supplemental Data 3).

159

160  We computed the percentage of expressing cells (the naïve approach) and PGC-RSMD for all

161  genes in all clusters. We removed genes expressing in less than 10% of the cluster cells. For

162  comparison, in the naïve approach, in each cluster, we selected the top genes sorted by the highest

163  percentage of expressing cells as the cluster markers. In the PGC-RSMD approach, we selected

164  the smallest-RSMD genes with its p-value < 0.05 as the cluster marker. In this work, we focused

165  on identifying the heart muscle cell clusters and their markers. We manually examine the

166  distribution of cells expressing the well-known heart muscle cell markers: *Myh7*, *Actc1*, and *Tnnt2*

167  [40-47].

168

169

170  **Seting up the re-identifying cluster ID problem**

171  To compare the robustness of our PGC-RSMD markers with other approaches, we setup the

172  re-identifying cluster ID as follow. From the visual coordinates and 9 clusters of 5075 cells in [29],

173  we randomly divided the dataset into the training set (4060 cells – 80%) and the test set (1015 cells

174  – 20%) such that set has samples of all 9 clusters. Using the training set and markers' expression

175 found by PGC-RSMD, in comparison with other approaches, we applied the neural network

176 algorithm [48] to train models that identify cluster ID. We evaluated these models in the test set

177 and recorded the classification accuracy and area-under-receiving characteristic curve (AUC).

178 Here, we hypothesized that the 'better' markers would yield higher classification accuracy and

179 AUC. The other approaches being compared with PGC-RSMD are:

180 - The baseline approach: in this approach, we would train the classification models using all

181 genes expression.

182 - The differential expression approach: in this approach, we use Fisher's exact test [49], which

183 computes the likelihood of a gene being expressed (raw expression > 0) in a cluster, compare to

184 the likelihood of the gene being expressed outside the cluster. In this work, we select the DEG

185 markers in each cluster according to the following criteria: odd ratio > 5 and the percentage of

186 expressing cell ($m$) > 50%.

187 - The SpatialDE [26] approach: SpatialDE finds the gene with high variance regarding the

188 distribution of 'point' on the spatial 2D space. The 'null' hypothesis in this approach is the gene

189 distribution in the 'spatial space' follows a multivariate normal distribution. The marker is selected

190 if the gene expression distribution is significantly different from the null distribution, recorded in

191 the p-value. In this work, we select the SpatialDE marker according to the following criteria: q-

192 value (adjusted p-value) < 0.05 and percentage of expressing cell ($m$) > 50%.

193 In both the DEG and the SpatialDE approach, we sort the markers according to the decreasing

194 order of $m$. To make a fair comparison, we use the same number of markers, ranging from 5 to

195 100, found by PGC-RSMD, DEG and SpatialDE to train the classification models.

196

197

198 **Results**

199 **PGC-RSMD strongly correlates to the percentage of expressing cell in a cluster**

200 In **Figure 2**, we show that the fitness between the cluster PGC and the sub-cluster PGC strongly

201 correlates to the percentage of expressing cell-points as the 'sub-cluster' $m$ in the circle-shaped

202 simulation. In addition, as $m$ increases, the RSMD variance decreases. We represented the fitness

203 by the **r**oot-**m**ean-**s**quare **d**eviation (RSMD) as showed in the method section. In this figure, for

204 each $m$ (from 5 to 95), we randomly generate 1000 sub-clusters and their PGCs. The detailed result

205 of this simulation could be found at the Supplemental Data 1.

206       In addition, we observed a similar correlation when experimenting with the mouse fetal lung

207    single-cell data [29]. **Figure 3a** shows the dataset clusters visualization using tSNE [33] and the

208    chosen cluster. To synthesize a 3000-point cluster with the same shape to the chosen cluster, we

209    still applied the random-by-rejection [34] as presented in the Material and Method section. **Figure**

210    **3b** still shows a strong correlation between $m$ and RSMD. The detailed result of this simulation

211    could be found at the Supplemental Data 2.

212

213       On the other hand, the PGC approach has the potential to answer whether the marker could

214    identify subpopulations of cells in a cluster. **Figure 4a** demonstrates the 30000-point cluster with

215    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

216    this example, RSMD = 0.033 (**Figure 4b**), which is greater than the RSMD distribution computed

217    from the random and uniformly-distributed cluster with the same $m$ (**Figure 4c**).

218

219    **Figure 5** shows a decrease of PGC-RSMD performance in the drop-out scenario. Briefly, the

220    synthetic data has 2 clusters, 250 distinct markers for each cluster. Each gene has a specific drop-

221    out rate as presented in the Material and Method section. Using the PGC-RSMD scores in each

222    cluster to differentiate these 500 the cluster-specific markers, we observed that PGC-RSMD

223    achieves very high area-under-receiver-characteristic curve (AUC) (>0.95) when the drop-out

224    probability is small ($\leq 5\%$). However, AUC decreases significantly with the probability of drop-

225    out (**Figure 5a**). This phenomenon further demonstrates the strong association between RSMD

226    and the percentage of expressing-cell. When the drop-out rate increases, the percentage of

227    expressing-cell decreases; therefore, RSMD may mischaracterize a high-dropout marker as non-

228    marker.

229

230

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

232    *PGC-RSMD detects markers to support cell-type identification in single-cell mouse*

233    *neonatal heart data*

234    **Figure 6** summarizes the neonatal mouse heart single-cell data [29] and its 9-cluster markers.

235    **Figure 6a** visualizes these 9 clusters with tSNE. The PGC-RSMD founds 258 genes, which are

236    the union of the smallest 100-PCG-RSMD genes found in each cluster, marking these clusters

237   (Supplemental Data 3). **Figures 6b** and **6c** showed that the gene-cluster marker-association reflects

238   the underlying gene expression in the single-cell data. In these heatmap figures, each row

239   corresponds to one gene.

240

241   We identified the muscle-cell clusters 1, 4 and 9 by the expression of *Myh7*, *Actc1*, and *Tnnt2*,

242   which strongly express in muscle cell type [40-47] (**Figure 7**). Compared to the naïve method

243   using the percentage of expressing cell, our PGC-RSMD is significantly better by detecting *Actc1*,

244   which are missed by the naïve approach (**Figure 8**). Furthermore, our approach identified *Mgrn1*

245   [50, 51], *Ifitm3* [52], *Myl6b* [53] marking cluster 1, which could play important roles in cardiac

246   muscle functionality, heart failure, and heart development. These genes are not identified using

247   the naïve approach (**Figure 8**). On the other hand, among genes having a high percentage of

248   expressing cell, our PGC-RSMD suggests that *Ndufa4l2*, *Mdh2*, and *Atp5g1* may not be heart

249   muscle cell markers. However, they could suggest a subtype of heart muscle cells (**Figure 9**). The

250   percentage of expressing cells, PGC-RSMD, statistical p-value and ranks for all genes could be

251   found in Supplemental Data 3.

252

253       *Re-identifying the cells' cluster ID from markers*

254   We observe that the markers found by the PGC-RSMD approach achieve better performance

255   than the similar SpatialDE [26] markers, and similar performance to the differentially-expressed-

256   gene (DEG) when being used to re-identify cell's cluster ID. Briefly, after computing the visual

257   coordinate and cluster ID of all cells, we randomly split the dataset [29] into the training (80%)

258   and test (20%) sets. We only applied the baseline PGC-RSMD, SpatialDE and DEG approaches

259   to find the markers and built machine learning models to predict the cells' cluster ID from these

260   markers in the training set. In this experiment, we used all genes to train the predictor in the

261   baseline approach. The detailed description of this experiment could be found in the method

262   section. Evaluating the prediction models in the test set, the PGC-RSMD approach performs

263   closely to the DEG; both have cluster ID prediction accuracy above 0.9 and AUC above 0.95 on

264   average (**Figure 10**). These two approaches significantly outperform SpatialDE, whose accuracy

265   is just above the baseline.

266

267

## Discussion

268

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

282

283     In this work, we primarily demonstrate how PGC detects markers for single cluster, which
284     does not need the gene expression from other clusters in the dataset. The approach could be
285     extended to incorporate the 'global' expression as follow. First, a PGC analysis can be performed
286     with marker cells as the foreground and all cells (regardless of their cluster assignments) as the
287     background. Second, a PGC analysis can be performed for each cluster in the dataset
288     independently and compare among the clusters' marker lists. In the neonatal mouse heart case-
289     study, this approach shows two types of marker: one expressing globally in all clusters, which are
290     likely heart-tissue specific; the other express locally in one or some specific cluster, which are
291     likely cell-type specific.

292

293     In addition to our proposed PGC approach, we could apply several alternative strategies to
294     integrate the gene expression and visual information to solve the single-cell biomarker discovery
295     problem. For example, the fractal dimension analysis strategies [57, 58], which focus on evaluating
296     the uniformity of cell-point distribution, could be applied to identify markers in which the
297     expressing cells distribute more densely than they are in the overall cluster. In addition, we could
298     also customize the statistical texture analysis in image processing, such as homogeneity and

299   integrity [59, 60], to analyze the difference between the overall cluster cell-point and cell-

300   expressing gene point as the metric to determine markers. On the other hand, choosing the

301   appropriate visual approach depends on the nature of the data and the problem. Our experiment

302   with the re-identifying cluster ID shows that the well-established SpatialDE [26] does not

303   outperform our approach and the DEG approach. One explanation is that in our problem, a good

304   marker for identifying cell type usually follows a good 'default' distribution over the visual space;

305   meanwhile, the SpatialDE aims to find markers that express significantly different from a default

306   distribution.

307

308   The major limitation of our proposed PGC-RSMD approach is the long computational time,

309   especially when comparing to the DEG approaches. This is similar to SpatialDE, which also used

310   visual information to detect marker genes. The DEG approaches may only need to compute one

311   statistical test to determine whether a gene is a marker for all clusters. Meanwhile, to draw the

312   curves, PGC-RSMD would need to compute hundreds to thousands, which depends on the desired

313   curve resolution, to characterize one gene in one cluster. Due to the long computational time, we

314   were not able to create multiple simulations, which is the ideal approach, run to compute the

315   statistical [32] p-value for the RSMD score. Therefore, we decided to reapply the estimation

316   presented to compute the p-value. This approach is computationally more efficient but may not

317   well-reflect the statistical characteristic of the single-cell data. In addition, we have not fully

318   tackled the problem of choosing the right threshold to determine whether a gene expresses in a

319   cell. Because of the strong association between PGC-RSMD and the percentage of expressing-

320   cell, we expect that the result would significantly different when choosing a different threshold to

321   determine whether a gene expresses in a cell. In this work, choosing 0 as the threshold still yields

322   good performance because of the high sparsity in the real dataset.

323

## Conclusions

325   In this work, we have presented Polar Gini Curve, a novel technique to detect markers from

326   the single-cell RNA expression data using visual information. In principle, our technique could

327   complement the state-of-the-art approach: the PGC technique finds markers such that the

328   expressing cells are evenly distributed throughout the cluster space; meanwhile, the state-of-the-

329   art approach finds markers assuming a multivariate normal distribution of gene expression in the

330    visual space. We have demonstrated that the PGC technique performs better in some tasks in
331    single-cell analysis.
332
333

## Authors' contribution

335    TN designed and implemented the core polar Gini algorithm, designed the sampling strategies
336    in the simulation, performed the neonatal heart single-cell case-study, and primarily prepared the
337    manuscript. JJ prepared the software package, preprocessed the single-cell data, and executed the
338    simulation designs. NX designed different simulation scenarios, interpreted the statistical
339    outcomes, and prepared the literature review. JC originated the idea of using Gini coefficient
340    curves to integrate gene expression, cell-point distribution, and cluster shape to solve the
341    biomarker discovery problem, designed the performance evaluation, and supervised the overall
342    technical development. All authors reviewed/revised and approved the manuscript.
343
344

## Competing interests

346    The authors have declared that no competing interests exist.
347
348

## Acknowledgments

352
353

## References

355    1.    Zhu, Z., et al., Single-cell transcriptome in the identification of disease biomarkers:
356          opportunities and challenges. J Transl Med, 2014. 12: p. 212.
357    2.    Finak, G., et al., MAST: a flexible statistical framework for assessing transcriptional
358          changes and characterizing heterogeneity in single-cell RNA sequencing data. Genome
359          Biol, 2015. 16: p. 278.

360   3.   McCarthy, D.J., et al., Scater: pre-processing, quality control, normalization and
361        visualization of single-cell RNA-seq data in R. Bioinformatics, 2017. 33(8): p. 1179-1186.

362   4.   Poirion, O.B., et al., Single-Cell Transcriptomics Bioinformatics and Computational
363        Challenges. Front Genet, 2016. 7: p. 163.

364   5.   Jaakkola, M.K., et al., Comparison of methods to detect differentially expressed genes
365        between single-cell populations. Brief Bioinform, 2017. 18(5): p. 735-743.

366   6.   Wang, T., et al., Comparative analysis of differential gene expression analysis tools for
367        single-cell RNA sequencing data. BMC Bioinformatics, 2019. 20(1): p. 40.

368   7.   Love, M.I., W. Huber, and S. Anders, Moderated estimation of fold change and dispersion
369        for RNA-seq data with DESeq2. Genome Biol, 2014. 15(12): p. 550.

370   8.   Robinson, M.D., D.J. McCarthy, and G.K. Smyth, edgeR: a Bioconductor package for
371        differential expression analysis of digital gene expression data. Bioinformatics, 2010.
372        26(1): p. 139-40.

373   9.   Pezzotti, N., et al., Approximated and User Steerable tSNE for Progressive Visual
374        Analytics. IEEE Trans Vis Comput Graph, 2017. 23(7): p. 1739-1752.

375   10.  Becht, E., et al., Dimensionality reduction for visualizing single-cell data using UMAP.
376        Nat Biotechnol, 2018.

377   11.  Qiu, X., et al., Reversed graph embedding resolves complex single-cell trajectories. Nat
378        Methods, 2017. 14(10): p. 979-982.

379   12.  Yang, Y., et al., SAFE-clustering: Single-cell Aggregated (from Ensemble) clustering for
380        single-cell RNA-seq data. Bioinformatics, 2019. 35(8): p. 1269-1277.

381   13.  Kiselev, V.Y., et al., SC3: consensus clustering of single-cell RNA-seq data. Nat Methods,
382        2017. 14(5): p. 483-486.

383   14.  Aibar, S., et al., SCENIC: single-cell regulatory network inference and clustering. Nat
384        Methods, 2017. 14(11): p. 1083-1086.

385   15.  Zheng, G.X., et al., Massively parallel digital transcriptional profiling of single cells. Nat
386        Commun, 2017. 8: p. 14049.

387   16.  Satija, R., et al., Spatial reconstruction of single-cell gene expression data. Nat Biotechnol,
388        2015. 33(5): p. 495-502.

389   17.  Kiselev, V.Y., T.S. Andrews, and M. Hemberg, Challenges in unsupervised clustering of
390        single-cell RNA-seq data. Nat Rev Genet, 2019. 20(5): p. 273-282.

391   18.   Yuan, G.C., et al., Challenges and emerging directions in single-cell analysis. Genome
392         Biol, 2017. 18(1): p. 84.
393   19.   Conover, W.J. and W.J. Conover, Practical nonparametric statistics. 1980.
394   20.   Hollander, M., D.A. Wolfe, and E. Chicken, Nonparametric statistical methods. Vol. 751.
395         2013: John Wiley & Sons.
396   21.   Guo, M., et al., SINCERA: A Pipeline for Single-Cell RNA-Seq Profiling Analysis. PLoS
397         Comput Biol, 2015. 11(11): p. e1004575.
398   22.   Birnbaum, Z. On a use of the Mann-Whitney statistic. in Proceedings of the Third Berkeley
399         Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the
400         Theory of Statistics. 1956. The Regents of the University of California.
401   23.   Korthauer, K.D., et al., A statistical approach for identifying differential distributions in
402         single-cell RNA-seq experiments. Genome Biol, 2016. 17(1): p. 222.
403   24.   Trapnell, C., et al., The dynamics and regulators of cell fate decisions are revealed by
404         pseudotemporal ordering of single cells. Nat Biotechnol, 2014. 32(4): p. 381-386.
405   25.   Kharchenko, P.V., L. Silberstein, and D.T. Scadden, Bayesian approach to single-cell
406         differential expression analysis. Nat Methods, 2014. 11(7): p. 740-2.
407   26.   Svensson, V., S.A. Teichmann, and O. Stegle, SpatialDE: identification of spatially
408         variable genes. Nat Methods, 2018. 15(5): p. 343-346.
409   27.   Edsgard, D., P. Johnsson, and R. Sandberg, Identification of spatial expression trends in
410         single-cell gene expression data. Nat Methods, 2018. 15(5): p. 339-342.
411   28.   Gini, C., Concentration and dependency ratios. Rivista di politica economica, 1997. 87: p.
412         769-792.
413   29.   Han, X., et al., Mapping the Mouse Cell Atlas by Microwell-Seq. Cell, 2018. 173(5): p.
414         1307.
415   30.   Strang, G., et al., Introduction to linear algebra. Vol. 3. 1993: Wellesley-Cambridge Press
416         Wellesley, MA.
417   31.   Hyndman, R.J. and A.B. Koehler, Another look at measures of forecast accuracy.
418         International journal of forecasting, 2006. 22(4): p. 679-688.
419   32.   Yue, Z., et al., WIPER: Weighted in-Path Edge Ranking for biomolecular association
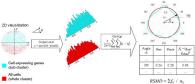420         networks. Quantitative Biology, 2019. 7(4): p. 313-326.

421   33.   Maaten, L.v.d. and G. Hinton, Visualizing data using t-SNE. Journal of machine learning
422         research, 2008. 9(Nov): p. 2579-2605.
423   34.   Bishop, C.M., Pattern recognition and machine learning. 2006: springer.
424   35.   MathWorks.    Matlab    -    inpolygon.    2019       2019/06/05];    Available    from:
425         https://www.mathworks.com/help/matlab/ref/inpolygon.html.
426   36.   MathWorks.    Matlab    -    boundary.    2019       2019/06/05];    Available    from:
427         https://www.mathworks.com/help/matlab/ref/boundary.html.
428   37.   Baruzzo, G., I. Patuzzi, and B. Di Camillo, SPARSim Single Cell: a count data simulator
429         for scRNA-seq data. Bioinformatics, 2019.
430   38.   Ester, M., et al. A density-based algorithm for discovering clusters in large spatial
431         databases with noise. in Kdd. 1996.
432   39.   MathWorks.       Matlab       -       dbscan.       2019;       Available       from:
433         https://www.mathworks.com/help/stats/dbscan.html.
434   40.   Bashyam, M.D., et al., Molecular genetics of familial hypertrophic cardiomyopathy (FHC).
435         J Hum Genet, 2003. 48(2): p. 55-64.
436   41.   Finsterer, J., C. Stollberger, and J.A. Towbin, Left ventricular noncompaction
437         cardiomyopathy: cardiac, neuromuscular, and genetic factors. Nat Rev Cardiol, 2017.
438         14(4): p. 224-237.
439   42.   Keren, A., P. Syrris, and W.J. McKenna, Hypertrophic cardiomyopathy: the genetic
440         determinants of clinical disease expression. Nat Clin Pract Cardiovasc Med, 2008. 5(3): p.
441         158-68.
442   43.   Morita, H., et al., Shared genetic causes of cardiac hypertrophy in children and adults. N
443         Engl J Med, 2008. 358(18): p. 1899-908.
444   44.   Jiang, H.K., et al., Reduced ACTC1 expression might play a role in the onset of congenital
445         heart disease by inducing cardiomyocyte apoptosis. Circ J, 2010. 74(11): p. 2410-8.
446   45.   Kwon, C., et al., A regulatory pathway involving Notch1/beta-catenin/Isl1 determines
447         cardiac progenitor cell fate. Nat Cell Biol, 2009. 11(8): p. 951-7.
448   46.   Wei, B. and J.P. Jin, TNNT1, TNNT2, and TNNT3: Isoform genes, regulation, and
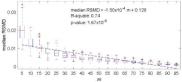449         structure-function relationships. Gene, 2016. 582(1): p. 1-13.

450    47.    Ju, Y., et al., Troponin T3 expression in skeletal and smooth muscle is required for growth
451           and postnatal survival: characterization of Tnnt3(tm2a(KOMP)Wtsi) mice. Genesis, 2013.
452           51(9): p. 667-75.

453    48.    Russell, S. and P. Norvig, Artifical Intelligence: A Modern Approach. 2003. Prentice Hall,
454           Upper Saddle River, New Jersey.

455    49.    Mehta, C.R. and N.R. Patel, A network algorithm for performing Fisher's exact test in r× c
456           contingency tables. Journal of the American Statistical Association, 1983. 78(382): p. 427-
457           434.

458    50.    Mukherjee, R. and O. Chakrabarti, Regulation of Mitofusin1 by Mahogunin Ring Finger-
459           1 and the proteasome modulates mitochondrial fusion. Biochim Biophys Acta, 2016.
460           1863(12): p. 3065-3083.

461    51.    Liu, X., et al., Differential microRNA Expression and Regulation in the Rat Model of Post-
462           Infarction Heart Failure. PLoS One, 2016. 11(8): p. e0160920.

463    52.    Lau, S.L., et al., Interferons induce the expression of IFITM1 and IFITM3 and suppress
464           the proliferation of rat neonatal cardiomyocytes. J Cell Biochem, 2012. 113(3): p. 841-7.

465    53.    Wang, L., et al., Mutations in myosin light chain kinase cause familial aortic dissections.
466           Am J Hum Genet, 2010. 87(5): p. 701-7.

467    54.    Shalek, A.K., et al., Single-cell RNA-seq reveals dynamic paracrine control of cellular
468           variation. Nature, 2014. 510(7505): p. 363-9.

469    55.    Rizzetto, S., et al., Impact of sequencing depth and read length on single cell RNA
470           sequencing data of T cells. Sci Rep, 2017. 7(1): p. 12781.

471    56.    McDavid, A., et al., Data exploration, quality control and testing in single-cell qPCR-based
472           gene expression experiments. Bioinformatics, 2013. 29(4): p. 461-7.

473    57.    Fortin, C., et al., Fractal dimension in the analysis of medical images. IEEE Engineering
474           in Medicine and Biology Magazine, 1992. 11(2): p. 65-71.

475    58.    Davies, S. and P. Hall, Fractal analysis of surface roughness by using spatial data. Journal
476           of the Royal Statistical Society: Series B (Statistical Methodology), 1999. 61(1): p. 3-37.

477    59.    Bharati, M.H., J.J. Liu, and J.F. MacGregor, Image texture analysis: methods and
478           comparisons. Chemometrics and intelligent laboratory systems, 2004. 72(1): p. 57-71.

479    60.    Kunimatsu, A., et al., Comparison between glioblastoma and primary central nervous

480           system lymphoma using MR image-based texture analysis. Magnetic Resonance in

481           Medical Sciences, 2017: p. mp. 2017-0044.

482

483

484    **Figure legends**

485    Figure 1. Overall workflow to compute the PGC-RSMD metric for one gene in one cluster of cells.

486    Here, the data points, histogram, and PGC for cells expressing the gene are cyan. The ones for the

487    whole cells in the cluster are red.

488

489    Figure 2. Boxplot showing a strong correlation between 'subcluster' percentage ($m$) and cluster-

490    subcluster PGC fitness (RSMD) in uniformly-distributed and a circular cluster.

491

492    Figure 3. a) The selected cluster for the experiment in [29]. b) Correlation between 'subcluster'

493    percentage ($m$) and cluster-subcluster PGC fitness (RSMD) in the selected cluster.

494

495    Figure 4. The ring-shape simulation study: a) Visualization of the cluster and ring-shape sub-

496    cluster ($m = 0.75$); b) PGC yield RSMD = 0.033; c) Distribution of RSMD, extracted from Figure

497    2 with $m = 75\%$, when the sub-cluster uniformly distributed on the cluster area.

498

499    Figure 5. PGC-RSMD performance in recalling cluster marker in drop-out simulation. a) heatmap

500    showing the simulation design of 500 markers and 4500 neutral genes, with drop out / percentage

501    of cell expressing between 5 and 100%; b) The simulation data 2D visualization; c) the AUC drops

502    when drop-out increases.

503

504    Figure 6. The result from mouse neonatal heart single-cell [29] analysis. a) the tSNE plot shows 9

505    clusters. b) gene-cluster marker relationship (from 258 genes) found by PGC-RSMD; ■ gene is

506    found as marker, ■ gene is found as non-marker. c) expression heatmap for these genes.

507

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

509

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

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

512    genes in each metric

513

514    Figure 9. PGC-RSMD shows that gene haves high percentage of expressing cells: *Ndufa4l2*, *Mdh2*,

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

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

517

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

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

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

521

2D visualization

Cell expressing genes (sub-cluster)

All cells (whole cluster)

Project on $y = \pm\cos(\theta) + x\sin(\theta)$

$$\frac{1}{2k^2 \beta} \sum_{i=1}^{k} \sum_{j=1}^{k} |z_i - z_j|$$

| Angle $\theta$ | $E_{min}$ | $S_{min}$ | $S^*_{min}(E_{min}/S_{min})$ |
|---|---|---|---|
| 0° | 0.26 | 0.29 | 0.65 |
| 30° | | | |

$$RSMD = \sum f_{ij} \; / \; n_{ij}$$

median RSMD = -1.50×10⁻⁴ m + 0.126
R-square: 0.74
p-value: 1.67×10⁻⁶

a)

b)

median RSMD = -1.65×10$^{-4}$ m + 0.015
R-square: 0.80
p-value: 2.06×10$^{-7}$

a) b) c)

a)

Cluster 1 marker    Cluster 2 marker    Neutral gene
250 genes           250 genes           (4500)

Cluster 1

Cluster 2

0   0.45   0      0   0.45   95%        5%
       drop out              % cell expressing ($m$)

Log10 expression: 3.5, 3, 2.5, 2, 1.5, 1, 0.5

b)

• Cluster 1
• Cluster 2

c)

cluster 1
cluster 2

AUC

Drop-out probability

a) (tSNE plot with axes tSNE 1 and tSNE 2, legend 1–9)

b) (heatmap with x-axis Cluster ID, labels 1–9)

c) (heatmap with x-axis Cluster ID, labels 1–9)

**Actc1, expression in cluster 1** — tSNE plot with Expressing cells (red) and Non-expressing cells (cyan). % Expressing cell: 41 (252)

**Actc1, PGCs in cluster 1** — polar plot. RSMD: 0.013 (73). All cells (blue), Expressing cells (red).

**Mgrn1, expression in cluster 1** — tSNE plot with Expressing cells (red) and Non-expressing cells (cyan). % Expressing cell: 59 (182)

**Mgrn1, PGCs in cluster 1** — polar plot. RSMD: 0.010 (43). All cells (blue), Expressing cells (red).

**Ifitm3, expression in cluster 1** — tSNE plot with Expressing cells (red) and Non-expressing cells (cyan). % Expressing cell: 37 (288)

**Ifitm3, PGCs in cluster 1** — polar plot. RSMD: 0.012 (56). All cells (blue), Expressing cells (red).

**Myl6b, expression in cluster 1** — tSNE plot with Expressing cells (red) and Non-expressing cells (cyan). % Expressing cell: 45 (238)

**Myl6b, PGCs in cluster 1** — polar plot. RSMD: 0.015 (87). All cells (blue), Expressing cells (red).

a) Accuracy vs # selected markers with legend: PGC-RSMD, SpatialDE, DEG, and baseline.

b) AUC vs # selected markers with legend: PGC-RSMD, SpatialDE, DEG, and baseline.