
Determining Cancerous Regions from Spatial Transcriptomics Data of Breast Cancer Tissue

Max Land¹, Rebecca Elyanow², Ben Raphael²

¹ Department of Mathematics, Princeton University, Princeton, NJ 08540, USA

² Department of Computer Science, Princeton University, Princeton, NJ 08540, USA

Abstract

Previous methods of determining normal cells vs cancer cells based on gene expression data have relied on analyzing a subset of tissue specific genes (e.g. BRCA2 mutation in breast tissue). However, the diversity of gene expression in cancer cells presents difficulties for this approach. In this paper, we present two novel classification methods that are not dependent on the presence of specific genes. Instead, these methods consider the entire transcriptome as a whole and use PCA and percent dropout to reduce our data to one dimensional feature vectors. Next, we apply a simple thresholding classification method to these feature vectors to distinguish normal cells from healthy cells. We find that this novel clustering approach according to the first principal component of PCA and percent dropout have high performance and present a promising new approach in identifying tumor regions in cancer tissue.

1 Introduction

Due to the development of single-cell sequencing technologies [10], it is well known that multi-cellular organisms can generate and maintain homogeneous populations of healthy cells in individual tissues [4]. Cancer, however, can disrupt homogeneity and cause large changes to gene expression of infected cells. This is because cancer often involves mutations in proto-oncogenes, tumor suppressor genes, and other regulatory regions of the genome. These mutations alter the transcription rate and genetic diversity of cancer cells. This is reflected in single-cell sequencing studies of breast [7], kidney [12], bladder[5], and colon tumors [3]. Thus, given the changes in gene expression caused by cancer, an apparent question arises: Can we distinguish healthy cells vs cancer cells based on gene expression data?

The answer, as it turns out, is not so simple. Current methods involve constructing gene expression profiles of cancer cells to identify specific genes prevalent in cancer cells [2]. Researchers then look for the expression of these genes in cells to identify which cells are cancerous. However, the problem with this approach is that cancer is actually known for its heterogeneity at both the inter- and intra-tumor levels [1, 2]. Within a tumor, different spatial regions have been found to contain different cancer cell clones [1]. In addition, tumor cells are also infiltrated with stromal, immune and other cell types. These cells form the basis of tumor heterogeneity [11, 6]. As a result, characterization of tumor cells from tissue-wide gene expression data is an inherently difficult challenge. For example, CIBERSORT, a method for characterizing cell composition of complex tissues from their gene expression profiles, performs considerably worse on samples containing solid tumors [8]. Additionally, to analyze different types of cancer tissue, researchers would need to construct separate gene expression profiles.

In this paper, we propose two novel methods to distinguish between normal cells and cancer cells without requiring prior characterization and establishment of gene expression profiles of various types of cancer cells. Instead, we focus on the transcriptome as a whole. Intuitively, we know there is a difference in gene expression data between normal and cancer cells. Our two models use PCA and percent dropout respectively to capture this difference in a one-dimensional feature

vector. Next, we apply a simple thresholding classification to classify our samples. To analyze the performance of our models, we utilize spatial transcriptomics data. Spatial transcriptomics (ST) is a new technology that allows researchers to measure gene activity in a tissue sample and map the location of the activity [9]. The spatial pattern of gene expression within a tissue is obtained by sequencing mRNA from a grid spots, which we refer to as ST spots, each containing a small number of cells.

In the case of cancer tissue, the nature of spatial transcriptomics allows us to view the location of our ST spots in the original tissue sample. We can visually see which ST spots come from cancerous vs non-cancerous regions of the tissue by plotting the spots on a cancer visualization image of the tissue. This creates a unique gene expression data set of cancer tissue ST spots with the true label known for each ST spot: cancerous or non-cancerous. This, in turn, allows us to perform exploratory cluster analysis to investigate novel methods to determine cancerous from non-cancerous tissue regions using the spatial transcriptomics data. In particular, we explored models based on PCA and percent dropout as mentioned previously in this paper.

2 Methods

2.1 Data/Data Pre-processing

We worked with a publicly available breast cancer data set published by the Spatial Research Lab [9]. The data set consists of four layers of breast cancer tissue sliced from the same sample. For each layer, we have:

1. Spatial Transcriptomics Data: Gene expression matrix, where genes are columns and ST spots (identified by coordinates) are rows. For all datasets, spots outside the tissue were removed.
2. Visualization of Cancerous Regions: Hematoxylin & eosin stained brightfield image of the tissue layer. In these images, the dark purple tissue corresponds to the cancerous regions and light purple tissue corresponds to the healthy, non-cancerous regions.

Next, we generated our set of true labels l_{true} for each layer. To do so, we plotted each ST spot on our stained brightfield image. Then, we manually labeled the ST spot as cancerous or non-cancerous depending on the location. We report a visualization of l_{true} below:

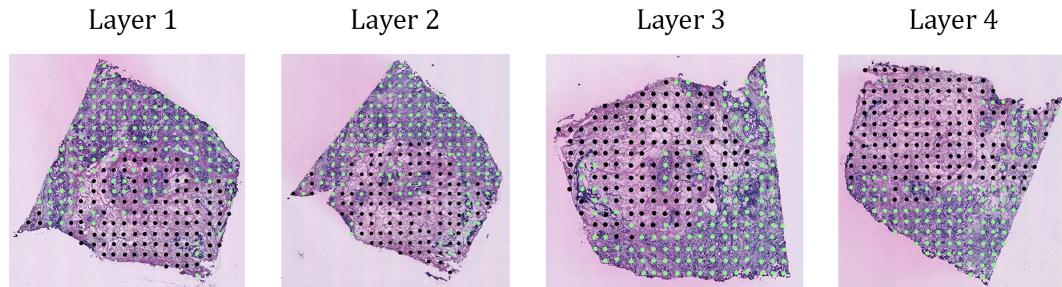


Figure 1: Visualization of our manually clustered ST spots.

Note: This visual classification of tumor regions of the cell is reflective of the practice used by oncologists. The light green represents a cancerous ST spot, and black represents a non-cancerous ST spot.

2.2 Method 1: PCA Classifier

We investigated the ability of the first principal component of PCA to differentiate between cancerous and non-cancerous regions of the tissue. To do so, we first preprocessed the spatial transcriptomics data using the following steps:

- 108 1. First, we filtered for genes with at least 15 nonzero entries. In others, genes with non-zero
109 counts in fewer than 15 ST spots were filtered out.
- 110 2. The data is then library-size normalized, where counts for each gene in a spot are divided by
111 the sum of counts for the ST spot. This normalization removes variance due to differences
112 in size or number of cells in each ST spot.
- 113 3. Next, the data is log normalized with pseudocount 1 ($\log(X + 1)$).
- 114 4. Finally, we Z-normalized each ST spots (by row).

117 Next, we applied PCA to generate our first principal component vector. For a cancer tissue layer
118 with n ST spots sequenced, we get the first principal component vector:

$$120 \quad pc = (pc_1, \dots, pc_n)$$

121 Then, we performed a thresholding classification to cluster the spots into cancerous vs non-
122 cancerous. We wanted to get a vector of predictive labels for each spot:

$$124 \quad l_{\text{pred}} = (l_{\text{pred}_1}, \dots, l_{\text{pred}_n})$$

126 To do so, we clustered each spots' respective principal component vector pc against a threshold
127 value t . For some t and $i = 1, \dots, n$,

$$129 \quad l_{\text{pred}} = \begin{cases} 1 & \text{if } pc_i \leq t \\ 0 & \text{if } pc_i > t \end{cases}$$

132 Note: Preliminary analysis showed that tumor cells tend to have lower pc values. So, we assigned
133 our prediction labels l_{pred} to be consistent with l_{true} which assigned cancerous ST spots to 1 and
134 non-cancerous ST spots to 0.

135 In summary, given a threshold t , we can generate a set of predictive labels l_{pred} based on PCA.

136 We will refer to the function as:

$$138 \quad \text{pca_classifier}(pc, t) = l_{\text{pred}}$$

141 2.3 Method 2: Percent Dropout Classifier

142 We investigated ability of using percent dropout to cluster cancerous vs non-cancerous spots. From
143 the spatial transcriptomics data, we calculated the percent of genes with zero counts for each ST
144 spot. For n ST spots, this resulted in a n -length percent dropout vector:

$$146 \quad pd = (pd_1, \dots, pd_n)$$

148 From here, we repeated the classification steps from our PCA Classifier. We performed a threshold-
149 ing classification to cluster the spots into cancerous vs non-cancerous. The goal was to get a vector
150 of predictive labels for each spot:

$$151 \quad l_{\text{pred}} = (l_{\text{pred}_1}, \dots, l_{\text{pred}_n})$$

153 To do so, we clustered each spots' respective percent dropout component against a threshold value
154 t . For some t and $i = 1, \dots, n$,

$$155 \quad l_{\text{pred}_i} = \begin{cases} 1 & \text{if } pd_i \leq t \\ 0 & \text{if } pd_i > t \end{cases}$$

158 Thus, given a threshold t , we can generate a set of predictive labels l_{pred} based on percent dropout.

160 We will refer to the function as:

$$161 \quad \text{pd_classifier}(pd, t) = l_{\text{pred}}$$

162 **2.4 Classifier Evaluation**
163

164 We tested the effectiveness of our PCA classifier and percent dropout classifier on the breast cancer
165 data set. We used two main metrics to evaluate performance:

- 166 1. Adjusted Rand Index (ARI): The Rand Index (RI) computes a similarity measure between
167 two clusterings by considering all pairs of samples and counting pairs that are assigned in
168 the same or different clusters in the predicted and true clusterings. The raw RI score is then
169 adjusted for chance by:

170
$$\text{ARI} = \frac{\text{RI} - \text{Expected_RI}}{\max(\text{RI}) - \text{Expected_RI}}$$

171 The ARI score is calculated using scikit-learn's: `adjusted_rand_score()`.

- 172 2. Precision-Recall Curve: First, we define precision and recall:

173
$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \text{ and } \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

174 where TP = # of True Positives, FP = # of False Positives, and FN = # of False Negatives.

175 Precision-Recall curves summarize the trade-off between the true positive rate and the positive predictive value for a predictive model using different probability thresholds.

176 **2.5 Choosing a threshold**
177

178 For each breast cancer tissue layer, we wanted to calculate ARI score and Precision-Recall score for
179 each classifier. However, each classifier is function of a threshold t . Thus the question becomes,
180 how do we pick a threshold t ?

181 To address this, we calculated ARI scores and Precision-Recall scores for various thresholds t . To
182 choose t , we took 200 different values of t evenly spaced between $[pc_{\min}, pc_{\max}]$ or $[pd_{\min}, pd_{\max}]$
183 (depending on the classifier) to generate 200 different predictions l_{pred} .

184 To clarify:

- 185 • Given principal component vector pc ,

186
$$pc_{\max} = \max\{pc_1, \dots, pc_n\} \text{ and } pc_{\min} = \min\{pc_1, \dots, pc_n\}$$

- 187 • Given percent dropout vector pd

188
$$pd_{\max} = \max\{pd_1, \dots, pd_n\} \text{ and } pd_{\min} = \min\{pd_1, \dots, pd_n\}$$

189 Here is the algorithm in pseudo-code for PCA:

200 **Algorithm :** Calculate ARI score and Precision-Recall for 200 different threshold values

201 1: $pc \leftarrow \text{get_Principal_Component}(\text{SpatialTranscriptomicsData})$
202 2: **procedure** EVALUATE_CLASSIFIER(pc, l_{true})
203 3: $t \leftarrow \text{linspace}(\min(pc), \max(pc), 200)$
204 4: **for** $i = 1 : 200$ **do**
205 5: $l_{\text{pred}} \leftarrow \text{pca_Classifier}(pc, t_i)$
206 6: $ari \leftarrow \text{calculate_ARI}(l_{\text{true}}, l_{\text{pred}})$
207 7: $\text{precision} \leftarrow \text{calculate_Precision}(l_{\text{true}}, l_{\text{pred}})$
208 8: $\text{recall} \leftarrow \text{calculate_Recall}(l_{\text{true}}, l_{\text{pred}})$
209 9: **end for**
210 10: **return** $ari, \text{precision}, \text{recall}$
211 11: **end procedure**

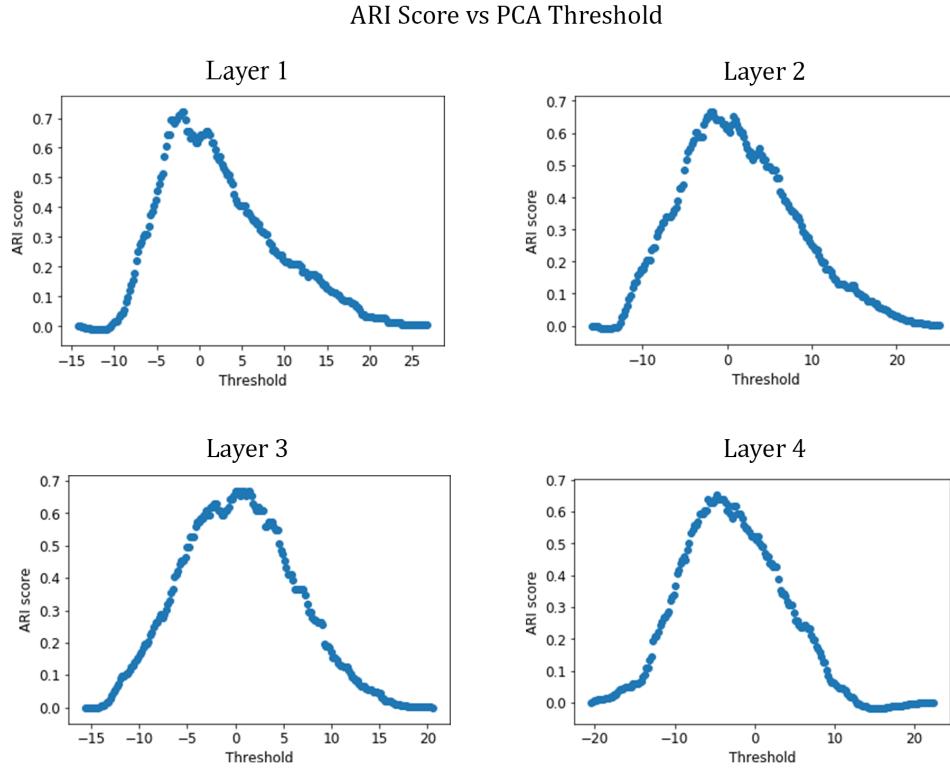
212 The algorithm for Percent Dropout follows the same design.

213

216 **3 Results**
 217

218 **3.1 Results of ARI Score: PCA Classifier**
 219

220 The publicly available breast cancer data set consists of four layers. For each layer, we tested the
 221 performance of our PCA classifier by calculating the ARI score as a function of a threshold t . We
 222 report the results below:
 223



252 Next, we report the maximum ARI score and the threshold t associated with that score. Since we are
 253 classifying by PCA, we also report the proportion of total variance explained by the first principal
 254 component:
 255

	Layer 1	Layer 2	Layer 3	Layer 4
Max ARI	0.722	0.666	0.668	0.653
Optimal t	-2.060	-1.932	1.478	-4.739
Prop of Total Variance Explained	0.012	0.015	0.013	0.015

260 As seen in Figure 2, the ARI score graphed as a function of t is a peaked graph. For optimal t , we see
 261 that the maximum ARI score ranges from 0.65 – 0.72 across all four layers, despite our first principal
 262 components explaining only 1 – 2% of the total variance of our data. This shows that clustering the
 263 ST spots according to the first principal component is an effective method to distinguish between
 264 cancerous and non-cancerous.
 265

Next, we created a visualization of the best prediction by the PCA Classifier. For the optimal t that resulted in the maximum ARI score for each layer, we plotted the ST spots on the tissue image color-coded by its predicted label given by the PCA Classifier. As a reminder, light green corresponds to a cancerous ST spot, and black corresponds to a non-cancerous ST spot. We report the results below:

Best Clustering According to PCA Threshold

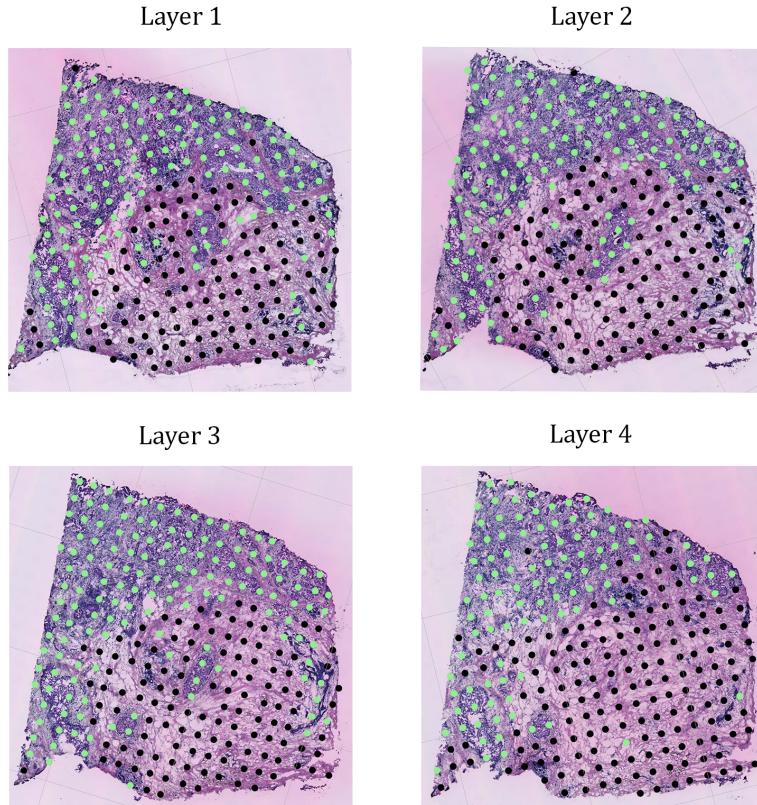
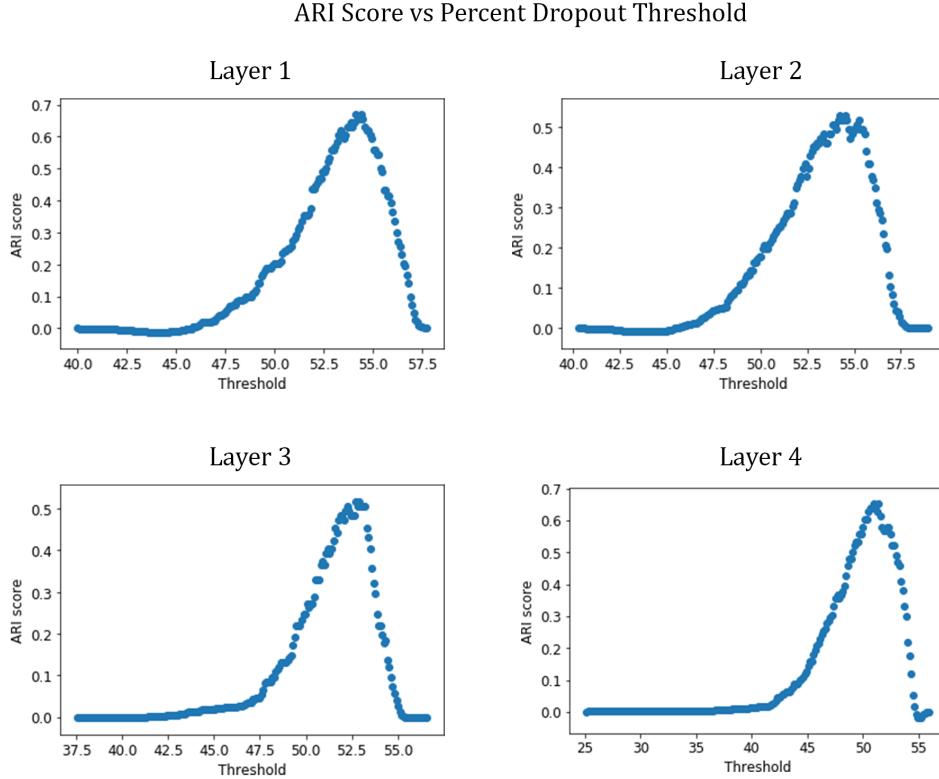


Figure 3: Each layer reflects the clustering corresponding to the highest ARI achieved in Figure 2.

The visualization of the predictions in Figure 3 show that the PCA classifier is effective at classifying the "obvious" ST spots— those that are clearly in a cancerous region or non-cancerous region. Next, comparing the visualization of our PCA classifier (Figure 3) to the true labels (Figure 1), we see the misclassifications mainly occur at the "no man's land" where the cancerous and non-cancerous regions meet.

324 **3.2 Results of ARI Score: Percent Dropout Classifier**
 325

326 Next, we repeated the same analysis for our Percent Dropout Classifier. For each layer, we tested
 327 the performance of our Percent Dropout classifier by calculating the ARI score as a function of a
 328 threshold t . We report the results below:



356 Figure 4: Performance Evaluation of the Percent Dropout Classifier: ARI Score
 357

358 Next, we report the maximum ARI score and the threshold t associated with that score:
 359

360
 361

	Layer 1	Layer 2	Layer 3	Layer 4
Max ARI	0.669	0.530	0.516	0.653
Optimal t	54.128	54.240	52.916	51.410

364 For the Percent Dropout Classifier, we see that the ARI score graphed as a function of t is also a
 365 peaked graph in Figure 4. Compared to the PCA Classifier, the maximum ARI score is slightly
 366 lower, but still ranges from 0.51 – 0.67. This result suggest that the PCA classifier is more effective
 367 than the Percent Dropout classifier.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Next, we created a visualization of the best prediction by the Percent Dropout Classifier. For the optimal t that resulted in the maximum ARI score for each layer, we plotted the ST spots on the tissue image color-coded by its predicted label given by the Percent Dropout Classifier. As a reminder, light green corresponds to a cancerous ST spot, and black corresponds to a non-cancerous ST spot. We report the results below:

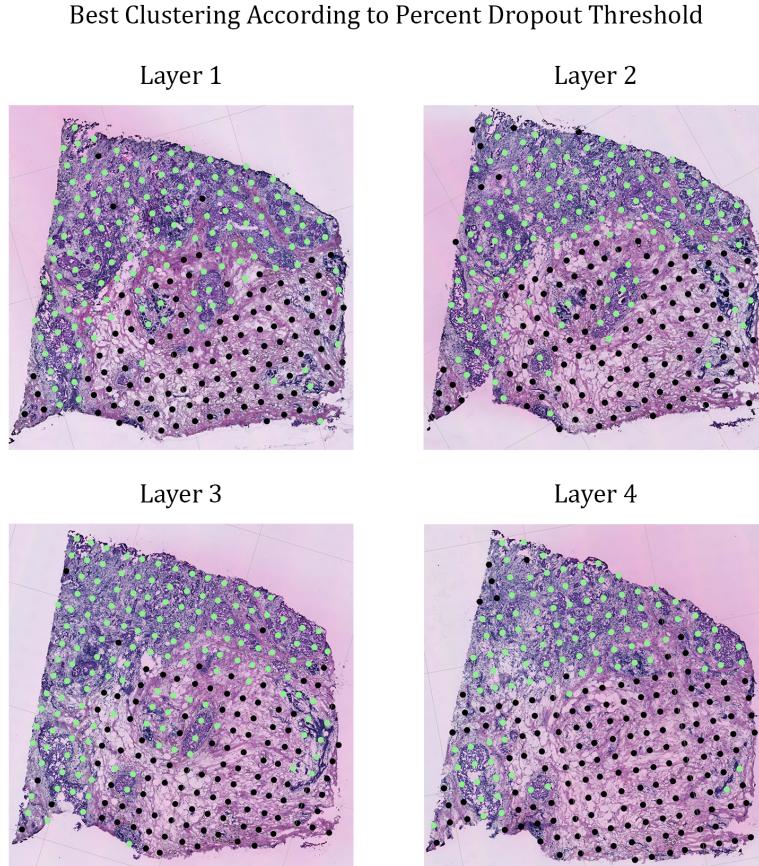


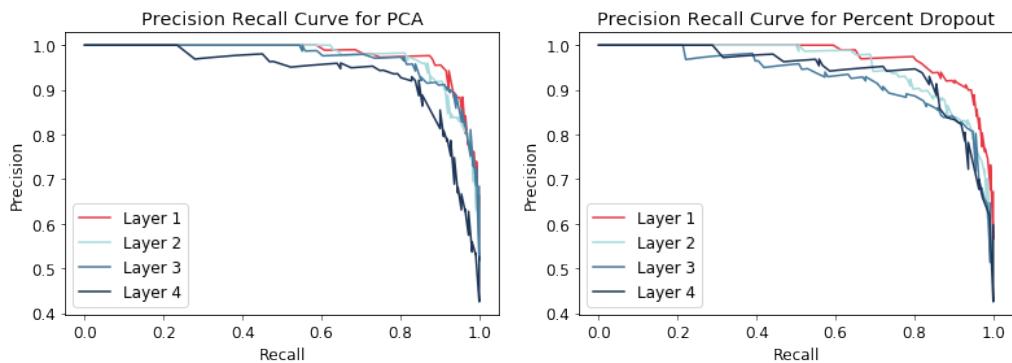
Figure 5: Each layer reflects the clustering corresponding to the highest ARI achieved in Figure 4.

Looking at the visualization of predictions in Figure 5, we see the predictions are slightly worse. For the most part, the Percent Dropout Classifier does classify the majority of ST spots correctly—the large regions of cancerous and non-cancerous are correct. However, compared to the PCA classifier in Figure 3, there are obvious misclassifications in the large cancerous regions. Again, this substantiates the claim from previous result that the PCA Classifier has better performance than the Percent Dropout Classifier.

432 **3.3 Precision-Recall Curves**
 433

434 Another way to assess our classifiers is through precision-recall curves. As seen in the algorithm in
 435 section 2.5, for each threshold t , we calculate a precision and recall score. Because we calculate 200
 436 values of t for each layer, this results in 200 data points of the form: (recall, precision)

437 We then used those points to create a Precision-Recall curve for each layer, and we report the results
 438 below:



453 Figure 6: Precision-Recall Curves for PCA
 454

455 Overall, the precision-recall curves for both the PCA Classifier and the Percent Dropout Classifier
 456 are curved appropriately. This indicates good performance.
 457

458 Furthermore, we see that for both classifiers, Layer 1 has the best performance. This is consistent
 459 with our ARI results seen in Figures 2 and 4. As a matter of fact, the order of performance of the
 460 four layers for each classifier reflects the order reported by our ARI results.

461 Next, comparing the two precision-recall curves against each other, we see that the precision-recall
 462 curve for PCA is curved more towards (1, 1). This substantiates our analysis that the PCA Classifier
 463 has better performance than the Percent Dropout Classifier.

464 **3.4 Distribution of Components**
 465

466 Both the PCA and Percent Dropout classifiers return a prediction of cancerous or non-cancerous
 467 for each ST spot dependent on a threshold t . For optimal t , we see that the our PCA and Percent
 468 Classifier perform well. However, in real applications, we don't know the true labels, thus we cannot
 469 experimentally determine the optimal threshold. Thus, this begs the following question: Is there a
 470 way to select a good threshold t ?
 471

472 To gain some insight into choosing a threshold t , we plot the distributions of our PCA components.
 473 Specifically, we split the PCA components into their underlying cancerous and non-cancerous dis-
 474 tributions based on the true labels. We also plot the experimentally calculated optimal threshold t
 475 for reference. As a reminder:

	Layer 1	Layer 2	Layer 3	Layer 4
Optimal t	-2.060	-1.932	1.478	-4.739

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

Next, for each layer, we investigate the distribution of (pc_1, \dots, pc_n) for our principal component vector pc . We report our distribution plots:

Density Plots of PCA Components

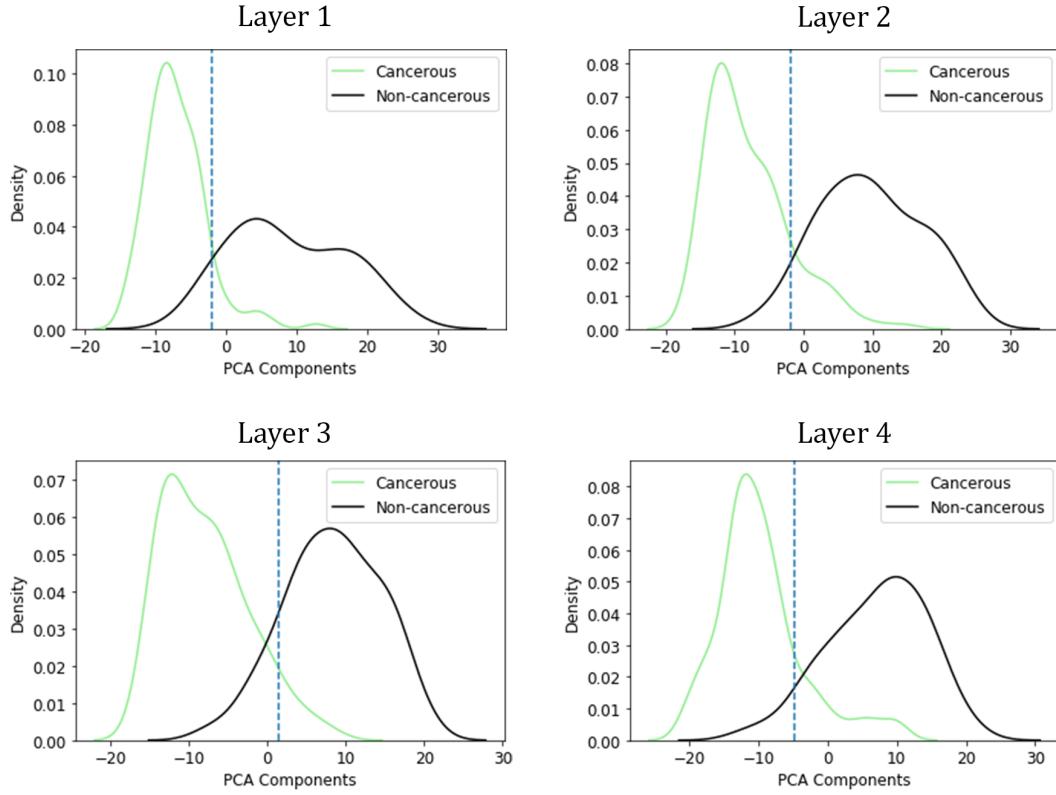


Figure 7: Distribution of PCA components for cancerous and non-cancerous ST spots.

From Figure 7, we immediately notice that the cancerous and non-cancerous ST spots are drawn from two separate underlying distributions. This further substantiates that the classifying by the first principal component can reasonably distinguish between cancerous and non-cancerous ST spots. In light of this, it appears picking the optimal t becomes as simple as picking the overlap between the two tails of the respective distributions. This is confirmed by our experimental results—the optimal ARI occurs at the overlap of these two distributions, as indicated by the dotted blue line.

Again, however, it is important to recognize that it is not this straightforward in practice. In real applications, we do not know the true set of cancerous and non-cancerous spots. Thus, we cannot split the pc values into the two underlying distributions. Instead, we will only have the combined overall distribution. Thus, it might make more sense to view our problem in a more general sense:

"Given that we know our vector of values pc are drawn from two underlying distributions, can we accurately classify each individual data points into the respective distribution it is drawn from?"

Also, it is important to keep in mind that the two distributions will most likely overlap as seen in Figure 7. Thus, our current method of simply classifying by an optimal t will inherently result in misclassifications (both false positives and false negatives). However, there are ways to improve prediction by taking advantage of the spatial location of our ST spots. Our results show that cancerous and non-cancerous cells are drawn from distinct distributions. The overlap, therefore, most likely represent ST spots that are in the process of becoming cancerous. In Figure 3, we see that most of the misclassifications are ST spots that lie in between cancerous and non-cancerous regions. Intuitively, the ST spots that lie closer to definitive cancerous ST spots are more likely to be cancer-

ous and vice versa. So, for ST spots that lie in that overlap region, future models can incorporate the spatial information to resolve misclassifications.

Next, we plot the distribution of our Percent Dropout components. Specifically, we split the Percent Dropout components into their underlying cancerous and non-cancerous distributions based on the true labels. We also plot the experimentally calculated optimal threshold t for reference. As a reminder:

	Layer 1	Layer 2	Layer 3	Layer 4
Optimal t	54.128	54.240	52.916	51.410

Next, for each layer, we investigate the distribution of (pd_1, \dots, pd_n) for our percent dropout component vector pd . We report our distribution plots:

Density Plots of Percent Dropout Components

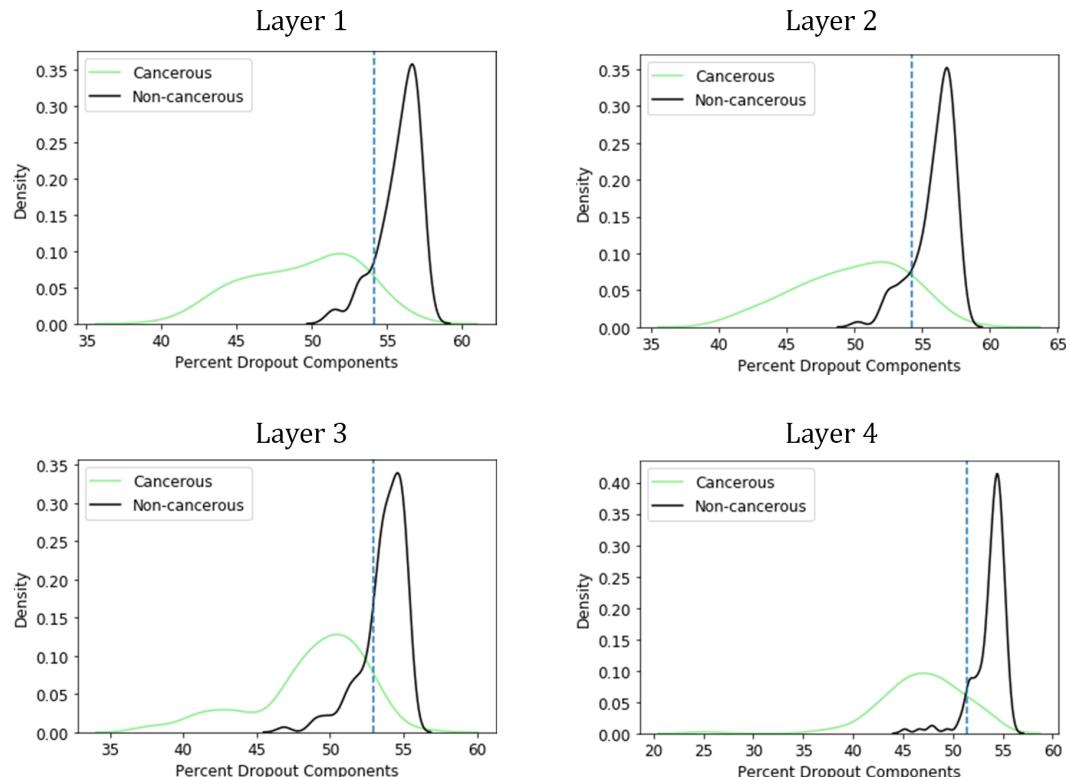


Figure 8: Distribution of Percent Dropout components for cancerous and non-cancerous ST spots.

Similarly to the PCA density plots, we see that cancerous and non-cancerous cells form two different underlying distributions in Figure 8. Most notably, the non-cancerous distribution is extremely peaked compared to the cancerous distribution.

The results for percent dropout are also more interpretable—cancerous ST spots tend to have a lower percent dropout. Percent dropout measures the percent of genes with zero expression. Thus, cancerous ST spots have a higher number of unique genes expressed. This makes sense, as cancerous cells arise from additional mutations in proto-oncogenes, tumor suppressor genes, and other regulatory regions of the genome. This biological understanding also explains the difference in kurtosis between the two distributions. Normal breast cancer cells drawn from the same tissue sample have homogeneous gene expressions [4]. This is reflected in the sharply peaked non-cancerous distribution. On the other hand, cancerous mutations are a result of random mutations. Different ST spots

594 are also at different stages of cancer (some spots have been infected longer than others). Thus, there
595 is a higher chance for variance in the diversity of gene expression in cancerous ST spots. This is
596 reflected in the distribution being less peaked and more distributed across percent dropout values.
597 Furthermore, previous research has shown cancer has heterogeneity within tumors [1].
598

599 4 Conclusion 600

601 In this paper, we proposed two novel methods to determine cancerous regions from spatial transcriptomics
602 data of cancer tissue. Specifically, the PCA Classifier clusters based on the first principal
603 component after applying PCA and the Percent Dropout Classifier clusters based on the proportion
604 of zero expression genes. We demonstrated that both classifiers show promising results with good
605 performance on breast cancer data. Out of the two, the PCA Classifier performs slightly better.
606 Further analysis revealed that both classifiers effectively reduce the transcriptome to one dimensional
607 feature vectors, and that the distribution of these feature vectors are not random, but rather a
608 combination two underlying distributions (cancerous vs non-cancerous). This allows a thresholding
609 classification technique to be effective at clustering accurately.

610 The results of this paper are an initial exploration of analyzing spatial transcriptomics data of cancer
611 tissue. There several directions for future work. The first is incorporating spatial information. In
612 general, tumors arise and proliferate rapidly—cancerous ST spots will often be found in large connected
613 regions. In the visualizations of our classifier results (Figures 3, 5), we see there are obvious
614 misclassifications (i.e. a supposedly non-cancer ST spot surrounded by cancerous ST spots). The
615 spatial information can be used ensure our classification results are more in line with our fundamental
616 understanding of cancer. We can also stack the layers together to create a 3-D spatial network.
617 Another important extension is to evaluate our classifiers on other types of spatial transcriptomics
618 data. For example, there is available transcriptomics data on prostate cancer tissue [2]. Furthermore,
619 it would be interesting to see what our PCA classifier would detect if there were no tumor regions
620 to detect.

621 5 Code Availability 622

624 The methods in this paper are implemented in Python and are available at: <https://github.com/mrland99/spatial-research>
625

627 References 628

- 629 [1] Noemi Andor, Trevor Graham, Marnix Jansen, Li Xia, C. Aktipis, Claudia Petritsch, Han-
630 lee Ji, and Carlo Maley. Pan-cancer analysis of the extent and consequences of intratumor
631 heterogeneity. *Nature Medicine*, 22, 11 2015.
- 632 [2] Emelie Berglund, Jonas Maaskola, Niklas Schultz, Stefanie Friedrich, Maja Marklund, Joseph
633 Bergenstråhle, Firas Tarish, Anna Tanoglidi, Sanja Vickovic, Ludvig Larsson, Fredrik Salmén,
634 Christoph Ogris, Karolina Wallenborg, Jens Lagergren, Patrik Ståhl, Erik Sonnhammer,
635 Thomas Helleday, and Joakim Lundeberg. Spatial maps of prostate cancer transcriptomes
636 reveal an unexplored landscape of heterogeneity. *Nature Communications*, 9, 12 2018.
- 637 [3] Ellen Heitzer, Martina Auer, Christin Gasch, Martin Pichler, Peter Ulz, Eva Hoffmann, Sigurd
638 Lax, Julie Waldispuehl-Geigl, Oliver Mauermann, Carolin Lackner, Gerald Höfler, Florian
639 Eisner, Heinz Sill, Hellmut Samonigg, Klaus Pantel, Sabine Riethdorf, Thomas Bauernhofer,
640 Jochen Geigl, and Michael Speicher. Complex tumor genomes inferred from single circulating
641 tumor cells by array-cgh and next-generation sequencing. *Cancer research*, 73, 03 2013.
- 642 [4] Hai Le, Monika Looney, Benjamin Strauss, Michael Bloodgood, and Antony Jose. Tissue
643 homogeneity requires inhibition of unequal gene silencing during development. *The Journal
644 of Cell Biology*, 214:jcb.201601050, 07 2016.
- 645 [5] Yingrui Li, Xun Xu, Luting Song, Yong Hou, Zesong Li, Shirley Tsang, Fuqiang Li, Kate
646 Im, Kui Wu, Hanjie Wu, Xiaofei Ye, Guibo Li, Linlin Wang, Bob Zhang, Jie Liang, Wei Xie,
647 Renhua Wu, Hui Jiang, Xiao Liu, and Jun Wang. Single-cell sequencing analysis characterizes

- 648 common and cell-lineage-specific mutations in a muscle-invasive bladder cancer. *GigaScience*,
649 1:12, 08 2012.
- 650 [6] Carlo Maley, C. Aktipis, Trevor Graham, Andrea Sottoriva, Amy Boddy, Michalina
651 Janiszewska, Ariosto Silva, Marco Gerlinger, Yinyin Yuan, Kenneth Pienta, Karen Anderson,
652 Robert Gatenby, Charles Swanton, David Posada, Chung-I Wu, Joshua Schiffman, E Hwang,
653 Kornelia Polyak, Alexander Anderson, and Darryl Shibata. Classifying the evolutionary and
654 ecological features of neoplasms. *Nature reviews. Cancer*, 17, 09 2017.
- 655 [7] Nicholas Navin, Jude Kendall, Jennifer Troge, Peter Andrews, Linda Rodgers, Jeanne McIn-
656 doo, Kerry Cook, Asya Stepansky, Dan Levy, Diane Esposito, Lakshmi Muthuswamy, Alexan-
657 der Krasnitz, W. Mccombie, James Hicks, and Michael Wigler. Tumor evolution inferred by
658 single cell sequencing. *Nature*, 472:90–4, 03 2011.
- 659 [8] Aaron Newman, Chih Liu, Michael Green, Andrew Gentles, Weigu Feng, Yue Xu, Chuong
660 Hoang, Maximilian Diehn, and Ash Alizadeh. Robust enumeration of cell subsets from tissue
661 expression profiles. *Nature methods*, 12, 03 2015.
- 662 [9] Patrik L Ståhl, Fredrik Salmén, Sanja Vickovic, Anna Lundmark, José Fernández Navarro,
663 Jens Magnusson, Stefania Giacomello, Michaela Asp, Jakub O Westholm, Mikael Huss, et al.
664 Visualization and analysis of gene expression in tissue sections by spatial transcriptomics.
665 *Science*, 353(6294):78–82, 2016.
- 666 [10] Fuchou Tang, Catalin Barbacioru, Yangzhou Wang, Ellen Nordman, Clarence Lee, Nanlan Xu,
667 Xiaohui Wang, John Bodeau, Brian Tuch, Asim Siddiqui, Kaiqin Lao, and M Surani. mrna-seq
668 whole-transcriptome analysis of a single cell. *Nature methods*, 6:377–82, 05 2009.
- 669 [11] Kenneth Valkenburg, Amber Groot, and Kenneth Pienta. Targeting the tumour stroma to im-
670 prove cancer therapy. *Nature Reviews Clinical Oncology*, 15, 04 2018.
- 671 [12] Xun Xu, Yong Hou, Xuyang Yin, Bao Li, Aifa Tang, Luting Song, Fuqiang Li, Shirley Tsang,
672 Kui Wu, Hanjie Wu, Weiming He, Li Zeng, Manjie Xing, Renhua Wu, Hui Jiang, Xiao Liu,
673 Dandan Cao, Guangwu Guo, Xueda Hu, and Jun Wang. Single-cell exome sequencing reveals
674 single-nucleotide mutation characteristics of a kidney tumor. *Cell*, 148:886–95, 03 2012.
- 675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701