

Manuscript Number:	GIGA-D-18-00086R2	
Full Title:	Clustering trees: a visualisation for evaluating clusterings at multiple resolutions	
Article Type:	Research	
Funding Information:	Department of Education, Australian Government	Mr Luke Zappia
	National Health and Medical Research Council (APP1126157)	Dr Alicia Oshlack
Abstract:	<p>Clustering techniques are widely used in the analysis of large data sets to group together samples with similar properties. For example, clustering is often used in the field of single-cell RNA-sequencing in order to identify different cell types present in a tissue sample. There are many algorithms for performing clustering and the results can vary substantially. In particular, the number of groups present in a data set is often unknown and the number of clusters identified by an algorithm can change based on the parameters used. To explore and examine the impact of varying clustering resolution we present clustering trees. This visualisation shows the relationships between clusters at multiple resolutions allowing researchers to see how samples move as the number of clusters increases. In addition, meta-information can be overlaid on the tree to inform the choice of resolution and guide in identification of clusters. We illustrate the features of clustering trees using a series of simulations as well as two real examples, the classical iris dataset and a complex single-cell RNA-sequencing dataset. Clustering trees can be produced using the clustree R package available from CRAN (https://CRAN.R-project.org/package=clustree) and developed on GitHub (https://github.com/lazappi/clustree).</p>	
Corresponding Author:	Alicia Oshlack AUSTRALIA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Luke Zappia	
First Author Secondary Information:		
Order of Authors:	Luke Zappia Alicia Oshlack	
Order of Authors Secondary Information:		
Response to Reviewers:	We have now made the changes to the formatting requested by the editor.	
Additional Information:		
Question	Response	
Are you submitting this manuscript to a special series or article collection?	No	
Experimental design and statistics	Yes	
Full details of the experimental design and		

<p>statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

[Click here to view linked References](#)

Clustering trees: a visualisation for evaluating clusterings at multiple resolutions

Luke Zappia (1, 2)

Alicia Oshlack* (1, 2)

1 Bioinformatics, Murdoch Children's Research Institute; 2 School of Biosciences, University of Melbourne

*corresponding author

Corresponding author email: alicia.oshlack@mcri.edu.au

Corresponding author ORCID ID: 0000-0001-9788-5690

Clustering techniques are widely used in the analysis of large data sets to group together samples with similar properties. For example, clustering is often used in the field of single-cell RNA-sequencing in order to identify different cell types present in a tissue sample. There are many algorithms for performing clustering and the results can vary substantially. In particular, the number of groups present in a data set is often unknown and the number of clusters identified by an algorithm can change based on the parameters used. To explore and examine the impact of varying clustering resolution we present clustering trees. This visualisation shows the relationships between clusters at multiple resolutions allowing researchers to see how samples move as the number of clusters increases. In addition, meta-information can be overlaid on the tree to inform the choice of resolution and guide in identification of clusters. We illustrate the features of clustering trees using a series of simulations as well as two real examples, the classical iris dataset and a complex single-cell RNA-sequencing dataset. Clustering trees can be produced using the `clustree` R package available from CRAN (<https://CRAN.R-project.org/package=clustree>) and developed on GitHub (<https://github.com/lazappi/clustree>).

Keywords: Clustering - Visualisation - scRNA-seq

Introduction

Clustering analysis is commonly used to group similar samples across a diverse range of applications. Typically, the goal of clustering is to form groups of samples that are more similar to each other than to samples in other groups. While fuzzy or soft clustering approaches assign each sample to every cluster with some probability, and hierarchical clustering forms a tree of samples, most methods form hard clusters where each sample is assigned to a single group. This goal can be achieved in a variety of ways, such as by considering the distances between samples (e.g. k -

32 means [1–3], PAM [4]), areas of density across the dataset (e.g. DBSCAN [5]) or relationships to
1
2 33 statistical distributions [6].
3
4

5 34 In many cases the number of groups that should be present in a dataset is not known in advance
6
7 35 and deciding the correct number of clusters to use is a significant challenge. For some algorithms,
8
9 36 such as *k*-means clustering, the number of clusters must be explicitly provided. Other methods
10
11 37 have parameters that, directly or indirectly, control the clustering resolution and therefore the
12
13 38 number of clusters produced. While there are methods and statistics (such as the elbow method
14
15 39 [7] or silhouette plots [8]) designed to help analysts decide which clustering resolution to use,
16
17 40 they typically produce a single score which only considers a single set of samples or clusters at a
18
19
20
21 41 time.
22

23 42 An alternative approach would be to consider clusterings at multiple resolutions and examine
24
25 43 how samples change groupings as the number of clusters increases. This has led to a range of
26
27 44 cluster stability measures [9], many of which rely on clustering of perturbed or sub-sampled
28
29 45 datasets. For example, the model explorer algorithm sub-samples a dataset multiple times,
30
31 46 clusters each sub-sampled dataset at various resolutions and then calculates a similarity between
32
33 47 clusterings at the same resolution to give a distribution of similarities which can inform the choice
34
35 48 of resolution [10]. One cluster stability measure that isn't based on perturbations is that
36
37 49 contained in the SC3 package for clustering single-cell RNA-sequencing data [11]. Starting with a
38
39 50 set of cluster labels at different resolutions each cluster is scored, with clusters awarded increased
40
41 51 stability if they share the same samples as a cluster at another resolution, but penalised for being
42
43
44 52 at a higher resolution.
45
46
47

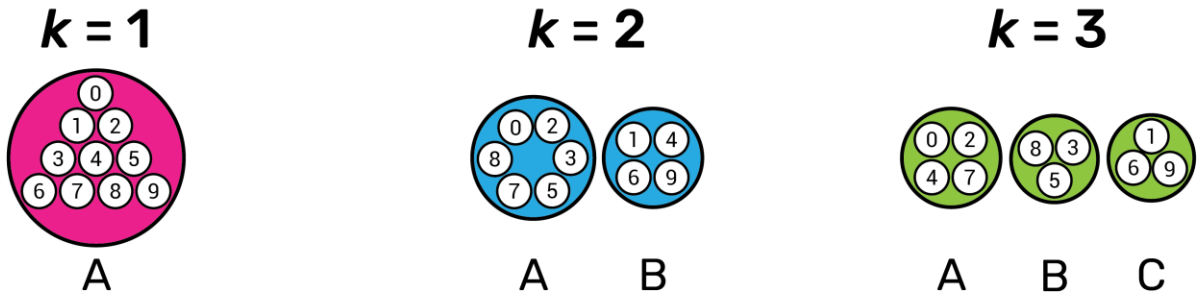
48 53 A similar simple approach is taken by the clustering tree visualisation we present here, without
49
50 54 calculating scores: (i) a dataset is clustered using any hard clustering algorithm at multiple
51
52 55 resolutions, producing sets of cluster nodes, (ii) the overlap between clusters at adjacent
53
54 56 resolutions is used to build edges, (iii) the resulting graph is presented as a tree. This tree can be
55
56 57 used to examine how clusters are related to each other, which clusters are distinct and which are
57
58 58 unstable. In the following sections we describe how we construct such a tree and present
59
60
61
62
63
64
65

59 examples of trees built from a classical clustering dataset and a complex single-cell RNA-
1
2 60 sequencing (scRNA-seq) dataset. The figures shown here can be produced in R using our publicly
3
4 61 available clustree package. Although clustering trees can not directly provide a clustering
5
6 62 resolution to use they can be a useful tool for exploring and visualising the range of possible
7
8
9 63 choices.

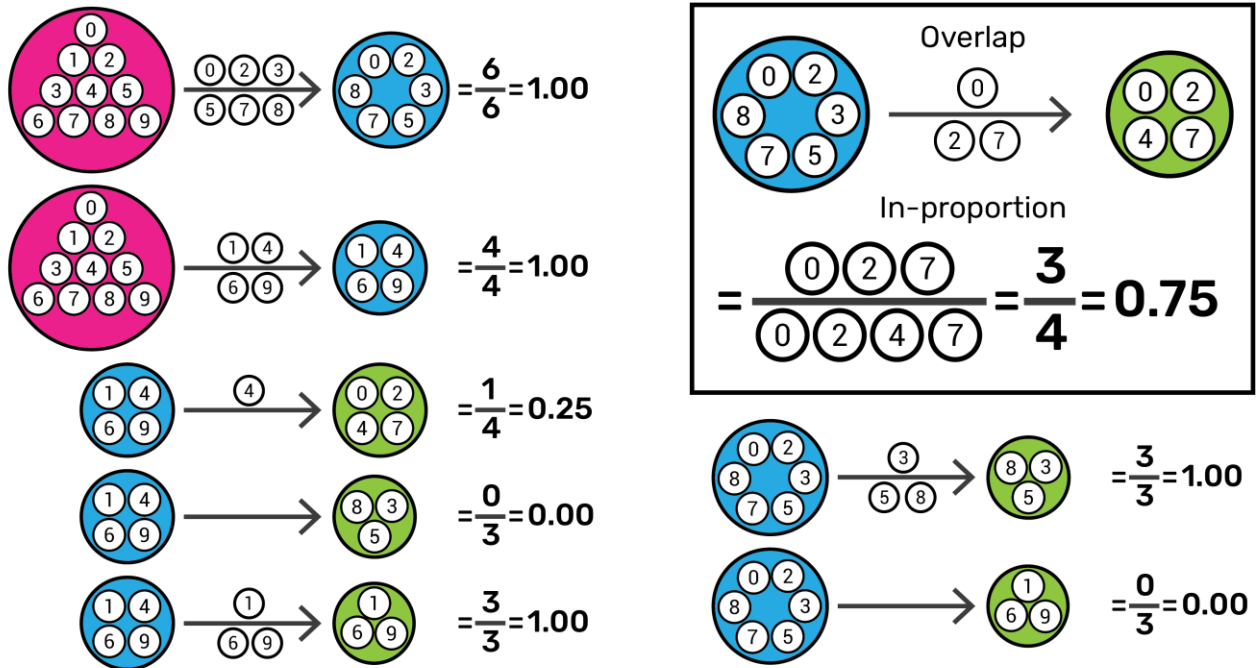
64 **Building a clustering tree**

10
11
12
13
14 65 To build a clustering tree, we start with a set of clusterings allocating samples to groups at several
15
16 66 different resolutions. These could be produced using any hard-clustering algorithm that allows
17
18 66 control of the number of clusters in some way. For example, this could be a set of samples
19
20 67 clustered using k -means with $k = 1, 2, 3$ as shown in Figure 1. We sort these clusterings so that
21
22 68 they are ordered by increasing resolution (k), then consider pairs of adjacent clusterings. Each
23
24 69 cluster $c_{k,i}$ (where $i = 1, \dots, n$ and n is the number of clusters at resolution k) is compared with
25
26 70 each cluster $c_{k+1,j}$ (where $j = 1, \dots, m$ and m is the number of clusters at resolution $k + 1$). The
27
28 71 overlap between the two clusters is computed as the number of samples that are assigned to both
29
30 72 $c_{k,i}$ and $c_{k+1,j}$. We next build a graph where each node is a cluster and each edge is an overlap
31
32 73 between two clusters. While we refer to this graph as a tree in this paper for simplicity it can more
33
34 74 correctly be described as a polytree, a special case of a directed acyclic graph where the
35
36 75 underlying undirected graph is a tree [12].
37
38 76
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1. Cluster at multiple resolutions



2. Find overlaps and calculate in-proportion



3. Filter edges and visualise tree

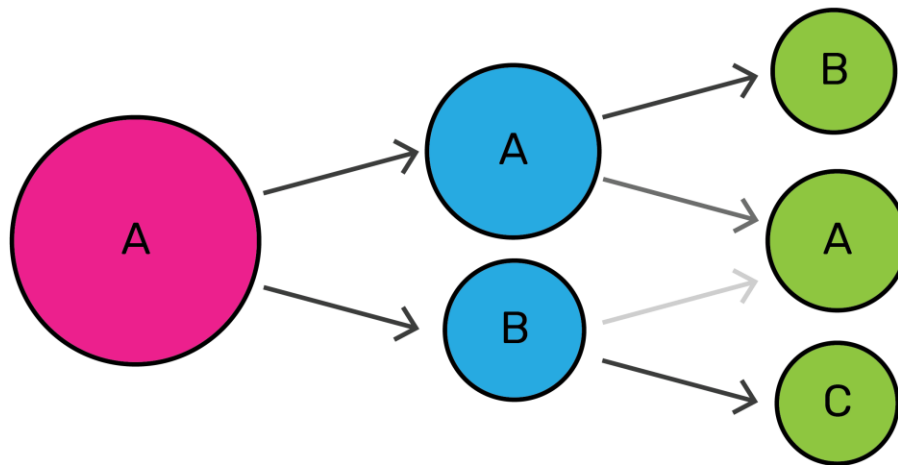


Figure 1 Illustration of the steps required to build a clustering tree. First a dataset must be clustered at different resolutions. The overlap in samples between clusters at adjacent resolutions is computed and used to calculate the in-proportion for each edge. Finally the edges are filtered and the graph visualised as a tree.

82 Many of the edges will be empty, for example in Figure 1 no samples in Cluster A at $k = 2$ end up
1
2 83 in Cluster B at $k = 3$. In some datasets there may also be edges that contain few samples. These
3
4 84 edges are not informative and result in a cluttered tree. An obvious solution for removing
5
6 85 uninformative, low-count edges is to filter them using a threshold on the number of samples they
7
8 86 represent. However, in this case the count of samples is not the correct statistic to use because it
9
10
11 87 favours edges at lower resolutions and those connecting larger clusters. Instead we define the in-
12
13 88 proportion metric as the ratio between the number of samples on the edge and the number of
14
15 89 samples in the cluster it goes towards. This metric shows the importance of the edge to the higher
16
17 90 resolution cluster independently of the cluster size. We can then apply a threshold to the in-
18
19
20 91 proportion in order to remove less informative edges.

21
22
23 92 The final graph can then be visualised. In theory any graph layout algorithm could be used but for
24
25 93 the clustree package we have made use of the two algorithms specifically designed for tree
26
27 94 structures available in the igraph package [13]. These are the Reingold-Tilford tree layout, which
28
29 95 places parent nodes above their children [14], and the Sugiyama layout which places nodes of a
30
31 96 directed acyclic graph in layers while minimising the number of crossing edges [15]. Both of these
32
33
34 97 algorithms can produce attractive layouts and as such we have not found the need to design a
35
36 98 specific layout algorithm for clustering trees. By default the clustree package uses only a subset of
37
38 99 edges when constructing a layout, specifically the highest in-proportion edges for each node. We
39
40
41 100 have found that this often leads to more interpretable visualisations, however users can choose to
42
43 101 use all edges if desired.

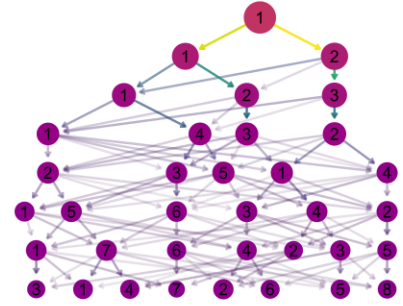
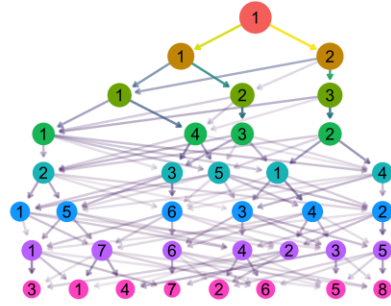
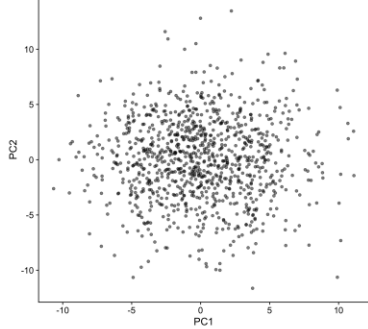
44
45
46 102 Whichever layout is used the final visualisation places the cluster nodes in a series of layers where
47
48 103 each layer is a different clustering resolution and edges show the transition of samples through
49
50 104 those resolutions. Edges are coloured according to the number of samples they represent and the
51
52 105 in-proportion metric is used to control the edge transparency, highlighting more important edges.
53
54
55 106 By default, the size of nodes is adjusted according to the number of samples in the cluster and
56
57 107 their colour indicates the clustering resolution. The clustree package also includes options for
58
59
60
61
62
63
64
65

108 controlling the aesthetics of nodes based on the attributes of samples in the clusters they
1
2109 represent as shown in the following examples.
3
4
5110 While a clustering tree is conceptually similar to the tree produced through hierarchical clustering
6
7111 there are some important differences. The most obvious are that a hierarchical clustering tree is
8
9112 the result of a particular clustering algorithm and shows the relationships between individual
10
11113 samples while the clustering trees described here are independent of clustering method and show
12
13114 relationships between clusters. The branches of a hierarchical tree show how the clustering
14
15115 algorithm has merged samples. In contrast, edges in a clustering tree show how samples move
16
17116 between clusters as the resolution changes and nodes may have multiple parents. While it is
18
19117 possible to overlay information about samples on a hierarchical tree this is not commonly done
20
21
22
23118 but is a key feature of the clustree package and how clustering trees could be used in practice.
24
25
26

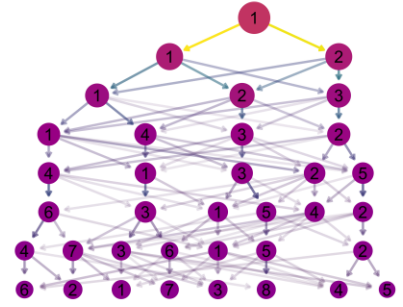
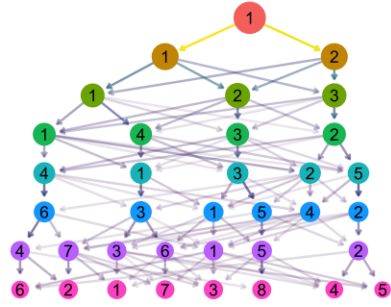
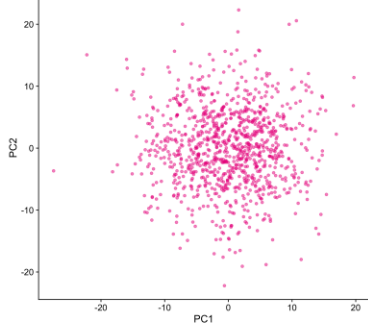
27119 **A demonstration using simulations**

28
29
30120 To demonstrate what a clustering tree can look like in different situations and how it behaves as a
31
32121 dataset is over-clustered we present some illustrative examples using simple simulations (see
33
34122 methods). We present five scenarios: random uniform noise (Simulation A), a single cluster
35
36123 (Simulation B), two clusters (Simulation C), three clusters (Simulation D) and four clusters
37
38124 (Simulation E). Each cluster consists of 1000 samples (points) generated from a 100 dimensional
39
40
41125 normal distribution and each synthetic dataset has been clustered using k -means clustering with
42
43126 $k = 1, \dots, 8$. We then use the clustree package to produce clustering trees for each dataset (Figure
44
45127 2).
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

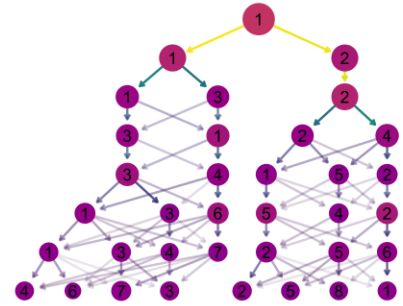
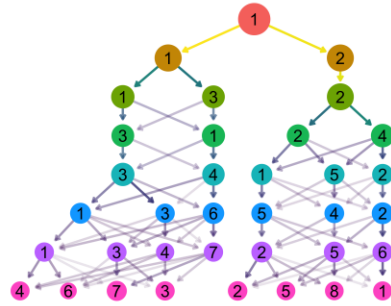
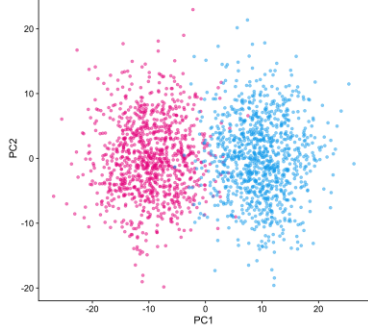
A - Uniform noise



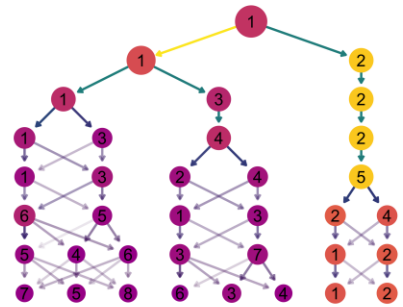
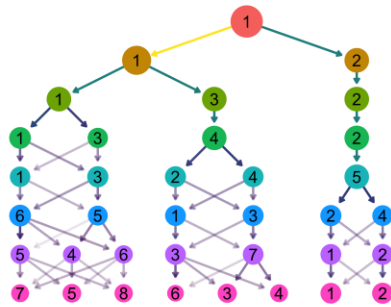
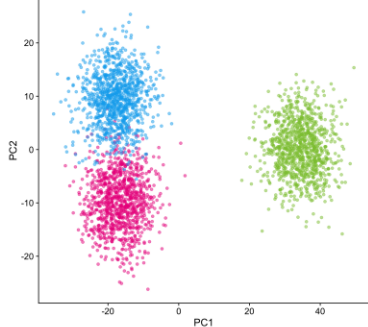
B - Single cluster



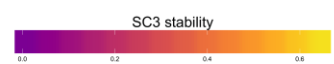
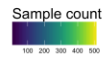
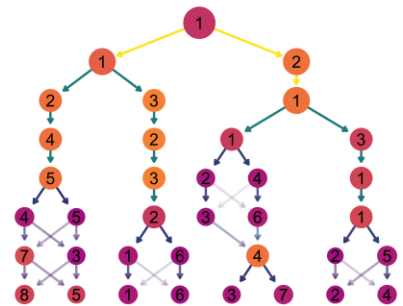
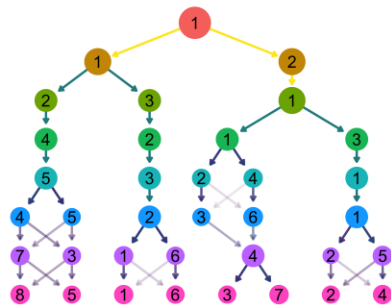
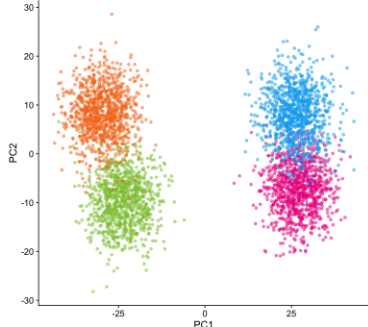
C - Two clusters



D - Three clusters



E - Four clusters



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

129 *Figure 2 Five synthetic datasets used to demonstrate clustering trees. For each dataset a scatter*
1130 *plot of the first two principal components, a default clustering tree and a clustering tree with*
2131 *nodes coloured by the SC3 stability index from purple (lowest) to yellow (highest) are shown.*
3132 *The five datasets contain: A) random uniform noise, B) a single cluster, C) two clusters, D) three*
4133 *clusters and E) four clusters.*

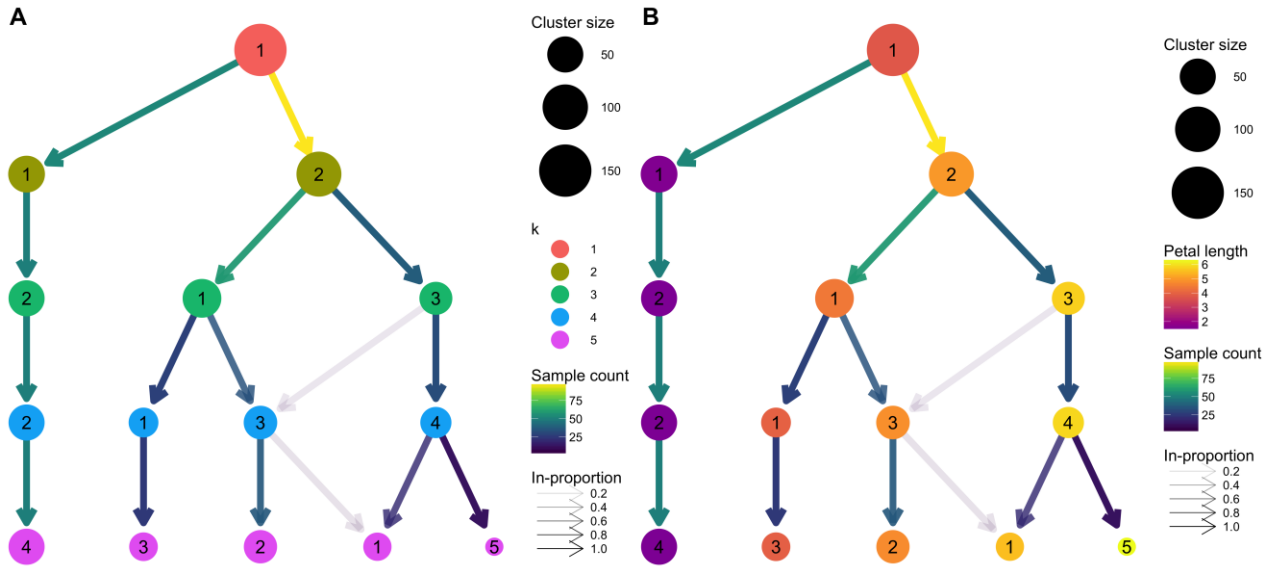
6134 Looking at the first two examples (uniform noise (Figure 2A) and a single cluster (Figure 2B)) we
7
8135 can clearly see how a clustering tree behaves when a clustering algorithm returns more clusters
9
10136 than are truly present in a dataset. New clusters begin to form from multiple existing clusters and
11
12137 many samples switch between branches of the tree resulting in low in-proportion edges. Unstable
13
14138 clusters may also appear then disappear as the resolution increases as seen in Figure 2E. As we
15
16139 add more structure to the datasets the clustering trees begin to form clear branches and low in-
17
18140 proportion edges tend to be confined to sections of the tree. By looking at which clusters are
19
20141 stable and where low in-proportion edges arise we can infer which areas of the tree are likely to be
21
22142 the result of true clusters and which are caused by over-clustering.

25
26
27143 The second clustering tree for each dataset shows nodes coloured according to the SC3 stability
28
29144 index for each cluster. As we would expect in the first two examples no cluster receives a high
30
31145 stability score. However, while we clearly see two branches in the clustering tree for the two
32
33146 cluster example (Simulation C) this is not reflected in the SC3 scores. No cluster receives a high
34
35147 stability score, most likely due to the high number of samples moving between clusters as the
36
37148 resolution increases. As there are more true clusters in the simulated datasets the SC3 stability
38
39149 scores become more predictive of the correct resolution to use, however it is important to look at
40
41
42150 the stability scores of all clusters at a particular resolution as taking the highest individual cluster
43
44151 stability score could lead to the incorrect resolution being used, as can be seen in the four cluster
45
46152 example (Simulation E). These examples show how clustering trees can be used to display
47
48
49153 existing clustering metrics in a way that can help to inform parameter choices.

54154 **A simple example**

55
56155 To further illustrate how a clustering tree is built, we will work through an example using the
57
58156 classical iris dataset [16]. This dataset contains measurements of the sepal length, sepal width,
59
60157 petal length and petal width from 150 iris flowers, 50 from each of three species: *Iris setosa*, *Iris*

158 *versicolor* and *Iris virginica*. The iris dataset is commonly used as example for both clustering
 1
 2159 and classification problems with the *Iris setosa* samples being significantly different to, and
 3
 4160 linearly separable from, the other samples. We have clustered this dataset using k -means
 5
 6161 clustering with $k = 1, \dots, 5$ and produced the clustering tree shown in Figure 3A.



162
 163
 164
 165
 166
 167
 168
 169
 170
 171
 172
 173
 174
 175
 176
 177
 178
 179
 180
 181
 182
 183
 184
 185
 186
 187
 188
 189
 190
 191
 192
 193
 194
 195
 196
 197
 198
 199
 200
 201
 202
 203
 204
 205
 206
 207
 208
 209
 210
 211
 212
 213
 214
 215
 216
 217
 218
 219
 220
 221
 222
 223
 224
 225
 226
 227
 228
 229
 230
 231
 232
 233
 234
 235
 236
 237
 238
 239
 240
 241
 242
 243
 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263
 264
 265
 266
 267
 268
 269
 270
 271
 272
 273
 274
 275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323
 324
 325
 326
 327
 328
 329
 330
 331
 332
 333
 334
 335
 336
 337
 338
 339
 340
 341
 342
 343
 344
 345
 346
 347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377
 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
 432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500

We see that there is one branch of the tree that is clearly distinct (presumably representing *Iris setosa*), remaining unchanged regardless of the number of clusters. On the other side we see the cluster at $k = 2$ cleanly splits into two clusters (presumably *Iris versicolor* and *Iris virginica*) at $k = 3$ but as we move to $k = 4$ and $k = 5$ we see clusters being formed from multiple branches with more low in-proportion edges. As we have seen in the simulated examples, this kind of pattern can indicate that the data has become over-clustered and we have begun to introduce artificial groupings.

We can check our assumption that the distinct branch represents the *Iris setosa* samples and the other two clusters at $k = 3$ are *Iris versicolor* and *Iris virginica* by overlaying some known information about the samples. In Figure 3B we have coloured the nodes by the mean petal length

180 of the samples they contain. We can now see that clusters in the distinct branch have the shortest
1 petals, with Cluster 1 at $k = 3$ having an intermediate length and Cluster 3 the longest petals. This
2181 feature is known to separate the samples into the expected species with *Iris setosa* having the
3
4182 shortest petals on average, *Iris versicolor* an intermediate length and *Iris virginica* the longest.
5
6183
7
8
9184 Although this is a very simple example it highlights some of the benefits of viewing a clustering
10 tree. We get some indication of the correct clustering resolution by examining the edges and we
11185
12 can overlay known information to assess the quality of the clustering. For example, if we observed
13
14186
15 that all clusters had the same mean petal length it would suggest that the clustering has not been
16187
17 successful as we know this is an important feature that separates the species. We could potentially
18188
19 learn more by looking at which samples follow low proportion edges or overlaying a series of
20
21189
22 features to try and understand what causes particular clusters to split.
23190
24
25
26

27191 **Clustering trees for single-cell RNA-seq data**

28192 One field that has begun to make heavy use of clustering techniques is the analysis of single-cell
29 RNA-sequencing (scRNA-seq) data. Single-cell RNA-sequencing is a recently developed
30193
31 technology that can measure how genes are expressed in thousands to millions of individual cells
32194
33 [18]. This technology has been rapidly adopted in fields like developmental biology and
34195
35 immunology where it is valuable to have information from single cells rather than measurements
36196
37 that are averaged across the many different cells in a sample using older RNA sequencing
38197
39 technologies. One of the key uses for scRNA-seq is to discover and interrogate the different cell
40
41198
42 types present in a sample of a complex tissue. In this situation, clustering is typically used to
43199
44 group similar cells based on their gene expression profiles. Differences in gene expression
45200
46 between groups can then be used to infer the identity or function of those cells [19]. The number
47201
48 of cell types (clusters) in an scRNA-seq dataset can vary depending on factors such as the tissue
49
50202
51 being studied, its developmental or environmental state and the number of cells captured. Often
52
53 the number of cells types is not known before the data is generated and some samples can contain
54203
55
56204
57
58
59
60
61
62
63
64
65

dozens of clusters. Therefore, deciding which clustering resolution to use is an important consideration in this application.

As an example of how clustering trees can be used in the scRNA-seq context we consider a commonly used Peripheral Blood Mononuclear Cell (PBMC) dataset. This dataset was originally produced by 10x Genomics and contains 2700 peripheral blood mononuclear cells, representing a range of well-studied immune cell types [20]. We have analysed this dataset using the Seurat package [21], a commonly used toolkit for scRNA-seq analysis, following the instructions in their tutorial with the exception of varying the clustering resolution parameter from zero to five (see methods). Seurat uses a graph-based clustering algorithm and the resolution parameter controls the partitioning of this graph, with higher values resulting in more clusters. The clustering trees produced from this analysis are shown in Figure 4.

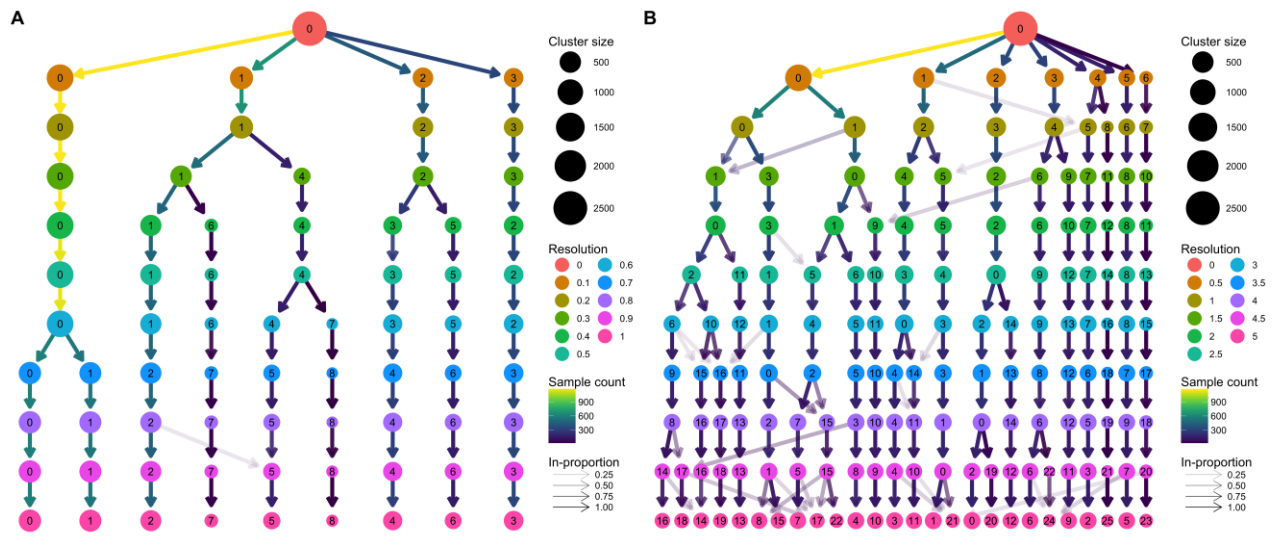


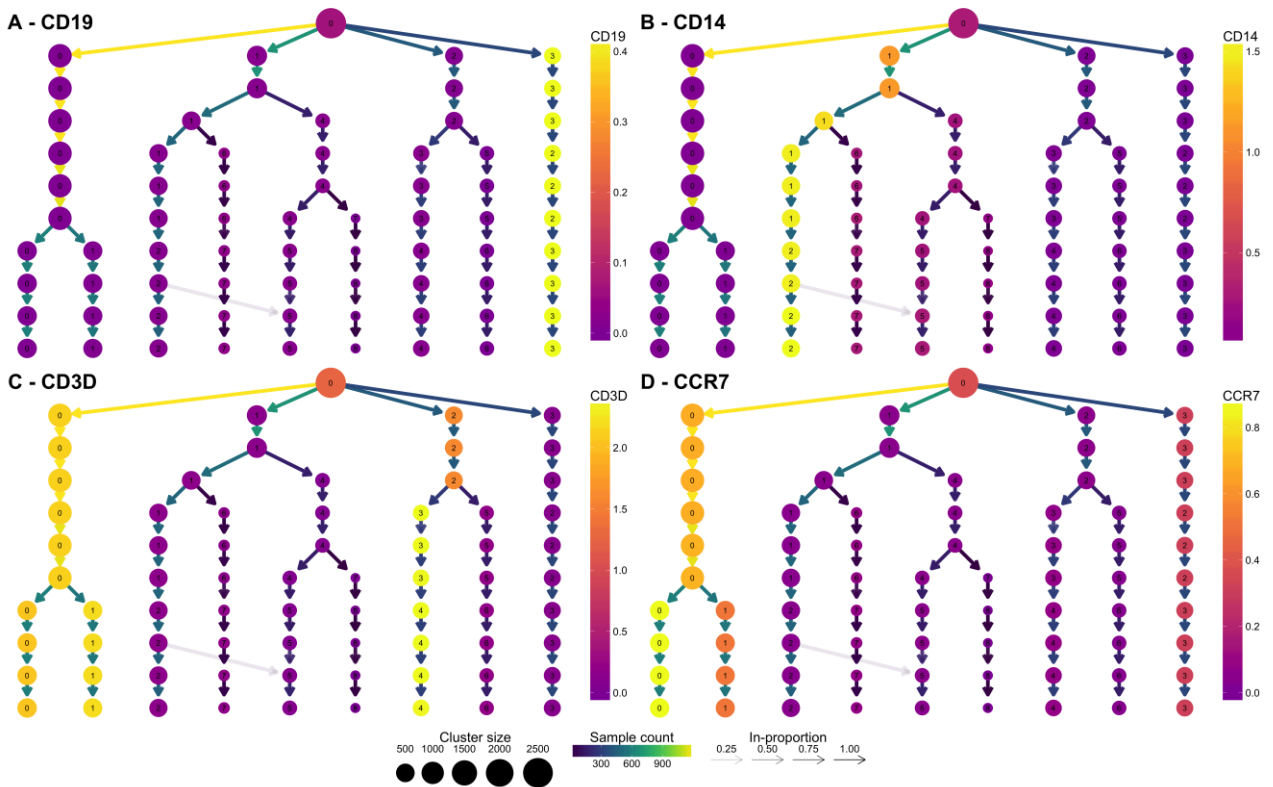
Figure 4 Two clustering trees of a dataset of 2700 Peripheral Blood Mononuclear Cells (PBMCs). A) results from clustering using Seurat with resolution parameters from zero to one. At a resolution of 0.1 we see the formation of four main branches, one of which continues to split up to a resolution of 0.4, after which there are only minor changes. B) resolutions from zero to five. At the highest resolutions we begin to see many low in-proportion edges indicating cluster instability. Seurat labels clusters according to their size with Cluster 0 being the largest.

The clustering tree covering resolutions zero to one in steps of 0.1 (Figure 4A) shows that four main branches form at a resolution of just 0.1. One of these branches, starting with Cluster 3 at resolution 0.1, remains unchanged while the branch starting with Cluster 2 splits only once at a resolution of 0.4. Most of the branching occurs in the branch starting with Cluster 1 which

consistently has sub-branches split off to form new clusters as the resolution increases. There are two regions of stability in this tree; at resolution 0.4-0.5 and resolution 0.7-1.0 where the branch starting at Cluster 0 splits in two.

Figure 4B shows a clustering tree with a greater range of resolutions, from zero to five in steps of 0.5. By looking across this range we can see what happens when the algorithm is forced to produce more clusters than are likely to be truly present in this dataset. As over-clustering occurs we begin to see more low in-proportion edges and new clusters forming from multiple parent clusters. This suggests that those areas of the tree are unstable and that the new clusters being formed are unlikely to represent true groups in the dataset.

Known marker genes are commonly used to identify the cell types that specific clusters correspond to. Overlaying gene expression information onto a clustering tree provides an alternative view that can help to indicate when clusters containing pure cell populations are formed. Figure 5 shows the PBMC clustering tree in Figure 4A overlaid with the expression of some known marker genes.



242 *Figure 5 Clustering trees of the PBMC dataset coloured according to the expression of known*
243 *markers. The node colours indicate the average of the \log_2 gene counts of samples in each*
244 *cluster. CD19 (A) identifies B cells, CD14 (B) shows a population of monocytes, CD3D (C) is a*
245 *marker of T cells and CCR7 (D) shows the split between memory and naive CD4 T cells.*

246 By adding this extra information, we can quickly identify some of the cell types. CD19 (Figure 5A)
247 is a marker of B cells and is clearly expressed in the most distinct branch of the tree. CD14 (Figure
248 5B) is a marker of a type of monocyte, which becomes more expressed as we follow one of the
249 central branches, allowing us to see which resolution identifies a pure population of these cells.
250 CD3D (Figure 5C) is a general marker of T cells and is expressed in two separate branches, one
251 which splits into low and high expression of CCR7 (Figure 5D), separating memory and naive
252 CD4 T cells. By adding expression of known genes to a clustering tree, we can see if more
253 populations can be identified as the clustering resolution is increased and if clusters are
254 consistent with known biology. For most of the Seurat tutorial a resolution of 0.6 is used, but the
255 authors note that by moving to a resolution of 0.8, a split can be achieved between memory and
256 naive CD4 T cells. This is a split that could be anticipated by looking at the clustering tree with the
257 addition of prior information.

258 **Discussion and conclusion**

259 Clustering similar samples into groups is a useful technique in many fields, but often analysts are
260 faced with the tricky problem of deciding which clustering resolution to use. Traditional
261 approaches to this problem typically consider a single cluster or sample at a time and may rely on
262 prior knowledge of sample labels. Here we present clustering trees, an alternative visualisation
263 that shows the relationships between clusterings at multiple resolutions. While clustering trees
264 cannot directly suggest which clustering resolution to use they can be a useful tool for helping to
265 make that decision, particularly when combined with other metrics or domain knowledge.

266 Clustering trees display how clusters are divided as resolution increases, which clusters are clearly
267 separate and distinct, which are related to each other and how samples change groups as more
268 clusters are produced. Although clustering trees can appear similar to the trees produced from
269 hierarchical clustering there are several important differences. Hierarchical clustering considers

270 the relationships between individual samples and doesn't provide an obvious way to form groups.
1
2271 In contrast, clustering trees are independent of any particular clustering method and show the
3
4272 relationships between clusters, rather than samples, at different resolutions, any of which could
5
6273 be used for further analysis.
8
9
10274 To illustrate the uses of clustering trees we presented a series of simulations and two examples of
11
12275 real analyses, one using the classical iris dataset and a second based on a complex scRNA-seq
13
14276 dataset. Both examples demonstrate how a clustering tree can help inform the decision of which
15
16277 resolution to use and how overlaying extra information can help to validate those clusters. This is
17
18278 of particular use to scRNA-seq analysis as these datasets are often large, noisy and contain an
19
20
21279 unknown number of cell types or clusters.
22
23280 Even when determining the number of clusters is not a problem, clustering trees can be a valuable
24
25
26281 tool. They provide a compact, information dense, visualisation that can display summarised
27
28282 information across a range of clusters. By modifying the appearance of cluster nodes based on
29
30283 attributes of the samples they represent, clusterings can be evaluated and identities of clusters
31
32284 established. Clustering trees potentially have applications in many fields and in the future could
33
34285 be adapted to be more flexible, such as by accommodating fuzzy clusterings. There may also be
35
36
37286 uses for more general clustering graphs to combine results from multiple sets of parameters or
38
39287 clustering methods.
40
41
42
43

44288 **Methods**

45 46289 **clustree**

47
48290 The clustree software package (v0.2.0) is built for the R statistical programming language
49
50291 (v3.5.0). It relies on the ggraph package (v1.0.1) [22], which is itself built on the ggplot2 (v2.2.1)
51
52292 [23] and tidygraph (v1.1.0) [24] packages. Clustering trees are displayed using the Reingold-
53
54
55293 Tilford tree layout or the Sugiyama layout, both available as part of the igraph package (v1.2.1).
56
57
58294 Figure panels shown here were produced using the cowplot package (v0.9.2) [25].
59
60
61
62
63
64
65

295 Simulations

1
296 Simulated datasets were constructed by generating points from statistical distributions. The first
3
4297 simulation (Simulation A) consists of 1000 points randomly generated from a 100 dimensional
5
6298 space using a uniform distribution between zero and 10. Simulation B consists of a single
7
8299 normally distributed cluster of 1000 points in 100 dimensions. The centre of this cluster was
9
10
11300 chosen from a normal distribution with mean zero and standard deviation 10. Points were then
12
13301 generated around this centre from a normal distribution with mean equal to the centre point and
14
15302 a standard deviation of five. The remaining three simulations were produced by adding additional
16
17303 clusters. In order to have a known relationship between clusters the centre for the new clusters
18
19304 was created by manipulating the centres of existing clusters. For Cluster 2 a random 100
20
21
22305 dimensional vector was generated from a normal distribution with mean zero and standard
23
24306 deviation two and added to the centre for Cluster 1. Centre 3 was the average of Centre 1 and
25
26307 Centre 2 plus a random vector from a normal distribution with mean zero and standard deviation
27
28308 five. To ensure a similar relationship between clusters 3 and 4 as between clusters 1 and 2, Centre
29
30
31309 4 was produced by adding half the vector used to produce Centre 2 to Centre 3 plus another
32
33310 vector from a normal distribution with mean zero and standard deviation two. Points for each
34
35311 cluster were generated in the same way as for Cluster 1. Simulation C consists of the points in
36
37312 clusters 1 and 2, Simulation D consists of clusters 1, 2 and 3, Simulation E consists of clusters 1, 2,
38
39
40313 3 and 4. Each simulated dataset was clustered using the “kmeans” function in the stats package
41
42314 with values of k from one to eight, a maximum of 100 iterations and 10 random starting positions.
43
44315 The clustering tree visualisations were produced using the clustree package with the tree layout.
45
46316 The simulated datasets and the code use to produce them are available from the repository for
47
48
49317 this paper [26].

51 Iris dataset

52
53
54319 The iris dataset is available as part of R. We clustered this dataset using the “kmeans” function in
55
56320 the stats package with values of k from one to five. Each value of k was clustered with a maximum
57
58
59321 of 100 iterations and with 10 random starting positions. The clustree package was used to
60
61
62
63
64
65

322 visualise the results using the Sugiyama layout. The clustered iris dataset is available as part of
1
323 the clustree package.

324 **PBMC dataset**

325 The PBMC dataset was downloaded from the Seurat tutorial page [27] and this tutorial was
326 followed for most of the analysis using Seurat version 2.3.1. Briefly cells were filtered based on the
327 number of genes they express and the percentage of counts assigned to mitochondrial genes. The
328 data was then log-normalised and 1838 variable genes identified. Potential confounding variables
329 (number of unique molecular identifiers and percentage mitochondrial expression) were
330 regressed from the dataset before performing principal component analysis on the identified
331 variable genes. The first 10 principal components were then used to build a graph which was
332 partitioned into clusters using Louvain modularity optimisation [28] with resolution parameters
333 in the range zero to five, in steps of 0.1 between zero and one and then in steps of 0.5. Clustree
334 was then used to visualise the results using the tree layout.

335 **Availability of source code and requirements**

336 Project name: clustree

337 Project home page: <https://github.com/lazappi/clustree>

338 Operating systems(s): Linux, MacOS, Windows

339 Programming language: R (≥ 3.4)

340 Other requirements: None

341 License: GPL-3

342 Any restrictions to use by non-academics: None

343 RRID: SCR_016293

344

345 **Availability of data and materials**

1
2 346 The clustree package is available from CRAN [29] and is being developed on GitHub [30]. The
3
4
5 347 code and datasets used for the analysis in this paper are also available from GitHub [26]. The
6
7 348 clustered iris dataset is included as part of clustree and the PBMC dataset can be downloaded
8
9 349 from the Seurat tutorial page [27] or the paper GitHub repository. Snapshots of the code are
10
11 350 available in the GigaScience repository, GigaDB [31].
12
13
14
15

16 351 **Declarations**

17 18 352 **Abbreviations**

19
20 353 PBMC: peripheral blood mononuclear cell
21
22

23 354 scRNA-seq: single-cell RNA-sequencing
24
25

26 27 355 **Consent for publication**

28
29 356 Not applicable.
30
31

32 357 **Competing interests**

33
34 358 The authors declare that they have no competing interests.
35
36
37

38 359 **Funding**

39
40 360 Luke Zappia is supported by an Australian Government Research Training Program (RTP)

41
42 361 Scholarship. Alicia Oshlack is supported through a National Health and Medical Research Council

43
44 362 Career Development Fellowship APP1126157. MCRI is supported by the Victorian Government's

45
46 363 Operational Infrastructure Support Program.
47
48

49 364 **Author's contributions**

50
51 365 LZ designed the clustering tree algorithm, wrote the clustree software package and drafted the

52
53 366 manuscript. AO supervised the project and commented on the manuscript.
54
55
56
57
58
59
60
61
62
63
64
65

Acknowledgements

Thank you to Marek Cmero for providing comments on a draft of the manuscript and the reviewers for their comments and suggestions.

References

1. Forgy WE. Cluster analysis of multivariate data : efficiency versus interpretability of classifications. *Biometrics* 1965;21:768–9. Available from: <https://ci.nii.ac.jp/naid/10009668881/>
2. Macqueen J. Some methods for classification and analysis of multivariate observations. In 5th Berkeley Symposium on Mathematical Statistics and Probability 1967. Available from: <http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.308.8619>
3. Lloyd S. Least squares quantization in PCM. *IEEE Trans Inf Theory*. 1982;28:129–37. Available from: <http://dx.doi.org/10.1109/TIT.1982.1056489>
4. Kaufman L, Rousseeuw PJ. Partitioning Around Medoids (Program PAM). *Finding Groups in Data*. John Wiley & Sons, Inc. 1990. pp. 68–125. Available from: <http://dx.doi.org/10.1002/9780470316801.ch2>
5. Ester M, Kriegel H-P, Sander J, Xu X. A Density-based Algorithm for Discovering Clusters a Density-based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*. Portland, Oregon: AAAI Press; 1996. pp. 226–31. Available from: <http://dl.acm.org/citation.cfm?id=3001460.3001507>
6. Fraley C, Raftery AE. Model-Based Clustering, Discriminant Analysis, and Density Estimation. *J Am Stat Assoc*. 2002;97:611–31. Available from: <http://www.tandfonline.com/doi/abs/10.1198/016214502760047131>
7. Thorndike RL. Who belongs in the family? *Psychometrika*. Springer-Verlag; 1953;18:267–76. Available from: <https://link.springer.com/article/10.1007/BF02289263>
8. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J Comput Appl Math*. 1987;20:53–65. Available from: <http://www.sciencedirect.com/science/article/pii/0377042787901257>
9. Luxburg U von. Clustering Stability: An Overview. *Foundations and Trends in Machine Learning*. Now Publishers; 2010;2:235–74. Available from: <http://dx.doi.org/10.1561/22000000008>
10. Ben-Hur A, Elisseeff A, Guyon I. A stability based method for discovering structure in clustered data. *Pac Symp Biocomput*. 2002;6–17. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/11928511>
11. Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, et al. SC3: consensus clustering of single-cell RNA-seq data. *Nat Methods*. 2017;14:483–6. Available from: <http://dx.doi.org/10.1038/nmeth.4236>
12. Rebane G, Pearl J. The Recovery of Causal Poly-Trees from Statistical Data. 2013; Available from: <http://arxiv.org/abs/1304.2736>

- 406 13. Csardi G, Nepusz T. The igraph software package for complex network research. *InterJournal,*
407 *Complex Systems.* 2006;1695:1–9.
- 2
3408 14. Reingold EM, Tilford JS. Tidier Drawings of Trees. *IEEE Trans Software Eng.* 1981;SE-7:223–
409 8. Available from: <http://dx.doi.org/10.1109/TSE.1981.234519>
- 5
6410 15. Sugiyama K, Tagawa S, Toda M. Methods for Visual Understanding of Hierarchical System
7411 Structures. *IEEE Trans Syst Man Cybern.* 1981;11:109–25. Available from:
8412 <http://dx.doi.org/10.1109/TSMC.1981.4308636>
- 9
10413 16. Anderson E. The Irises of the Gaspe Peninsula. *Bulletin of the American Iris Society.*
11414 1935;59:2–5.
- 13
14415 17. Fisher RA. The use of multiple measurements in taxonomic problems. *Ann Eugen.* Blackwell
15416 Publishing Ltd; 1936;7:179–88. Available from: [http://dx.doi.org/10.1111/j.1469-](http://dx.doi.org/10.1111/j.1469-1809.1936.tb02137.x)
16417 [1809.1936.tb02137.x](http://dx.doi.org/10.1111/j.1469-1809.1936.tb02137.x)
- 17
18418 18. Tang F, Barbacioru C, Wang Y, Nordman E, Lee C, Xu N, et al. mRNA-Seq whole-
19419 transcriptome analysis of a single cell. *Nat Methods.* 2009;6:377–82. Available from:
20420 <http://dx.doi.org/10.1038/nmeth.1315>
- 21
22421 19. Stegle O, Teichmann SA, Marioni JC. Computational and analytical challenges in single-cell
23422 transcriptomics. *Nat Rev Genet.* Nature Publishing Group; 2015;16:133–45. Available from:
24423 <http://dx.doi.org/10.1038/nrg3833>
- 25
26424 20. Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel
27425 digital transcriptional profiling of single cells. *Nat Commun.* 2017;8:14049. Available from:
28426 <http://dx.doi.org/10.1038/ncomms14049>
- 29
30427 21. Satija R, Farrell JA, Gennert D, Schier AF, Regev A. Spatial reconstruction of single-cell gene
31428 expression data. *Nat Biotechnol.* Nature Publishing Group; 2015;33:495–502. Available from:
32429 <http://dx.doi.org/10.1038/nbt.3192>
- 33
34430 22. Pedersen TL. ggraph: An Implementation of Grammar of Graphics for Graphs and Networks.
35431 2018. Available from: <https://CRAN.R-project.org/package=ggraph>
- 36
37432 23. Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer New York; 2010.
- 38
39433 24. Pedersen TL. tidygraph: A Tidy API for Graph Manipulation. 2018. Available from:
40434 <https://CRAN.R-project.org/package=tidygraph>
- 41
42435 25. Wilke CO. cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'. 2017.
43436 Available from: <https://CRAN.R-project.org/package=cowplot>
- 44
45437 26. Zappia L, Oshlack A. clustree-paper. GitHub; 2018. Available from:
46438 <https://github.com/Oshlack/clustree-paper>
- 48
49439 27. Satija Lab. Seurat PBMC3K tutorial. https://satijalab.org/seurat/pbmc3k_tutorial.html;
50440 2018. Available from: https://satijalab.org/seurat/pbmc3k_tutorial.html
- 51
52441 28. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E. Fast unfolding of communities in large
53442 networks. *J Stat Mech.* IOP Publishing; 2008;2008:P10008. Available from:
54443 <http://iopscience.iop.org/article/10.1088/1742-5468/2008/10/P10008/meta>
- 55
56444 29. Zappia L, Oshlack A. clustree: Visualise Clusterings at Different Resolutions. 2018. Available
57445 from: <https://CRAN.R-project.org/package=clustree>
- 58
59446 30. Zappia L, Oshlack A. clustree. GitHub; 2018. Available from:
60447 <https://github.com/lazappi/clustree>
- 61
62
63
64
65

448 31. Zappia L, Oshlack A: Supporting data for "Clustering trees: a visualisation for evaluating
449 clusterings at multiple resolutions" GigaScience Database. 2018. <http://dx.doi.org/100478>

2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65