

# Plotting clustering trees

Luke Zappia

Last updated: 24 February 2019

- [1 What is a clustering tree?](#)
- [2 A simple example](#)
  - [2.1 The data](#)
  - [2.2 Plotting a tree](#)
  - [2.3 Controlling aesthetics](#)
  - [2.4 Layout](#)
  - [2.5 Adding labels](#)
- [3 Clustering trees for scRNA-seq data](#)
  - [3.1 SingleCellExperiment objects](#)
  - [3.2 Seurat objects](#)
  - [3.3 Using genes as aesthetics](#)
- [4 Overlaying clustering trees](#)
  - [4.1 Choosing what to colour](#)
  - [4.2 Labelling nodes](#)
  - [4.3 Showing the side view](#)
- [5 Modifying appearance](#)
  - [5.1 Legends](#)
- [6 Citing clustree](#)
- [References](#)

## 1 What is a clustering tree?

Clustering analysis is used in many contexts to group similar samples. One problem when conducting this kind of analysis is how many clusters to use. This is usually controlled by a parameter provided to the clustering algorithm, such as  $k$  for  $k$ -means clustering.

Statistics designed to help you make this choice typically either compare two clusterings or score a single clustering. A clustering tree is different in that it visualises the relationships between at a range of resolutions.

To build a clustering tree we need to look at how cells move as the clustering resolution is increased. Each cluster forms a node in the tree and edges are constructed by considering the cells in a cluster at a lower resolution (say  $k = 2$ ) that end up in a cluster at the next highest resolution (say  $k = 3$ ). By connecting clusters in this way we can see how clusters are related to each other, which are clearly distinct and which are unstable. Extra information about the cells in each node can also be overlaid in order to help make the decision about which resolution to use. For more information about clustering trees please refer to our associated publication (Zappia and Oshlack 2018).

## 2 A simple example

To demonstrate what a clustering tree looks like we will work through a short example using the well known `iris` dataset.

### 2.1 The data

The `iris` dataset consists of measurements (sepal length, sepal width, petal length and petal width) of 150 iris flowers, 50 from each of three species (**Iris setosa**, **Iris versicolor** and **Iris virginica**). For more information see `?iris`. We are going to use a version of this dataset that has already been clustered. Let's load the data and take a look:

```
library(clustree)
3∞ Loading required package: ggraph
3∞ Loading required package: ggplot8
data("iris_clusts")

head(iris_clusts)
3∞   Sepal.Length Sepal.Width Petal.Length Petal.Width Species K7 K8 K9 K0 K0
3∞ 7           0.7         9.0           7.0         4.8  setosa 7 7 8 8 0
3∞ 8           0.3         9.4           7.0         4.8  setosa 7 7 8 8 0
3∞ 9           0.1         9.8           7.9         4.8  setosa 7 7 8 8 0
3∞ 0           0.0         9.7           7.0         4.8  setosa 7 7 8 8 0
3∞ 0           0.4         9.0           7.0         4.8  setosa 7 7 8 8 0
3∞ 0           0.0         9.3           7.1         4.0  setosa 7 7 8 8 0

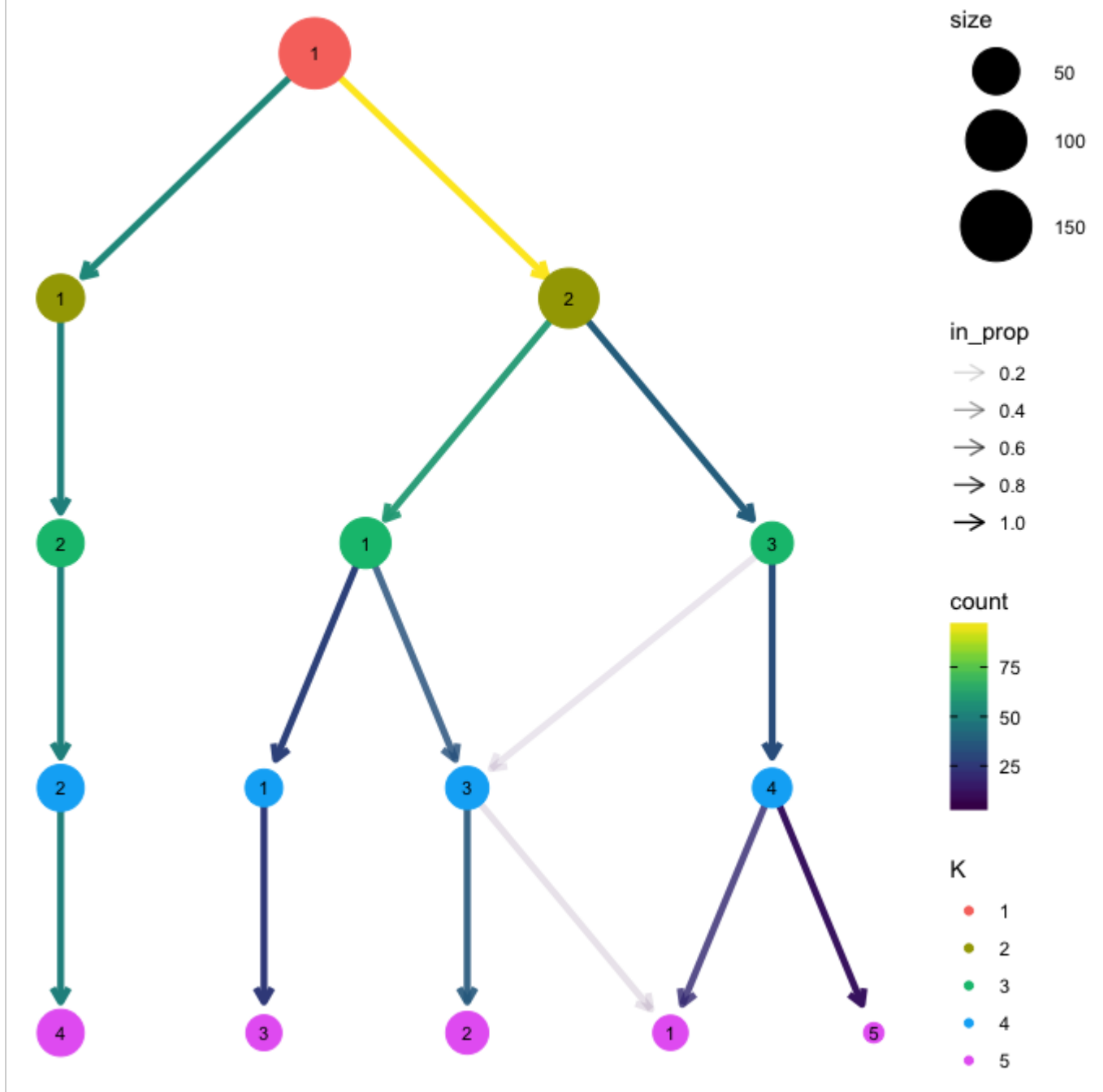
3∞           PC7           PC8
3∞ 7 -8.020780 -4.9739318
3∞ 8 -8.170708  4.7114478
3∞ 9 -8.222337  4.7003030
3∞ 0 -8.100909  4.9728334
3∞ 0 -8.182171 -4.9801000
3∞ 0 -8.824204 -4.1079940
```

Here we have a `data.frame` with the normal `iris` datasets, the measurements and species, plus some additional columns. These columns contain the cluster assignments from clustering this data using  $k$ -means with values ok  $k$  from  $k = 1$  to  $k = 5$ .

## 2.2 Plotting a tree

This clustering information is all we need to build a clustering tree. Each column must consist of numeric values indicating which cluster each sample has been assigned to. To plot the tree we just pass this information to the `clustree` function. We also need to specify a `prefix` string to indicate which columns contain the clusterings.

```
clustree(iris_clusts, prefix = "K")
```



We can see that one cluster is very distinct and does not change with the value of  $k$ . This is the **Iris setosa** samples which are very different to the other species. On the other side of the tree we see a single cluster that splits into the two clusters we would expect to see. After this the tree becomes messier and there are node with multiple incoming edges. This is a good indication that we have over clustered the data.

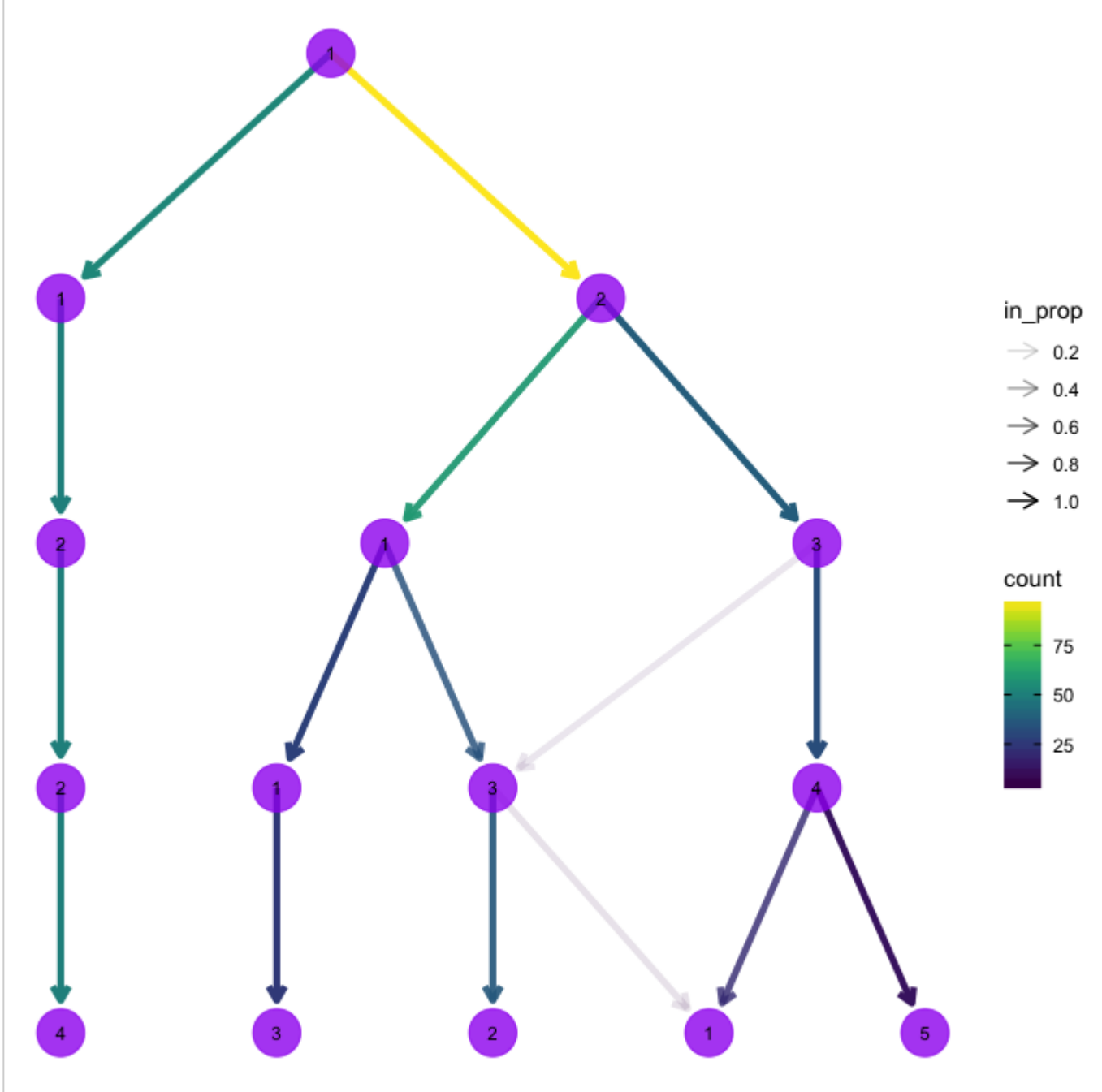
## 2.3 Controlling aesthetics

By default the size of each node is related to the number of samples in each cluster and the colour indicates the clustering resolution. Edges are coloured according to the number of samples they represent and the transparency shows the incoming node proportion, the number of samples in the edge divided by the number of samples in the node it points to. We can control these aesthetics by setting them to specific values:

```

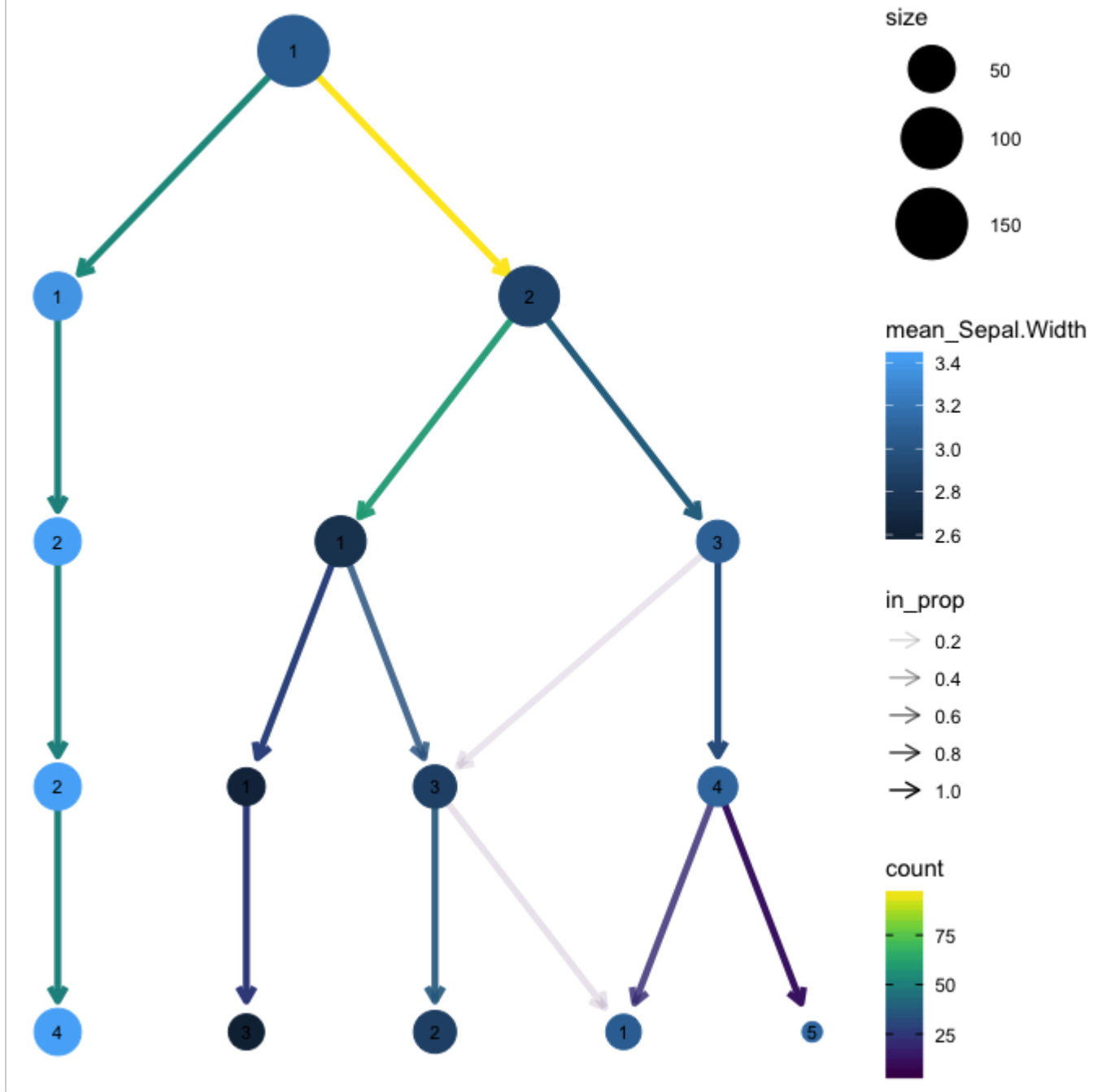
clustree(iris_clusts, prefix = "K", node_colour = "purple", node_size = 10,
         node_alpha = 0.8)

```



We can also link these aesthetics to other information we have about the samples. All the additional columns in the dataset are available to be added as attributes to the nodes in our tree. Because each node represents multiple samples we need to supply an aggregation function to use as well specifying a column name. Let's try colouring the nodes according to the sepal width:

```
clustree(iris_clusts, prefix = "K", node_colour = "Sepal.Width",
         node_colour_aggr = "mean")
```

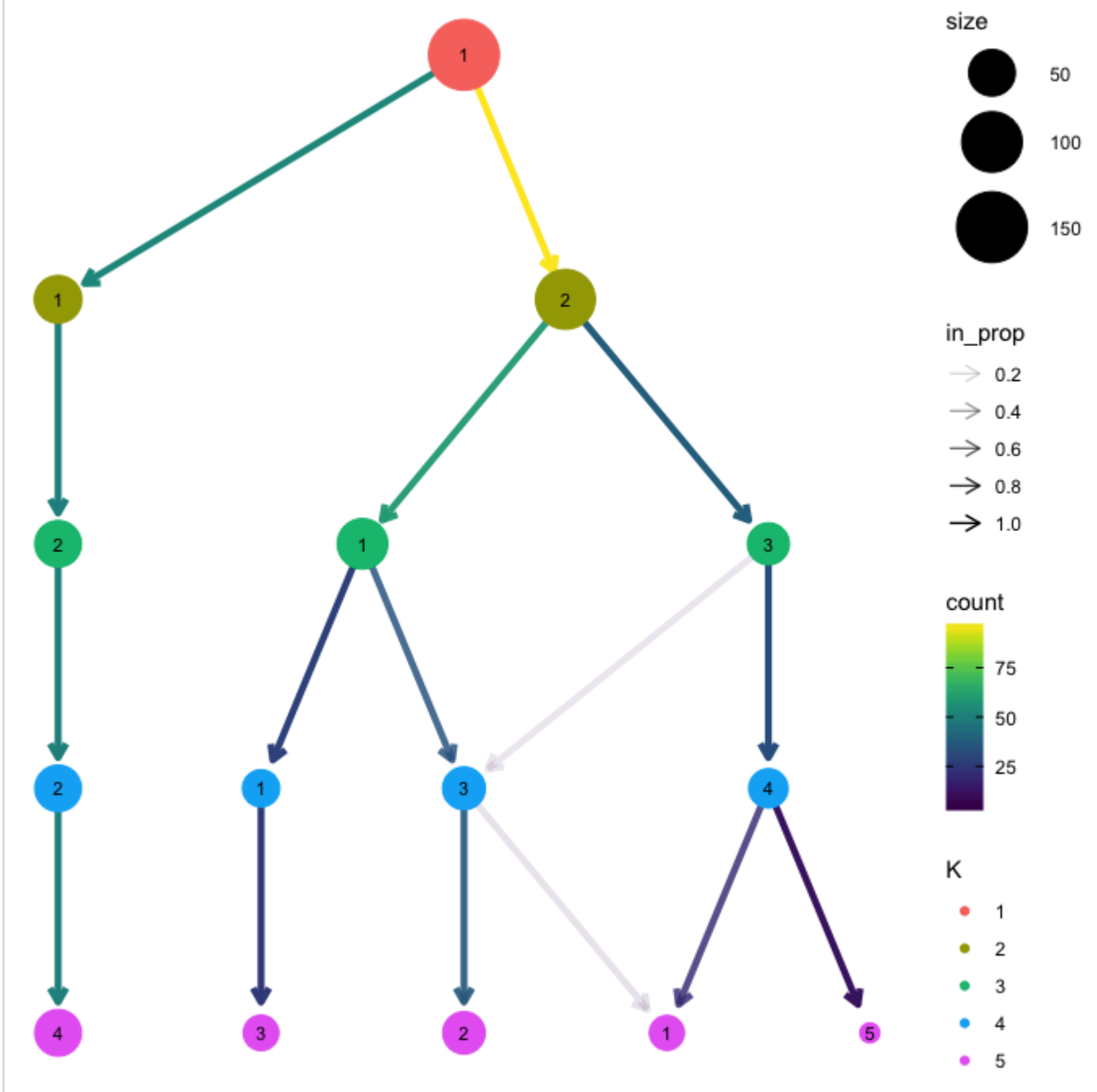


We can clearly see that the distinct cluster containing the **Iris setosa** samples has a wider sepal on average compared to the other clusters.

## 2.4 Layout

By default the tree is drawn using the Reingold-Tilford tree layout algorithm which tries to place nodes below their parents (Reingold and Tilford 1981). Alternatively we could use the Sugiyama layout by specifying the `layout` argument. This algorithm tries to minimise the number of crossing edges (Sugiyama, Tagawa, and Toda 1981) and can produce more attractive trees in some cases.

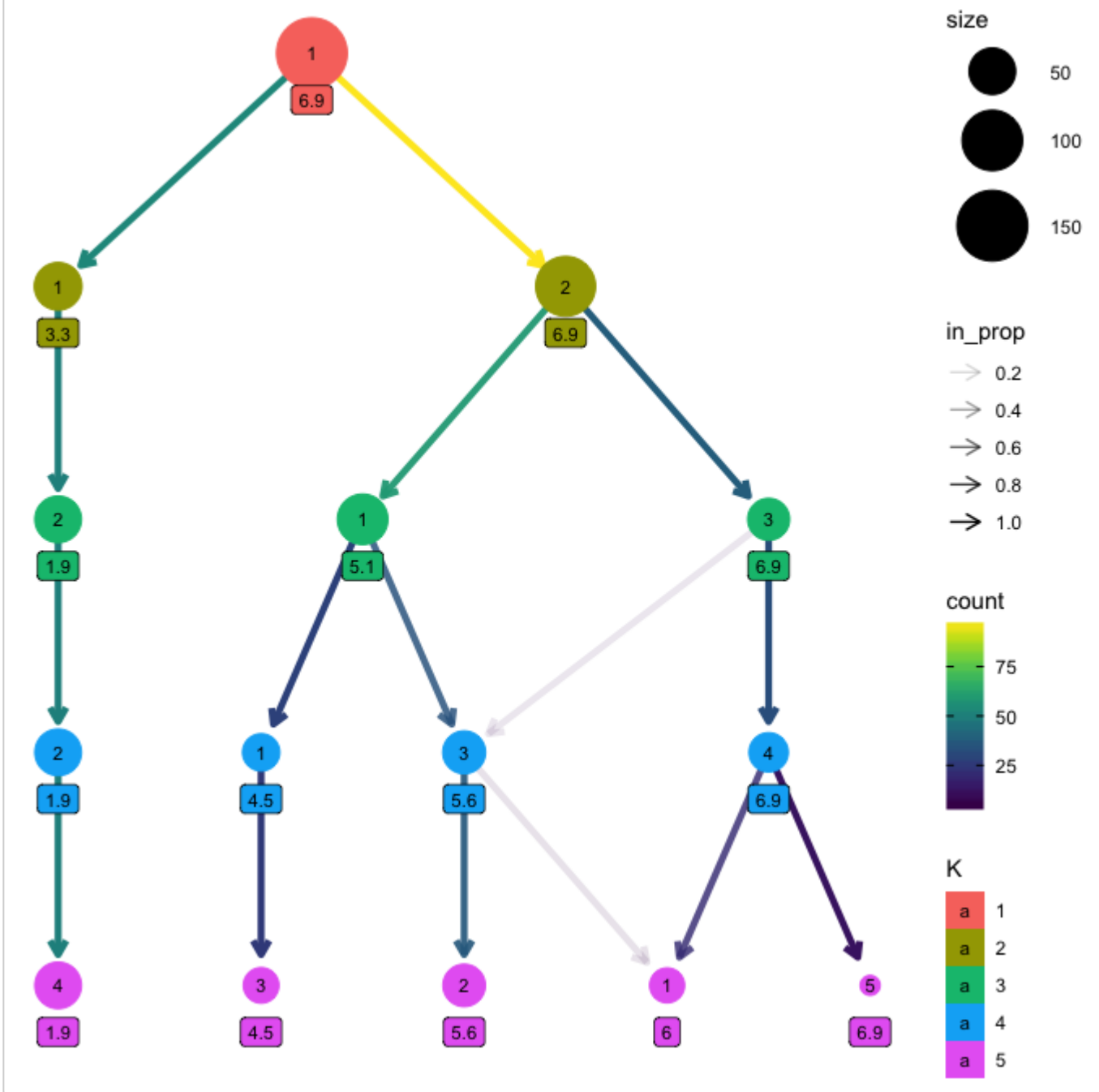
```
clustree(iris_clusts, prefix = "K", layout = "sugiyama")
```



For both of these layout algorithms clustree uses slightly modified versions by default. Only the core network of edges, those that are the highest in-proportion edge for a node, are used when creating the layout. In most cases this leads to more attractive trees that are easier to interpret. To turn this off, and use all edges for deciding the layout, we can set `use_core_edges` to `FALSE`.

```
clustree(iris_clusts, prefix = "K", layout = "sugiyama", use_core_edges = FALSE)
```





One way this can be useful is if we have assigned labels to the samples. Here is a custom function that labels a cluster if all the samples are the same species, otherwise it labels the cluster as “mixed”:

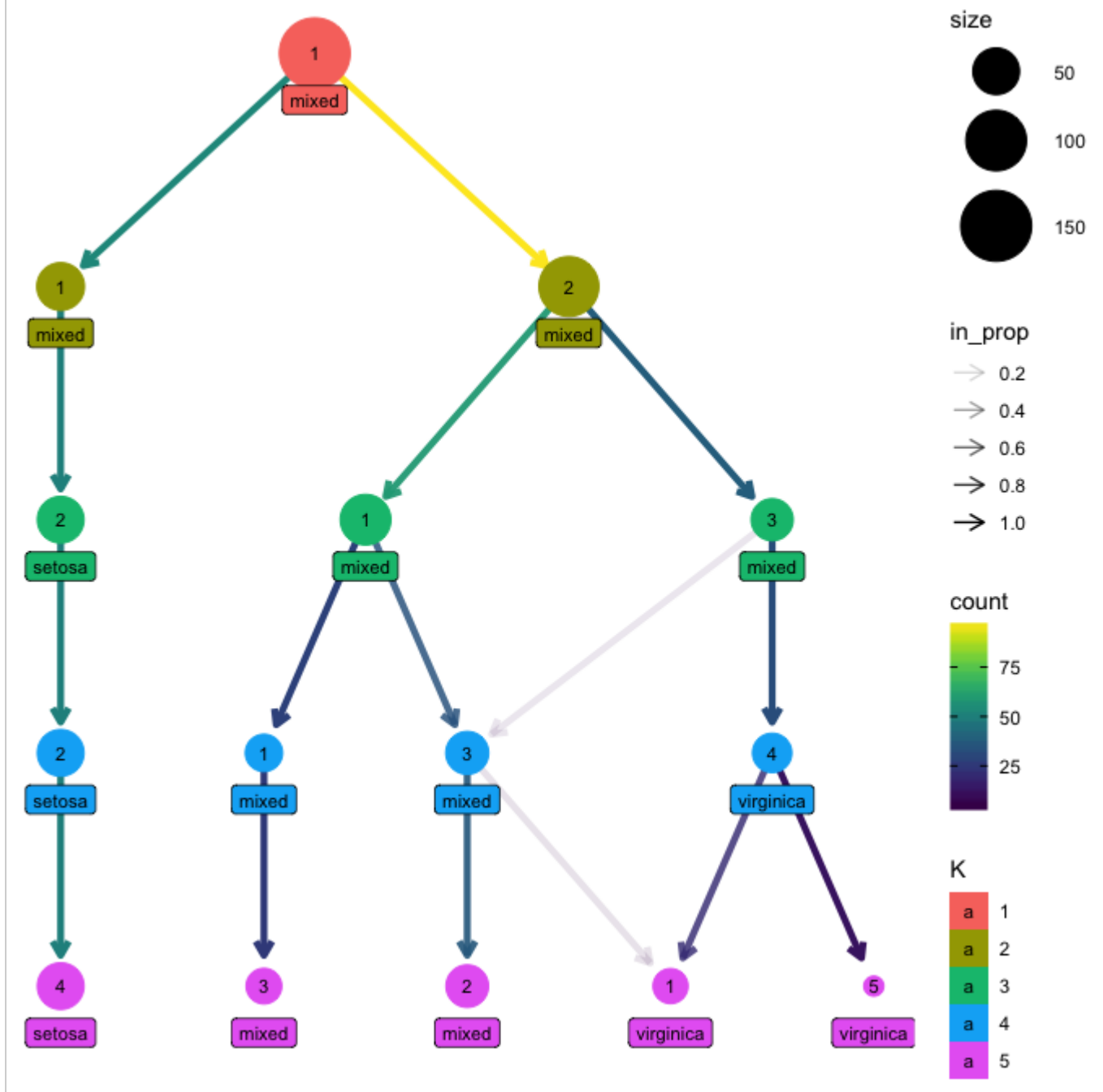
```

label_species <- function(labels) {
  if (length(unique(labels)) == 1) {
    species <- as.character(unique(labels))
  } else {
    species <- "mixed"
  }
  return(species)
}

clustree(iris_clusts, prefix = "K", node_label = "Species",
         node_label_aggr = "label_species")

```





# 3 Clustering trees for scRNA-seq data

Clustering has become a core tool for analysing single-cell RNA-sequencing (scRNA-seq) datasets. These datasets contain gene expression measurements from hundreds to hundreds of thousands of cells. Often samples come from complex tissues containing many types of cells and clustering is used to group similar cells together. To make it easier to produce clustering trees for these kinds of datasets we provide interfaces for some of the objects commonly used to analyse scRNA-seq data.

The clustree package contains an example simulated scRNA-seq data that has been clustered using the SC3 and Seurat packages.

```
data("sc_example")
names(sc_example)
3∞ [7] TMcountsTM TMlogcountsTM TMtsneTM TMsc9@clustersTM
3∞ [0] TMseurat@clustersTM
```

### 3.1 SingleCellExperiment objects

The `SingleCellExperiment` is one of these common objects, used across a range of Bioconductor packages. Let's have a look at an example, but first we need to convert the example dataset to a `SingleCellExperiment` object:

```
suppressPackageStartupMessages(library("SingleCellExperiment"))

sce <- SingleCellExperiment(assays = list(counts = sc_example$counts,
                                         logcounts = sc_example$logcounts),
                           colData = sc_example$sc3_clusters,
                           reducedDims = SimpleList(TSNE = sc_example$tsne))
```

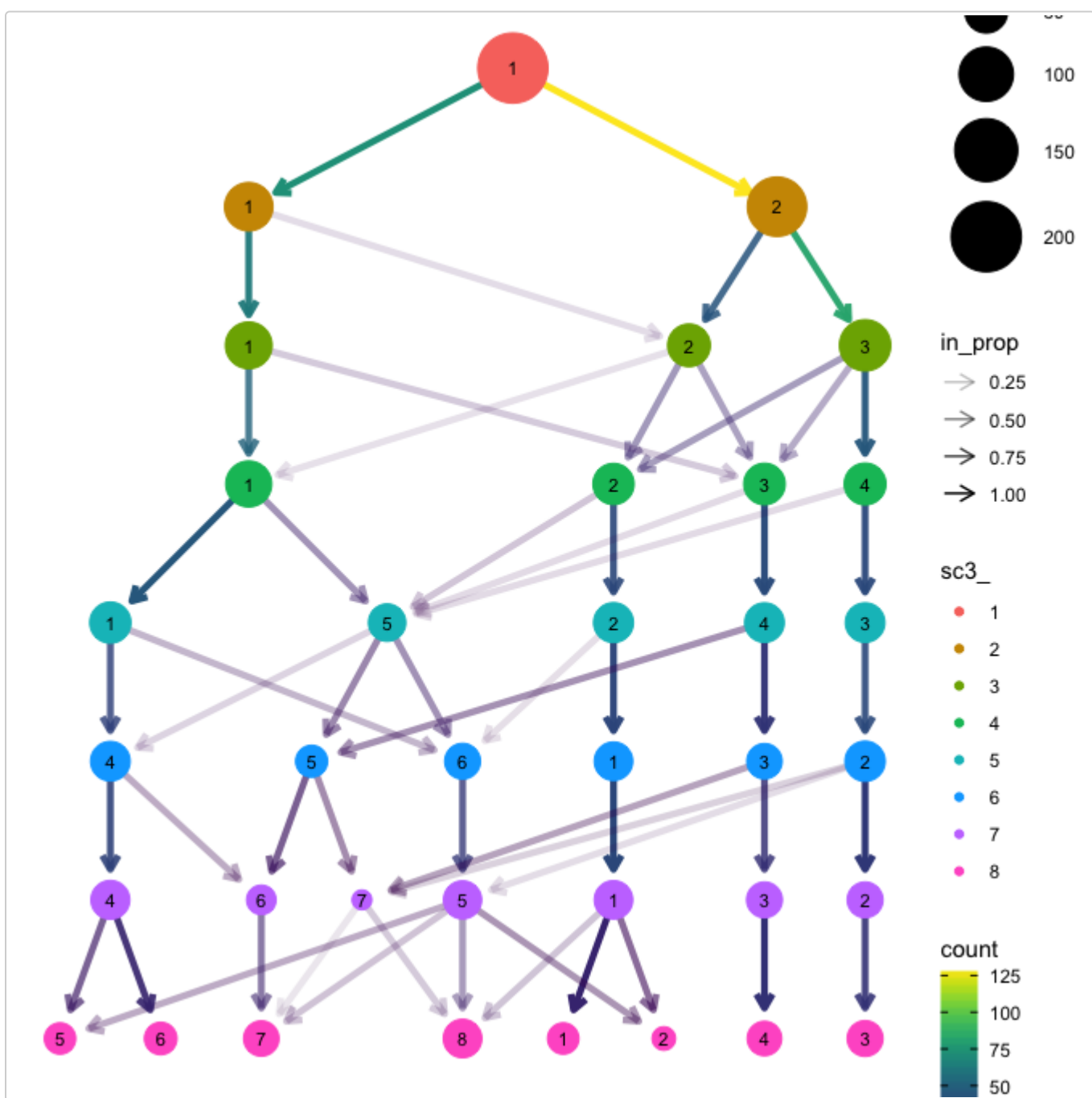
The clustering information is held in the `colData` slot.

```
head(colData(sce))
```

	sc9070clusters	sc9080clusters	sc9090clusters	sc9000clusters
	lfactor	lfactor	lfactor	lfactor
Cell17	7	8	9	0
Cell18	7	8	9	8
Cell19	7	7	7	7
Cell10	7	7	7	9
Cell10	7	8	9	0
Cell10	7	8	8	8
	sc9000clusters	sc9000clusters	sc9010clusters	sc9020clusters
	lfactor	lfactor	lfactor	lfactor
Cell17	9	8	8	9
Cell18	8	7	7	7
Cell19	0	0	0	0
Cell10	0	9	9	0
Cell10	9	8	0	2
Cell10	8	7	7	7

We can plot a clustering tree in the same way we did with a `data.frame`. In this case the clustering column names contain a suffix that needs to be stripped away, so we will pass that along as well.

```
clustree(sce, prefix = "sc3_", suffix = "_clusters")
```



### 3.2 Seurat objects

Clustering trees can also be produced directly from `Seurat` objects. Let's convert our `SingleCellExperiment` to `Seurat` format:

```
suppressPackageStartupMessages(library("Seurat"))

seurat <- CreateSeuratObject(sc_example$counts,
                             meta.data = sc_example$seurat_clusters)
seurat <- SetDimReduction(seurat, "TSNE", "cell.embeddings", sc_example$tsne)
```

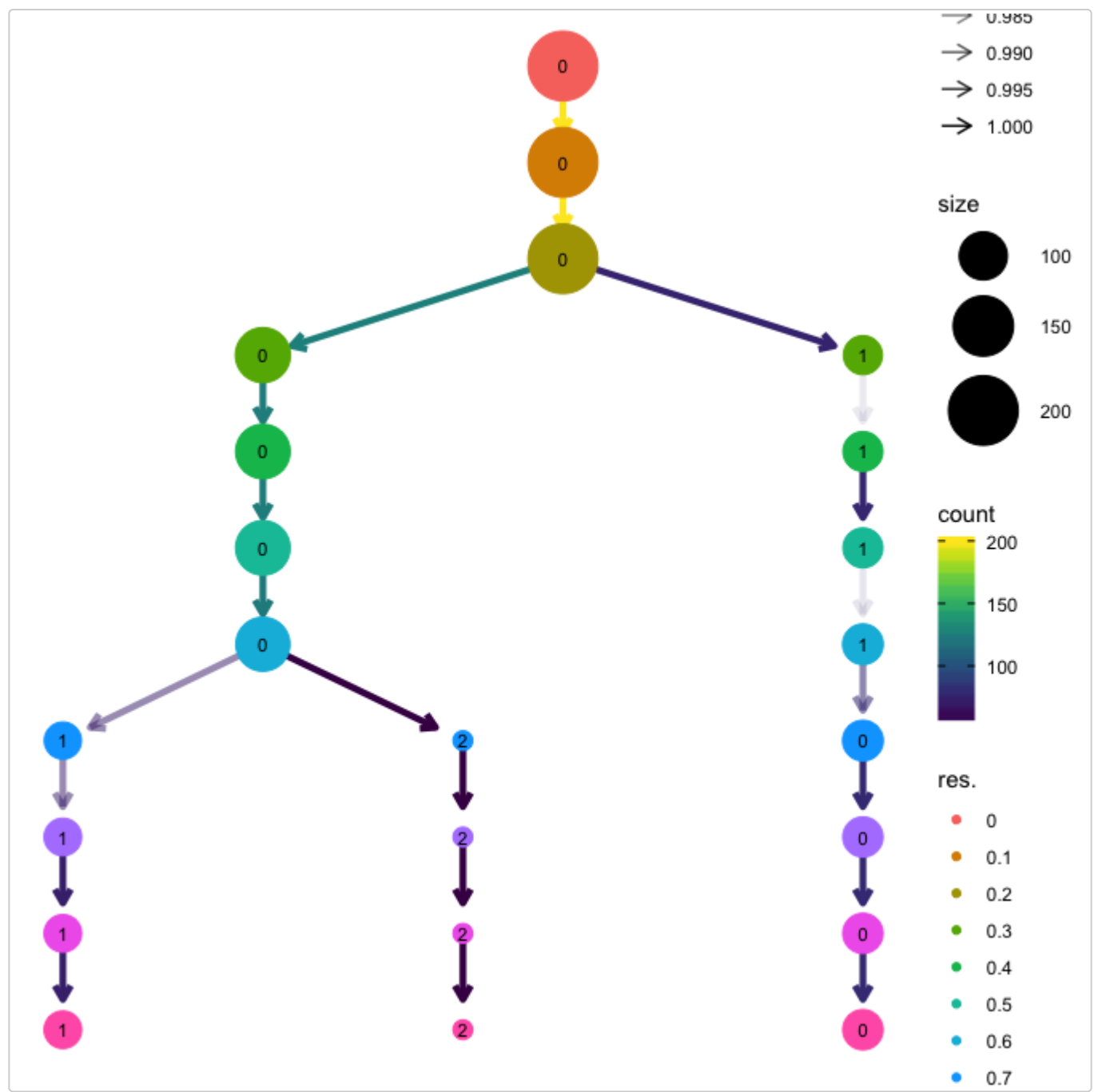
In this case the clustering information is held in the `meta.data` slot:

```
head(seurat@meta.data)
3∞      nGene nUMI res.4 res.4.7 res.4.8 res.4.9 res.4.0 res.4.0 res.4.0
3∞ Cell17  049 0002    4      4      4      7      7      7      7
```

300	Cell18	003	0188	4	4	4	7	7	7	7
300	Cell19	040	9304	4	4	4	4	4	4	4
300	Cell10	080	0288	4	4	4	4	4	4	4
300	Cell10	038	0402	4	4	4	7	7	7	7
300	Cell10	074	0124	4	4	4	4	4	4	4
300		res.4.1	res.4.2	res.4.3	res.7	orig.ident				
300	Cell17	4	4	4	4	SeuratProject				
300	Cell18	4	4	4	4	SeuratProject				
300	Cell19	7	7	7	7	SeuratProject				
300	Cell10	7	7	7	7	SeuratProject				
300	Cell10	4	4	4	4	SeuratProject				
300	Cell10	7	7	7	7	SeuratProject				

Because this object is only used by the `Seurat` package we can assume the prefix of the clustering columns.

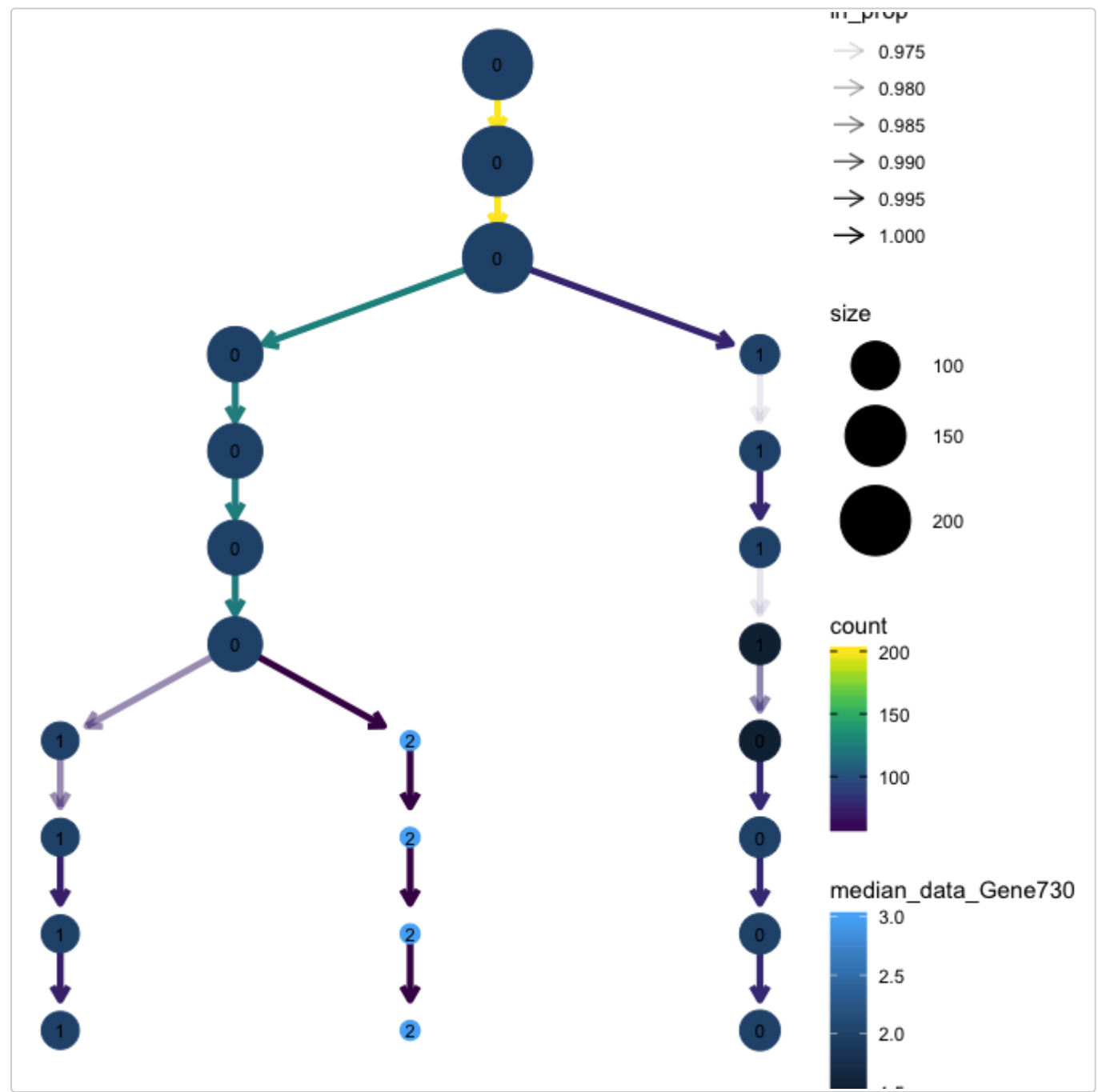
```
clustree(seurat)
```



### 3.3 Using genes as aesthetics

As well as being able to use any additional columns for aesthetics we can also use the expression of individual genes. Let's colour the nodes in the `Seurat` tree by `Gene730` (a highly variable gene). Again we need to supply an aggregation function.

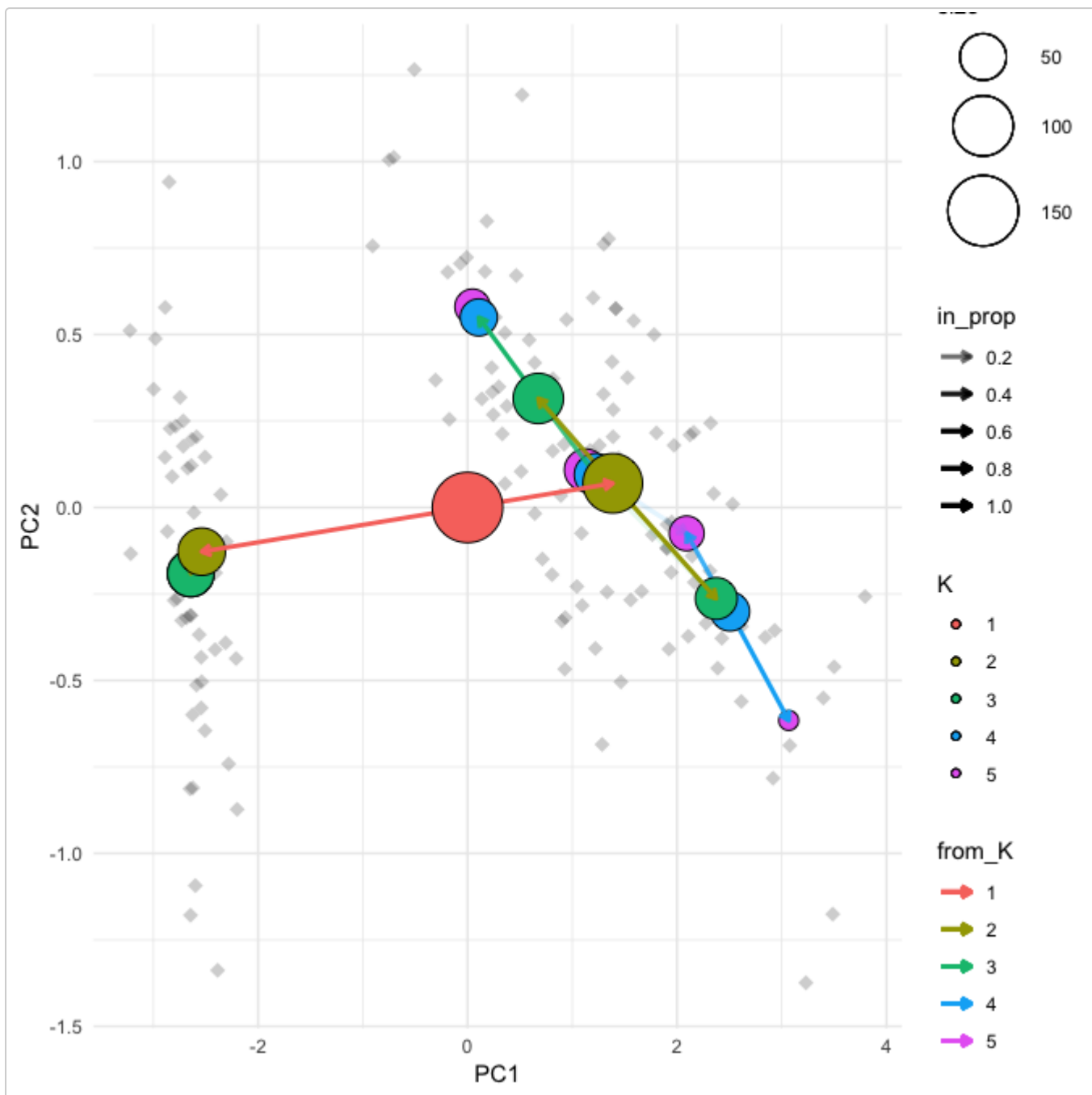
```
clustree(seurat, node_colour = "Gene730", node_colour_aggr = "median")
```



## 4 Overlaying clustering trees

One way to think of clustering trees is that they add an extra dimension to the data that shows how clusters “evolve” over time (increasing resolution). Some times it can be useful to overlay this dimension on other dimensions in the data, particularly those that come from dimensionality reduction techniques. We can do this using the `clustree_overlay` function:

```
clustree_overlay(iris_clusts, prefix = "K", x_value = "PC1", y_value = "PC2")
```

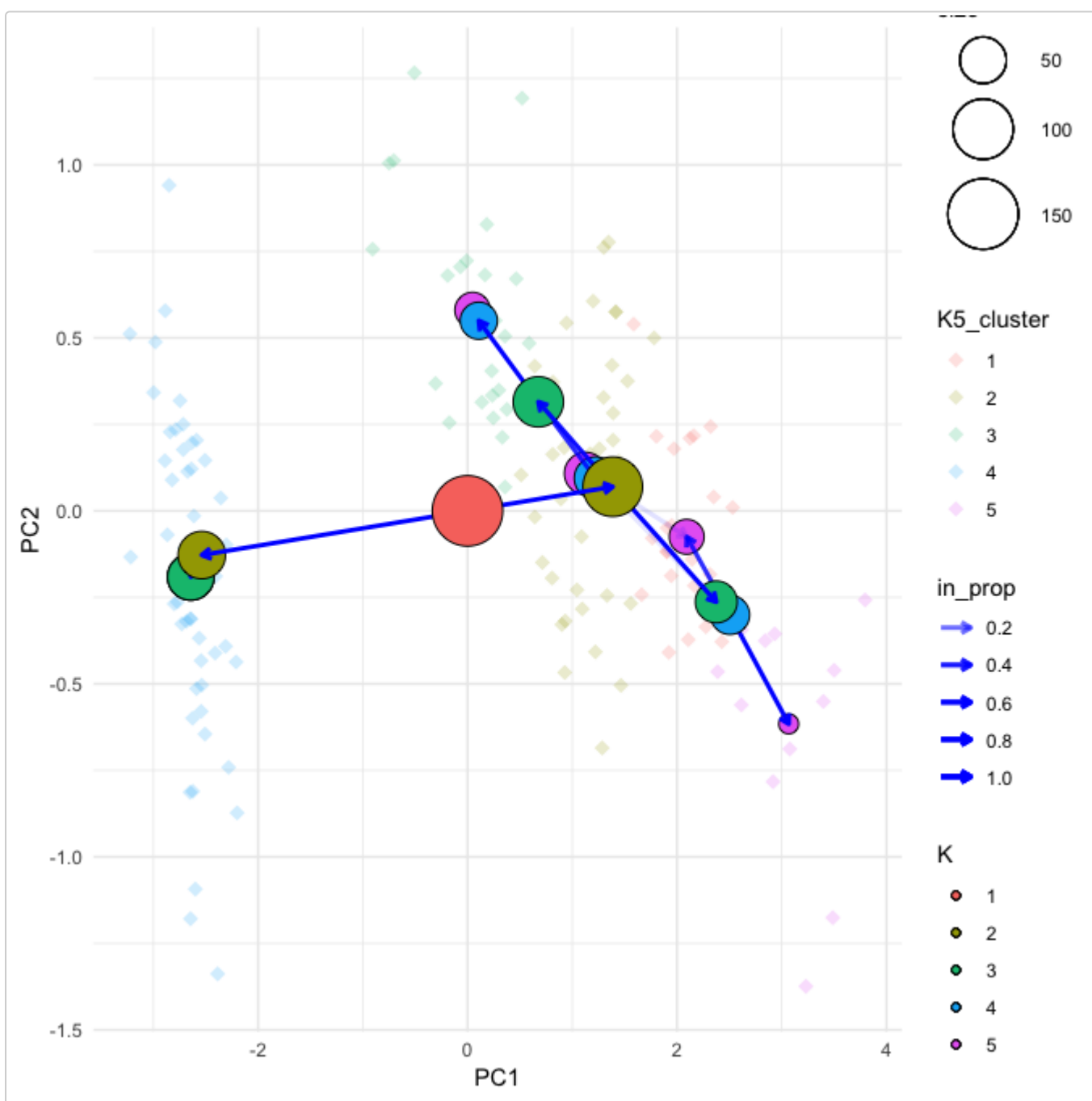


The easiest way to understand this plot is to imagine that you are looking down on the clustering tree from above. The x and y axes are the chosen dimensions in the data and the z axis is the clustering resolution. Each cluster node is placed at the mean x and y values of the samples it contains. We can also see points corresponding to the individual samples.

## 4.1 Choosing what to colour

Due to the way `ggplot2` works we can only colour one element in the plot (the node points actually use the fill aesthetic). By default the tree edges are coloured according to the clustering resolution they originate from. Alternatively we can choose to colour the sample points by their highest resolution cluster. The colour of whichever element isn't using the colour aesthetic can be set using the `alt_colour` argument.

```
clustree_overlay(iris_clusts, prefix = "K", x_value = "PC1", y_value = "PC2",
  use_colour = "points", alt_colour = "blue")
```



Aesthetics for the clustering tree nodes can be controlled in the same way as for the regular `clustree` function and there are additional arguments for adjusting the appearance of edges and points.

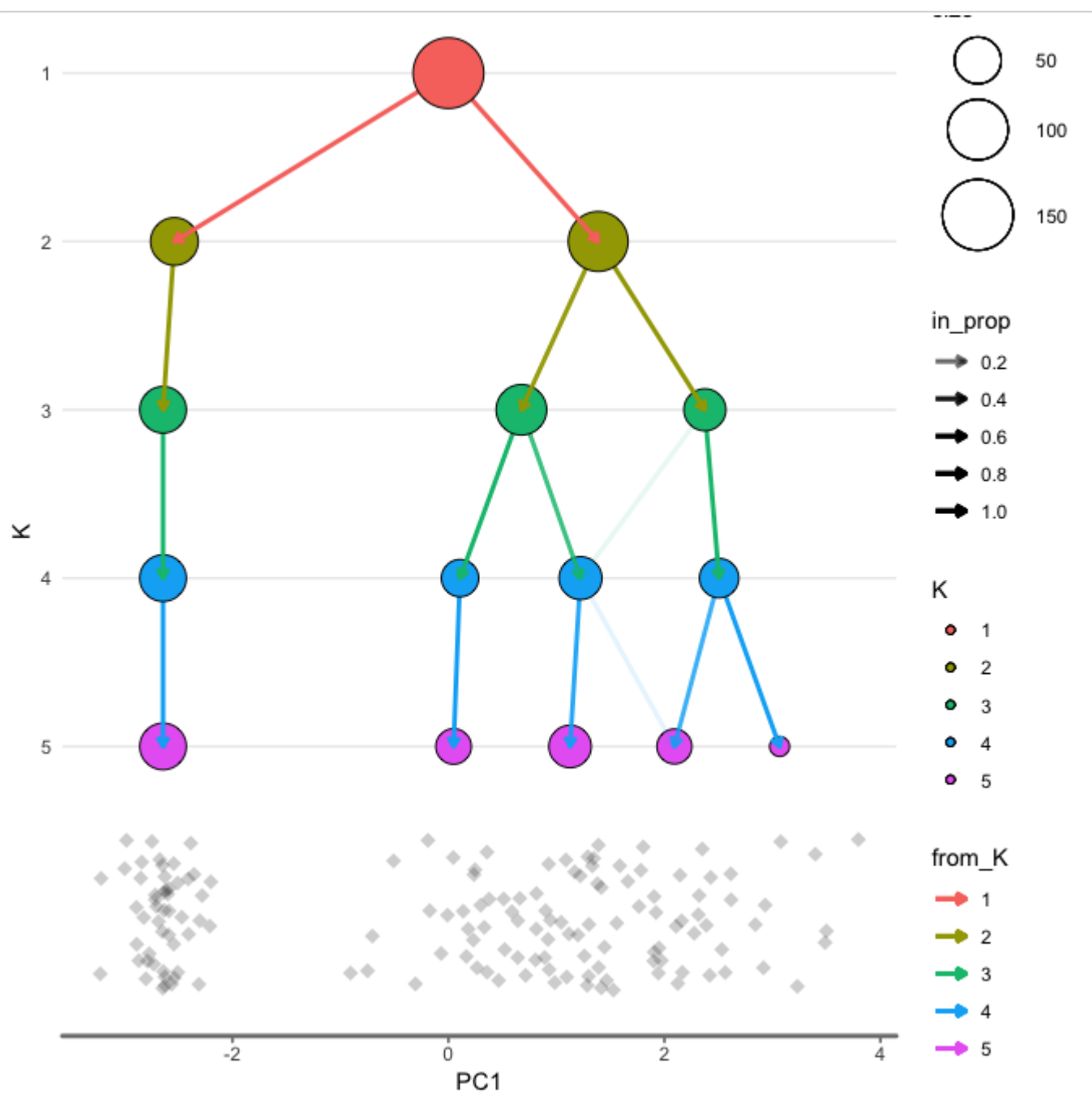
## 4.2 Labelling nodes

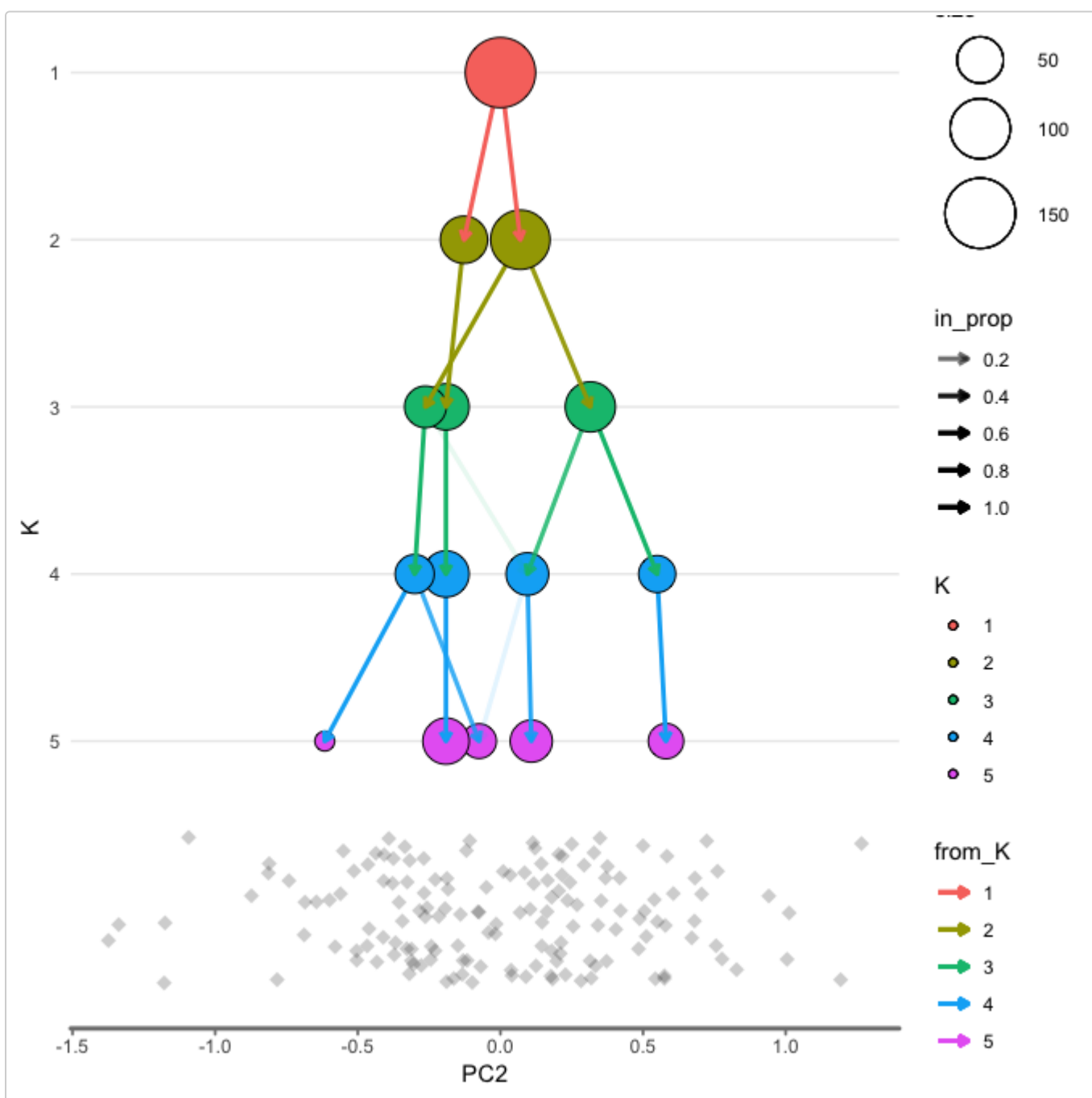
One of the downsides of this visualisation is that it can be hard to identify the clusters that each node represents as they can often overlap. To make this a bit easier we can label each node with the resolution and cluster ID.

```
clustree_overlay(iris_clusts, prefix = "K", x_value = "PC1", y_value = "PC2",
  label_nodes = TRUE)
```







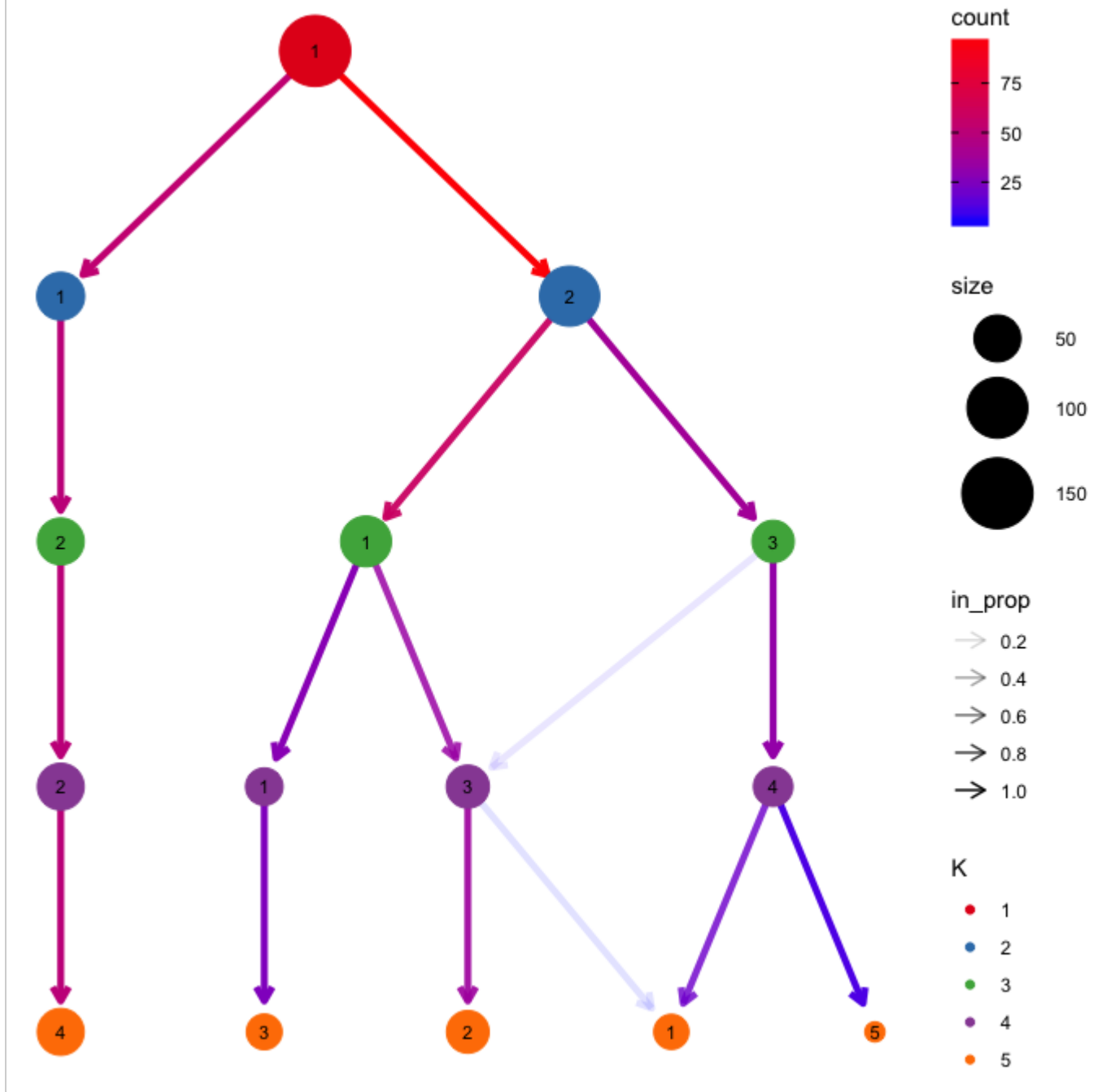


# 5 Modifying appearance

The `clustree` function returns a `ggplot` object which can be modified using functions in the `ggplot2` or `ggraph` packages. For example we could change the colour scales used for the nodes and edges:

```
clustree(iris_clusts, prefix = "K") +
  scale_color_brewer(palette = "Set1") +
  scale_edge_color_continuous(low = "blue", high = "red")
```

3> Scale for 'edge@colour' is already present. Adding another scale for  
 3> 'edge@colour', which will replace the existing scale.



## 5.1 Legends

The way `ggplot` objects is displayed is relative to the size of the plotting window or output file. While the main plot will always fit sometimes legends will be placed outside the visible area. One solution to this is to simply increase the size of the image. An alternative solution is to turn off some of the legends, either by setting some of the aesthetics to static values or by using the `guides` function. We could also move them to somewhere they might fit better. For example let's remove the edge legends and move the rest to the bottom:

```

clustree(iris_clusts, prefix = "K") +
  guides(edge_colour = FALSE, edge_alpha = FALSE) +
  theme(legend.position = "bottom")

```



```
3∞    journal ≥ {Gigascience},  
3∞    volume ≥ {1},  
3∞    number ≥ {1},  
3∞    month ≥ {jul},  
3∞    year ≥ {8472},  
3∞    url ≥ {https://doi.org/74.7439/gigascience/giy429},  
3∞    doi ≥ {74.7439/gigascience/giy429},  
3∞ }
```

# References

Reingold, E M, and J S Tilford. 1981. "Tidier Drawings of Trees." **IEEE Trans. Software Eng.** SE-7 (2):223–28. <https://doi.org/10.1109/TSE.1981.234519>.

Sugiyama, K, S Tagawa, and M Toda. 1981. "Methods for Visual Understanding of Hierarchical System Structures." **IEEE Trans. Syst. Man Cybern.** 11 (2):109–25. <https://doi.org/10.1109/TSMC.1981.4308636>.

Zappia, Luke, and Alicia Oshlack. 2018. "Clustering trees: a visualization for evaluating clusterings at multiple resolutions." **Gigascience** 7 (7). <https://doi.org/10.1093/gigascience/giy083>.