

Seurat - Clustering Tutorial

Dennis Dimitri Krutkin

2023-06-16

```
#Load necessary libraries
library(dplyr)

##

## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##      filter, lag

## The following objects are masked from 'package:base':
##
##      intersect, setdiff, setequal, union

library(Seurat)

## Attaching SeuratObject

library(patchwork)
```

Import 10X dataset

```
# Load the PBMC data set:
pbmc.data <- Read10X(data.dir = "./filtered_gene_bc_matrices/hg19/")

# Initialize the Seurat object with the raw (non-normalized data):
pbmc <- CreateSeuratObject(counts = pbmc.data, project = "pbmc3k", min.cells = 3, min.features = 200)

## Warning: Feature names cannot have underscores ('_'), replacing with dashes
## ('-')

pbmc

## An object of class Seurat
## 13714 features across 2700 samples within 1 assay
## Active assay: RNA (13714 features, 0 variable features)
```

Examine the loaded dataset

```
# Examine a 3 genes in the first 25 cells:
pbmc.data[c("CD3D", "TCL1A", "MS4A1"), 1:25]

## 3 x 25 sparse Matrix of class "dgCMatrx"

##      [[ suppressing 25 column names 'AAACATACACCAC-1', 'AAACATTGAGCTAC-1', 'AAACATTGATCAGC-1' ...
]]

##
## CD3D   4 . 10 . . . 1 2 3 1 . . 2 7 1 . . 1 3 . 2   3 . . . .
## TCL1A . . . . . . . . 1 . . . . . . . . . . . 1 . . .
## MS4A1 . 6 . . . . . . 1 1 1 . . . . . . . . . 36 1 2 . .
```

The “.” values in the matrix represent 0s (no molecules detected).

Most values in a scRNA-seq matrix are 0 - Seurat uses a sparse-matrix representation whenever possible. This results in significant memory and speed savings for Drop-seq/inDrop/10x data.

```
#Check the full size of object:
dense.size <- object.size(as.matrix(pbmc.data))
dense.size

## 709591472 bytes

#Check the sparse-represented size:
sparse.size <- object.size(pbmc.data)
sparse.size

## 29905192 bytes

#Check the ratio of full-size:sparse-size:
dense.size/sparse.size

## 23.7 bytes
```

Standard pre-processing workflow

The steps below demonstrate the standard pre-processing workflow for scRNA-seq data using Seurat.

The steps include selection and filtration of cells based on QC metrics, data normalization and scaling, and the detection of highly variable features.

Quality Control and selecting cells to further analyze

In Seurat, QC metrics can be explored and cells can be filtered according to any user-defined criteria.

Some QC metrics used by the community include:

- The number of unique genes detected in each cell
 - Low-quality cells or empty droplets often have very few genes
 - Cell doublets or multiplets may have unusually high gene counts
- The total number of molecules detected within a cell, which correlate strongly with unique genes
- The percentage of reads which map to the mitochondrial genome
 - Poor quality and dying cells frequently show extensive mitochondrial contamination
 - The PercentageFeatureSet() function calculates the percentage of counts originating from a set of features to calculate mitochondrial QC metrics
 - Genes starting with MT- are pooled as a set of mitochondrial genes

```
#Add QC stats for mitochondrial genes to the object metadata:
pbmc[["percent.mt"]] = PercentageFeatureSet(pbmc, pattern = "^MT-")
#pbmc[["percent.mt"]]
```

Where are QC metrics stored in Seurat?

- The number of unique genes and total molecules are automatically calculated during CreateSeuratObject()
 - They are stored in the object metadata

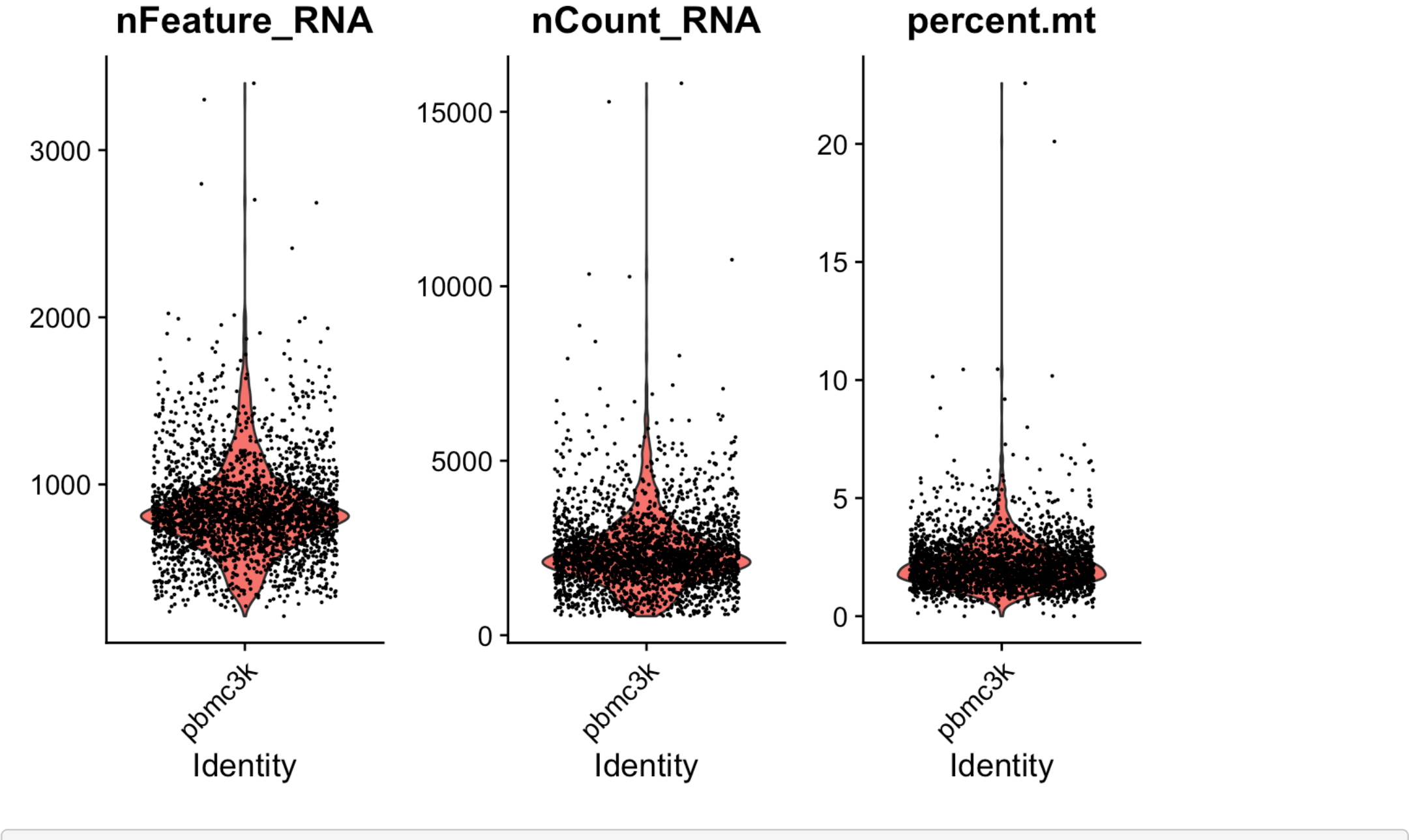
```
#Show QC metrics for the first 5 cells:
head(pbmc@meta.data, 5)

##      orig.ident nCount_RNA nFeature_RNA percent.mt
## AAACATACACCAC-1 pbmc3k      2419         779   3.0177759
## AAACATTGAGCTAC-1 pbmc3k      4903        1352   3.7935958
## AAACATTGATCAGC-1 pbmc3k      3147        1129   0.8897363
## AAACCGTGCTTCCG-1 pbmc3k      2639         960   1.7430845
## AAACCGTGTATGCG-1 pbmc3k       980         521   1.2244898
```

Below, QC metrics are visualized and then used to filter cells

- Cells are filtered that have unique feature counts over 2,500 or less than 200
- Cells are filtered that have >5% mitochondrial counts

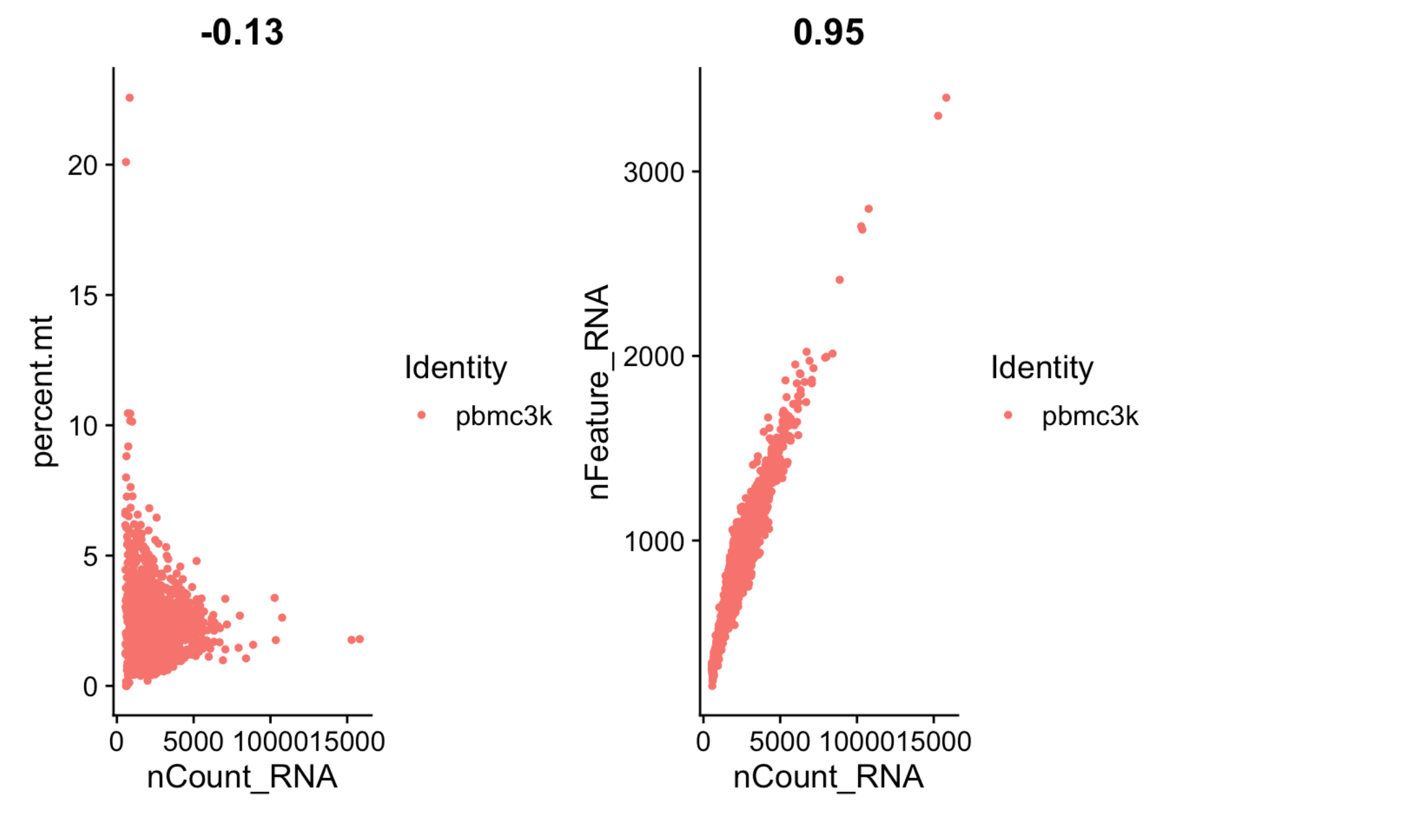
```
#Visualize QC metrics as violin plots:
VlnPlot(pbmc, features = c("nFeature_RNA", "nCount_RNA", "percent.mt"), ncol = 3)
```



```
#FeatureScatter is typically used to visualize feature-feature relationships

#It can be used for anything calculated by the object, such as columns in object metadata, PC scores, etc.

plot1 = FeatureScatter(pbmc, feature1 = "nCount_RNA", feature2 = "percent.mt")
plot2 = FeatureScatter(pbmc, feature1 = "nCount_RNA", feature2 = "nFeature_RNA")
plot1 + plot2
```



```
#Filter the PBMC data set:
pbmc_filtered <- subset(pbmc, subset = nFeature_RNA > 200 & nFeature_RNA < 2500 & percent.mt < 5)
```