# R Notebook

Rahi Shah (rs8579@nyu.edu (mailto:rs8579@nyu.edu))

downloding and loading packages:

Hide

```
install.packages('patchwork')
```

```
Error in install.packages : Updating loaded packages
```

Step 1: loading the TSV file

| | cell_091.133 <int> | cell_177.113 <int> | cell_289.088 <int> | cell_205.268 <int> | cell_162.063 <int> | cell_183.039 <int> |
|---|---|---|---|---|---|---|
| A1BG | 0 | 0 | 0 | 0 | 0 | 0 |
| A1CF | 0 | 0 | 0 | 1 | 0 | 0 |
| A2M | 0 | 0 | 0 | 0 | 0 | 0 |
| A2ML1 | 0 | 0 | 0 | 0 | 0 | 0 |
| A3GALT2 | 0 | 0 | 0 | 0 | 0 | 0 |
| A4GALT | 0 | 0 | 0 | 0 | 0 | 0 |

6 rows | 1-7 of 2000 columns

Step 2: creating Seurat object

project is titled 'pdac1'. the Seurat object created here retains genes whcih are expressed in at least 3 cells, and it retains cells that express at least 200 features (genes)

Hide

```
pdac1 <- CreateSeuratObject(counts = data, project = "pdac1", min.cells = 3, min.features = 200)
pdac1
```

```
An object of class Seurat
13248 features across 633 samples within 1 assay
Active assay: RNA (13248 features, 0 variable features)
```

step 3: label mitochondrial genes

This code creates a column named "percent.mt" and isolates all the mitochondrial expression
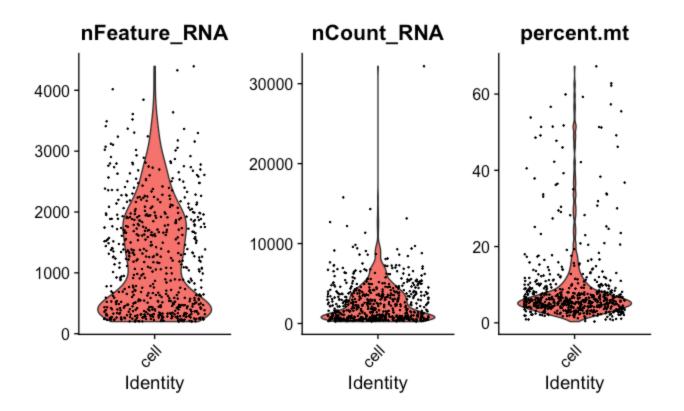
Hide

```
pdac1[["percent.mt"]] <- PercentageFeatureSet(object = pdac1, pattern = "^MT-")
```

step 4: visualize the distribution:

the code below creates a violin plot, which visualizes the feature data, count data, and percent of mitochondrial genes

Hide

```
VlnPlot(object = pdac1, features = c("nFeature_RNA", "nCount_RNA", "percent.mt"), ncol =
3)
```



step 5: filter data

the code below uses the 'subset' function to select data which has features between 200 and 2500, while the percent mitochondrial expression is less than 5.

Hide

```
pdac1 <- subset(x = pdac1, subset = nFeature_RNA > 200 & nFeature_RNA < 2500 & percent.m
t < 5)
```

step 6: normalize data

normalzing the data using Log method, and scaling the data by a vector of 10000

Hide

```
pdac1 <- NormalizeData(object = pdac1, normalization.method = "LogNormalize", scale.fact
or = 10000)
```

```
Performing log-normalization
0%    10   20   30   40   50   60   70   80   90   100%
[----|----|----|----|----|----|----|----|----|----|
**************************************************|
```

step 7: calculate gene variation

the "FindVariableFeatures" function allows R to look through the seurat obect and use the vst method to filter out the top 2000 features/gene which have the most amount of variation across all the cells in the data

Hide

```
pdac1 <- FindVariableFeatures(object = pdac1, selection.method = "vst", nfeatures = 200
0)
```

```
Calculating gene variances
0%    10   20   30   40   50   60   70   80   90   100%
[----|----|----|----|----|----|----|----|----|----|
**************************************************|
Calculating feature variances of standardized and clipped values
0%    10   20   30   40   50   60   70   80   90   100%
[----|----|----|----|----|----|----|----|----|----|
**************************************************|
```

step 8: scale data

scaling all the count data stored in the seurat object for all the genes (rows)

Hide

```
all.genes <- rownames(x = pdac1)
pdac1 <- ScaleData(object = pdac1, features = all.genes)
```

```
Centering and scaling data matrix

  |
  |                                                          |    0%
  |
  |=====                                                     |    7%
  |
  |=========                                                 |   14%
  |
  |==============                                            |   21%
  |
  |===================                                       |   29%
  |
  |=======================                                   |   36%
  |
  |============================                              |   43%
  |
  |=================================                         |   50%
  |
  |=====================================                     |   57%
  |
  |==========================================                |   64%
  |
  |===============================================           |   71%
  |
  |====================================================      |   79%
  |
  |=========================================================  |   86%
  |
  |=============================================================|   93%
  |
  |=================================================================| 100%
```

step 9: PCA

performing principal component analysis on the seurat object, specifically using the features that show high
variablility

Hide

```
pdac1 <- RunPCA(object = pdac1, features = VariableFeatures(object = pdac1))
```

```
PC_ 1
Positive:  SPINK1, TM4SF1, DUOXA2, AQP3, APOL1, SOX4, GC, TSPAN8, SLPI, MMP7
       DMBT1, DSTN, AGR2, SLC2A1, GABRP, ITGA2, EFNA1, AKR1B10, CYB5A, F3
       CES1, TFF3, AKR1C3, ANKRD36C, MPZL2, FXYD3, FHL2, BAIAP2L1, EIF4A2, LIPH
Negative:  LAPTM5, TYROBP, ALOX5AP, FPR1, S100A9, BASP1, RGS2, S100A8, SLC2A3, LCP1
       HCLS1, BCL2A1, FCGR3A, CSF3R, G0S2, FCGR3B, ARRB2, ST20, CD53, C5AR1
       ARHGDIB, FCGR2A, ITGB2, OSM, FCER1G, SAMSN1, TNFRSF1B, AIF1, SLA, CD37
PC_ 2
Positive:  G0S2, CXCL8, CXCR2, S100A8, ST20, BCL2A1, S100A9, HCAR3, FFAR2, IL1R2
       HCAR2, CMTM2, CSF3R, PROK2, FPR2, MMP9, SOD2, RIPOR2, FPR1, TLR4
       CXCR1, MME, PLEK, ALPL, FAM177B, RHOH, MGAM, MMP25, AQP9, ALOX5AP
Negative:  MS4A7, SLCO2B1, HLA-DQB1, HLA-DPA1, TGFBI, HLA-DQA1, HLA-DRA, CD163, HLA-DQA
2, HLA-DRB1
       HLA-DPB1, HLA-DRB5, LGALS1, C1QC, VIM, HLA-DQB2, NPC2, FCGR1A, C1QB, SDS
       APOE, GPR183, CTSB, C1QA, CSF1R, CD74, MS4A4A, CPVL, CD4, GPX1
PC_ 3
Positive:  AQP3, DUOXA2, GC, TFF3, CES1, DMBT1, HLA-DRB1, AKR1B10, HLA-DRB5, CRP
       HK2, CRISP3, NR4A1, SPINK1, CD74, CPM, ANKRD36, CTSS, FCGBP, REG1A
       ANKRD36C, SGK1, F5, CYB5A, TSPAN8, AKR1C3, ANKRD36B, AGR2, CA12, ALDOB
Negative:  COL6A1, SERPINB5, COL6A2, COL17A1, MIA, MGP, CRABP2, SNCG, WNT11, PLAU
       TIMP3, IGFBP2, S100P, DCBLD2, PCDH7, HMGA2, PODXL, LY6D, FSTL1, HIST1H2BD
       VCAM1, TOMM70, TRIM29, RNASE1, DMKN, LYPD3, KRT16, PRODH, ZSCAN9, MUC16
PC_ 4
Positive:  CD2, SEPT6, CCL5, SPOCK2, IL2RB, CD3D, CD8A, CD7, CD8B, KLRC1
       GZMB, RGS1, KLRC2, CD247, NPIPA7, IRF4, NPIPA3, NPIPA8, NPIPA2, SLAMF7
       PRF1, FYN, CXCR3, RPL3, ZNF683, SIT1, TRAF3, LCK, CTLA4, NPIPA1
Negative:  S100A8, S100A9, PLAUR, SOD2, BCL2A1, FPR1, BASP1, ST20, IL1RN, FCN1
       SLC11A1, FCGR3B, AQP9, S100P, C5AR1, FCGR3A, SLC2A3, G0S2, S100A12, CSF3R
       PTGS2, MMP9, CXCL8, FAM177B, EMP1, RAB31, FFAR2, FCGR2A, SMAP2, GCA
PC_ 5
Positive:  ATP8A1, PALM, SLCO2A1, FGF2, REEP4, COL4A2, RBM43, RAD51D, ACE, HEG1
       RASIP1, CYP1B1, SPARC, PCDH17, IGF2, PLEKHG1, SLC4A7, HTRA1, TCF4, AQP1
       PTPRG, ZNF419, EFEMP1, MKL2, SLFN11, PODXL, ADGRL1, TTC27, ERMARD, FN1
Negative:  MSX2, MKI67, LPXN, ADCY7, TOP2A, RNF214, NUSAP1, SNAI2, UBE2C, CCNA2
       CGRRF1, KIF14, CXCL14, MIS18BP1, CEP55, NCAPD2, NEURL1B, RACGAP1, CAD, RFWD3
       HMMR, PRODH, UBE2T, UBFD1, NPIPA1, DIAPH3, NGEF, LGALS1, NPIPA3, MCM7
```

step 10: visualize PCA data

visualizing the 5 PC generated above by using 'dims 1:5'

Hide

```
VizDimLoadings(object = pdac1, dims = 1:5, reduction = "pca")
```

**the results of visualizing the variation from PC1 indicate that there could potentially be at least 2 different types of cells that are present as the data clusters in 2 groups**

step 11: PCA heatmaps

dims = 1:10 specfies that I want 10 PCs and cells = 200 is for 200 cells

Hide

```
DimHeatmap(object = pdac1, dims = 1:10, cells = 200, balanced = TRUE)
```

step 12: dimensionality

num.replicate = 100 specifies that I want 100 replicates, and dims = 1:20 computes p-values for first 20 PCs

| | orig.ident<br><fctr> | nCount_R...<br><dbl> | nFeature_RNA<br><int> | percent.mt<br><dbl> | RNA_snn_res.0.5<br><fctr> | seurat_clus<br><fctr> |
|---|---|---|---|---|---|---|
| cell_091.133 | cell | 32176 | 1993 | 0.3170065 | 1 | 1 |
| cell_143.259 | cell | 6459 | 2233 | 4.2885896 | 1 | 1 |
| cell_057.153 | cell | 5026 | 2247 | 3.0839634 | 0 | 0 |
| cell_169.242 | cell | 5416 | 2398 | 4.8929099 | 0 | 0 |
| cell_058.010 | cell | 6806 | 2301 | 2.7622686 | 3 | 3 |
| cell_189.209 | cell | 4386 | 2119 | 2.8727770 | 0 | 0 |
| cell_023.156 | cell | 5060 | 2273 | 4.8616601 | 0 | 0 |
| cell_018.217 | cell | 5391 | 2212 | 2.5041736 | 0 | 0 |
| cell_312.307 | cell | 5615 | 2432 | 4.3811220 | 0 | 0 |
| cell_318.085 | cell | 8606 | 2477 | 4.5317221 | 3 | 3 |

1-10 of 10 rows

plotting the results from JackStraw

This code below plots the results from JackStraw and gives the resulting p-value, which can tell us which PCs are showing significant variation in our data

Hide

```
JackStrawPlot(object = pdac1, dims = 1:20)
```



- PC 2: 3.4e-33
- PC 3: 7.1e-11
- PC 4: 1.21e-05
- PC 5: 0.479
- PC 6: 0.0005
- PC 7: 0.00759
- PC 8: 0.0735
- PC 9: 0.00147
- PC 10: 1
- PC 11: 1
- PC 12: 0.0232
- PC 13: 6.33e-08
- PC 14: 7.14e-06
- PC 15: 0.133
- PC 16: 1
- PC 17: 0.0132
- PC 18: 0.133
- PC 19: 0.479
- PC 20: 1

graphing the results with Elbow plot:

Allows us to look at the amount of SD explained by each PC

Hide

```
ElbowPlot(object = pdac1)
```

step 13: clustering

first 9 PCs capture most of the variation, so we can use that by specifying 1:9. The resolution = 0.5 can be increased if we want more clusters (or vice versa)

Hide

```
pdac1 <- FindNeighbors(object = pdac1, dims = 1:9)
```

```
Computing nearest neighbor graph
Computing SNN
```

Hide

```
pdac1 <- FindClusters(object = pdac1, resolution = 0.5)
```

```
Modularity Optimizer version 1.3.0 by Ludo Waltman and Nees Jan van Eck

Number of nodes: 212
Number of edges: 5475

Running Louvain algorithm...
```

```
0%   10   20   30   40   50   60   70   80   90   100%
[----|----|----|----|----|----|----|----|----|----|
**************************************************|
```

```
Maximum modularity in 10 random starts: 0.7875
Number of communities: 4
Elapsed time: 0 seconds
```

step 14: UMAP on first 9 dimensions

installing UMAP; running first 9 dimensions; plotting the first 9 dimensions



**4 distinct clusters are visible in the UMAP (visibly there are 5 but two of them are grouped in 1)**

step 15: identify markers

the minimum number of cells that express the gene in at least 1 cluster is set at 25% (as per the Seurat markdown - I wasn't sure if this is something that needs to be altered). The log fold change threshold for the gene to be considered is also set at 0.25. In slice_max, I put n=1, specifying that I want 1 marker per cluster, which will then be used for plotting.

Hide

```
pdac1.markers <- FindAllMarkers(object = pdac1, only.pos = TRUE, min.pct = 0.25, logfc.t
hreshold = 0.25)
```

```
Calculating cluster 0
```

```
|                                              | 0 % ~calculating
|+                                             | 1 % ~04s
|++                                            | 2 % ~03s
|++                                            | 3 % ~03s
|+++                                           | 4 % ~03s
|+++                                           | 5 % ~03s
|++++                                          | 6 % ~03s
|++++                                          | 7 % ~03s
|+++++                                         | 8 % ~03s
|+++++                                         | 9 % ~03s
|++++++                                        | 10% ~03s
|++++++                                        | 11% ~03s
|+++++++                                       | 12% ~03s
|+++++++                                       | 14% ~03s
|++++++++                                      | 15% ~03s
|++++++++                                      | 16% ~03s
|+++++++++                                     | 17% ~03s
|+++++++++                                     | 18% ~03s
|++++++++++                                    | 19% ~03s
|++++++++++                                    | 20% ~03s
|+++++++++++                                   | 21% ~03s
|+++++++++++                                   | 22% ~02s
|++++++++++++                                  | 23% ~02s
|++++++++++++                                  | 24% ~02s
|+++++++++++++                                 | 25% ~02s
|+++++++++++++                                 | 26% ~02s
|+++++++++++++                                 | 27% ~02s
|++++++++++++++                                | 28% ~02s
|++++++++++++++                                | 29% ~02s
|+++++++++++++++                               | 30% ~02s
|+++++++++++++++                               | 31% ~02s
|++++++++++++++++                              | 32% ~02s
|++++++++++++++++                              | 33% ~02s
|+++++++++++++++++                             | 34% ~02s
|+++++++++++++++++                             | 35% ~02s
|++++++++++++++++++                            | 36% ~02s
|++++++++++++++++++                            | 38% ~02s
|+++++++++++++++++++                           | 39% ~02s
|+++++++++++++++++++                           | 40% ~02s
|++++++++++++++++++++                          | 41% ~02s
|++++++++++++++++++++                          | 42% ~02s
|+++++++++++++++++++++                         | 43% ~02s
|+++++++++++++++++++++                         | 44% ~02s
|++++++++++++++++++++++                        | 45% ~02s
|++++++++++++++++++++++                        | 46% ~02s
|+++++++++++++++++++++++                       | 47% ~02s
|+++++++++++++++++++++++                       | 48% ~02s
|++++++++++++++++++++++++                      | 49% ~02s
|++++++++++++++++++++++++                      | 50% ~02s
|+++++++++++++++++++++++++                     | 51% ~01s
|+++++++++++++++++++++++++++                   | 52% ~01s
```

```
|++++++++++++++++++++++++++++                      | 53% ~01s
|+++++++++++++++++++++++++++++                     | 54% ~01s
|+++++++++++++++++++++++++++++                     | 55% ~01s
|++++++++++++++++++++++++++++++                    | 56% ~01s
|++++++++++++++++++++++++++++++                    | 57% ~01s
|+++++++++++++++++++++++++++++++                   | 58% ~01s
|+++++++++++++++++++++++++++++++                   | 59% ~01s
|++++++++++++++++++++++++++++++++                  | 60% ~01s
|++++++++++++++++++++++++++++++++                  | 61% ~01s
|+++++++++++++++++++++++++++++++++                 | 62% ~01s
|+++++++++++++++++++++++++++++++++                 | 64% ~01s
|++++++++++++++++++++++++++++++++++                | 65% ~01s
|++++++++++++++++++++++++++++++++++                | 66% ~01s
|+++++++++++++++++++++++++++++++++++               | 67% ~01s
|+++++++++++++++++++++++++++++++++++               | 68% ~01s
|++++++++++++++++++++++++++++++++++++              | 69% ~01s
|++++++++++++++++++++++++++++++++++++              | 70% ~01s
|+++++++++++++++++++++++++++++++++++++             | 71% ~01s
|+++++++++++++++++++++++++++++++++++++             | 72% ~01s
|++++++++++++++++++++++++++++++++++++++            | 73% ~01s
|++++++++++++++++++++++++++++++++++++++            | 74% ~01s
|+++++++++++++++++++++++++++++++++++++++           | 75% ~01s
|+++++++++++++++++++++++++++++++++++++++           | 76% ~01s
|++++++++++++++++++++++++++++++++++++++++          | 77% ~01s
|++++++++++++++++++++++++++++++++++++++++          | 78% ~01s
|+++++++++++++++++++++++++++++++++++++++++         | 79% ~01s
|+++++++++++++++++++++++++++++++++++++++++         | 80% ~01s
|++++++++++++++++++++++++++++++++++++++++++        | 81% ~01s
|++++++++++++++++++++++++++++++++++++++++++        | 82% ~01s
|+++++++++++++++++++++++++++++++++++++++++++       | 83% ~01s
|+++++++++++++++++++++++++++++++++++++++++++       | 84% ~00s
|++++++++++++++++++++++++++++++++++++++++++++      | 85% ~00s
|++++++++++++++++++++++++++++++++++++++++++++      | 86% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++     | 88% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++     | 89% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++    | 90% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++    | 91% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++   | 92% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++   | 93% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++  | 94% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++  | 95% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++ | 96% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++ | 97% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++++| 98% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++++| 99% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++++| 100% elapsed=03s
```

```
Calculating cluster 1
```

```
  |                                                 |   0 % ~calculating
  |+                                                |   1 % ~02s
  |++                                               |   2 % ~01s
  |++                                               |   3 % ~01s
  |+++                                              |   5 % ~01s
  |+++                                              |   6 % ~01s
  |++++                                             |   7 % ~01s
  |++++                                             |   8 % ~01s
  |+++++                                            |   9 % ~01s
  |++++++                                           |  10% ~01s
  |++++++                                           |  11% ~01s
  |+++++++                                          |  12% ~01s
  |+++++++                                          |  14% ~01s
  |++++++++                                         |  15% ~01s
  |++++++++                                         |  16% ~01s
  |+++++++++                                        |  17% ~01s
  |++++++++++                                       |  18% ~01s
  |++++++++++                                       |  19% ~01s
  |+++++++++++                                      |  20% ~01s
  |+++++++++++                                      |  22% ~01s
  |++++++++++++                                     |  23% ~01s
  |++++++++++++                                     |  24% ~01s
  |+++++++++++++                                    |  25% ~01s
  |++++++++++++++                                   |  26% ~01s
  |++++++++++++++                                   |  27% ~01s
  |+++++++++++++++                                  |  28% ~01s
  |+++++++++++++++                                  |  30% ~01s
  |++++++++++++++++                                 |  31% ~01s
  |++++++++++++++++                                 |  32% ~01s
  |+++++++++++++++++                                |  33% ~01s
  |+++++++++++++++++                                |  34% ~01s
  |++++++++++++++++++                               |  35% ~01s
  |++++++++++++++++++                               |  36% ~01s
  |+++++++++++++++++++                              |  38% ~01s
  |+++++++++++++++++++                              |  39% ~01s
  |++++++++++++++++++++                             |  40% ~01s
  |++++++++++++++++++++                             |  41% ~01s
  |+++++++++++++++++++++                            |  42% ~00s
  |+++++++++++++++++++++                            |  43% ~00s
  |++++++++++++++++++++++                           |  44% ~00s
  |++++++++++++++++++++++                           |  45% ~00s
  |+++++++++++++++++++++++                          |  47% ~00s
  |+++++++++++++++++++++++                          |  48% ~00s
  |++++++++++++++++++++++++                         |  49% ~00s
  |++++++++++++++++++++++++                         |  50% ~00s
  |+++++++++++++++++++++++++                         |  51% ~00s
  |++++++++++++++++++++++++++                        |  52% ~00s
  |++++++++++++++++++++++++++                        |  53% ~00s
  |+++++++++++++++++++++++++++                       |  55% ~00s
  |+++++++++++++++++++++++++++                       |  56% ~00s
  |++++++++++++++++++++++++++++                      |  57% ~00s
```

```
|++++++++++++++++++++++++++++++                       | 58% ~00s
|++++++++++++++++++++++++++++++                       | 59% ~00s
|++++++++++++++++++++++++++++++++                     | 60% ~00s
|++++++++++++++++++++++++++++++++                     | 61% ~00s
|+++++++++++++++++++++++++++++++++                    | 62% ~00s
|+++++++++++++++++++++++++++++++++                    | 64% ~00s
|++++++++++++++++++++++++++++++++++                   | 65% ~00s
|++++++++++++++++++++++++++++++++++                   | 66% ~00s
|+++++++++++++++++++++++++++++++++++                  | 67% ~00s
|+++++++++++++++++++++++++++++++++++                  | 68% ~00s
|+++++++++++++++++++++++++++++++++++                  | 69% ~00s
|++++++++++++++++++++++++++++++++++++                 | 70% ~00s
|++++++++++++++++++++++++++++++++++++                 | 72% ~00s
|+++++++++++++++++++++++++++++++++++++                | 73% ~00s
|+++++++++++++++++++++++++++++++++++++                | 74% ~00s
|++++++++++++++++++++++++++++++++++++++               | 75% ~00s
|++++++++++++++++++++++++++++++++++++++               | 76% ~00s
|++++++++++++++++++++++++++++++++++++++               | 77% ~00s
|+++++++++++++++++++++++++++++++++++++++              | 78% ~00s
|+++++++++++++++++++++++++++++++++++++++              | 80% ~00s
|++++++++++++++++++++++++++++++++++++++++             | 81% ~00s
|++++++++++++++++++++++++++++++++++++++++             | 82% ~00s
|+++++++++++++++++++++++++++++++++++++++++            | 83% ~00s
|+++++++++++++++++++++++++++++++++++++++++            | 84% ~00s
|+++++++++++++++++++++++++++++++++++++++++            | 85% ~00s
|++++++++++++++++++++++++++++++++++++++++++           | 86% ~00s
|++++++++++++++++++++++++++++++++++++++++++           | 88% ~00s
|+++++++++++++++++++++++++++++++++++++++++++          | 89% ~00s
|+++++++++++++++++++++++++++++++++++++++++++          | 90% ~00s
|++++++++++++++++++++++++++++++++++++++++++++         | 91% ~00s
|++++++++++++++++++++++++++++++++++++++++++++         | 92% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++        | 93% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++        | 94% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++       | 95% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++       | 97% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++      | 98% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++    | 99% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++++++| 100% elapsed=01s
```

```
Calculating cluster 2
```

```
|                                                                    |  0 % ~calculating
|+                                                                   |  1 % ~01s
|++                                                                  |  3 % ~00s
|++                                                                  |  4 % ~00s
|+++                                                                 |  5 % ~00s
|++++                                                                |  6 % ~00s
|++++                                                                |  8 % ~01s
|+++++                                                               |  9 % ~01s
|++++++                                                              | 10% ~01s
|++++++                                                              | 12% ~01s
|+++++++                                                             | 13% ~01s
|++++++++                                                            | 14% ~01s
|++++++++                                                            | 15% ~01s
|+++++++++                                                           | 17% ~01s
|+++++++++                                                           | 18% ~01s
|++++++++++                                                          | 19% ~01s
|+++++++++++                                                         | 21% ~01s
|+++++++++++                                                         | 22% ~01s
|++++++++++++                                                        | 23% ~01s
|+++++++++++++                                                       | 24% ~01s
|+++++++++++++                                                       | 26% ~01s
|++++++++++++++                                                      | 27% ~01s
|+++++++++++++++                                                     | 28% ~01s
|+++++++++++++++                                                     | 29% ~00s
|++++++++++++++++                                                    | 31% ~00s
|+++++++++++++++++                                                   | 32% ~00s
|+++++++++++++++++                                                   | 33% ~00s
|++++++++++++++++++                                                  | 35% ~00s
|++++++++++++++++++                                                  | 36% ~00s
|+++++++++++++++++++                                                 | 37% ~00s
|++++++++++++++++++++                                                | 38% ~00s
|++++++++++++++++++++                                                | 40% ~00s
|+++++++++++++++++++++                                               | 41% ~00s
|++++++++++++++++++++++                                              | 42% ~00s
|++++++++++++++++++++++                                              | 44% ~00s
|+++++++++++++++++++++++                                             | 45% ~00s
|+++++++++++++++++++++++                                             | 46% ~00s
|++++++++++++++++++++++++                                            | 47% ~00s
|+++++++++++++++++++++++++                                           | 49% ~00s
|+++++++++++++++++++++++++                                           | 50% ~00s
|++++++++++++++++++++++++++                                          | 51% ~00s
|+++++++++++++++++++++++++++                                         | 53% ~00s
|+++++++++++++++++++++++++++                                         | 54% ~00s
|++++++++++++++++++++++++++++                                        | 55% ~00s
|+++++++++++++++++++++++++++++                                       | 56% ~00s
|+++++++++++++++++++++++++++++                                       | 58% ~00s
|++++++++++++++++++++++++++++++                                      | 59% ~00s
|+++++++++++++++++++++++++++++++                                     | 60% ~00s
|+++++++++++++++++++++++++++++++                                     | 62% ~00s
|++++++++++++++++++++++++++++++++                                    | 63% ~00s
|+++++++++++++++++++++++++++++++++                                   | 64% ~00s
```

```
|++++++++++++++++++++++++++++++++++               | 65% ~00s
|++++++++++++++++++++++++++++++++++               | 67% ~00s
|++++++++++++++++++++++++++++++++++               | 68% ~00s
|+++++++++++++++++++++++++++++++++++              | 69% ~00s
|+++++++++++++++++++++++++++++++++++++            | 71% ~00s
|+++++++++++++++++++++++++++++++++++++            | 72% ~00s
|++++++++++++++++++++++++++++++++++++++           | 73% ~00s
|+++++++++++++++++++++++++++++++++++++++          | 74% ~00s
|+++++++++++++++++++++++++++++++++++++++          | 76% ~00s
|++++++++++++++++++++++++++++++++++++++++         | 77% ~00s
|++++++++++++++++++++++++++++++++++++++++++       | 78% ~00s
|++++++++++++++++++++++++++++++++++++++++++       | 79% ~00s
|+++++++++++++++++++++++++++++++++++++++++++      | 81% ~00s
|++++++++++++++++++++++++++++++++++++++++++++     | 82% ~00s
|++++++++++++++++++++++++++++++++++++++++++++     | 83% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++    | 85% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++    | 86% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++   | 87% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++   | 88% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++  | 90% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++  | 91% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++ | 92% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++ | 94% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 95% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 96% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 97% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 99% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 100% elapsed=01s
```

```
Calculating cluster 3
```

```
|                                                        |  0 % ~calculating
|+                                                       |  1 % ~02s
|++                                                      |  2 % ~02s
|++                                                      |  3 % ~02s
|+++                                                     |  4 % ~02s
|+++                                                     |  5 % ~02s
|++++                                                    |  6 % ~02s
|++++                                                    |  7 % ~02s
|+++++                                                   |  8 % ~02s
|+++++                                                   |  9 % ~02s
|++++++                                                  | 10% ~02s
|++++++                                                  | 11% ~02s
|+++++++                                                 | 12% ~02s
|+++++++                                                 | 13% ~02s
|++++++++                                                | 14% ~02s
|++++++++                                                | 15% ~02s
|+++++++++                                               | 16% ~02s
|+++++++++                                               | 17% ~02s
|++++++++++                                              | 18% ~02s
|++++++++++                                              | 19% ~02s
|+++++++++++                                             | 20% ~02s
|+++++++++++                                             | 21% ~02s
|++++++++++++                                            | 22% ~02s
|++++++++++++                                            | 23% ~02s
|+++++++++++++                                           | 24% ~02s
|+++++++++++++                                           | 25% ~02s
|++++++++++++++                                          | 26% ~02s
|++++++++++++++                                          | 27% ~02s
|+++++++++++++++                                         | 28% ~02s
|+++++++++++++++                                         | 29% ~02s
|++++++++++++++++                                        | 30% ~02s
|++++++++++++++++                                        | 31% ~02s
|+++++++++++++++++                                       | 32% ~02s
|+++++++++++++++++                                       | 33% ~02s
|++++++++++++++++++                                      | 34% ~02s
|++++++++++++++++++                                      | 35% ~02s
|+++++++++++++++++++                                     | 36% ~02s
|+++++++++++++++++++                                     | 37% ~02s
|++++++++++++++++++++                                    | 38% ~02s
|++++++++++++++++++++                                    | 39% ~02s
|+++++++++++++++++++++                                   | 40% ~02s
|+++++++++++++++++++++                                   | 41% ~02s
|++++++++++++++++++++++                                  | 42% ~01s
|++++++++++++++++++++++                                  | 43% ~01s
|+++++++++++++++++++++++                                 | 44% ~01s
|+++++++++++++++++++++++                                 | 45% ~01s
|++++++++++++++++++++++++                                | 46% ~01s
|++++++++++++++++++++++++                                | 47% ~01s
|+++++++++++++++++++++++++                               | 48% ~01s
|+++++++++++++++++++++++++                               | 49% ~01s
|++++++++++++++++++++++++++                              | 51% ~01s
```

```
|+++++++++++++++++++++++++++                      |  52% ~01s
|+++++++++++++++++++++++++++                      |  53% ~01s
|+++++++++++++++++++++++++++                      |  54% ~01s
|++++++++++++++++++++++++++++                     |  55% ~01s
|++++++++++++++++++++++++++++                     |  56% ~01s
|+++++++++++++++++++++++++++++                    |  57% ~01s
|+++++++++++++++++++++++++++++                    |  58% ~01s
|++++++++++++++++++++++++++++++                   |  59% ~01s
|++++++++++++++++++++++++++++++                   |  60% ~01s
|+++++++++++++++++++++++++++++++                  |  61% ~01s
|+++++++++++++++++++++++++++++++                  |  62% ~01s
|++++++++++++++++++++++++++++++++                 |  63% ~01s
|++++++++++++++++++++++++++++++++                 |  64% ~01s
|+++++++++++++++++++++++++++++++++                |  65% ~01s
|+++++++++++++++++++++++++++++++++                |  66% ~01s
|++++++++++++++++++++++++++++++++++               |  67% ~01s
|++++++++++++++++++++++++++++++++++               |  68% ~01s
|+++++++++++++++++++++++++++++++++++              |  69% ~01s
|+++++++++++++++++++++++++++++++++++              |  70% ~01s
|++++++++++++++++++++++++++++++++++++             |  71% ~01s
|++++++++++++++++++++++++++++++++++++             |  72% ~01s
|+++++++++++++++++++++++++++++++++++++            |  73% ~01s
|+++++++++++++++++++++++++++++++++++++            |  74% ~01s
|++++++++++++++++++++++++++++++++++++++           |  75% ~01s
|++++++++++++++++++++++++++++++++++++++           |  76% ~01s
|+++++++++++++++++++++++++++++++++++++++          |  77% ~01s
|+++++++++++++++++++++++++++++++++++++++          |  78% ~01s
|++++++++++++++++++++++++++++++++++++++++         |  79% ~01s
|++++++++++++++++++++++++++++++++++++++++         |  80% ~01s
|+++++++++++++++++++++++++++++++++++++++++        |  81% ~00s
|+++++++++++++++++++++++++++++++++++++++++        |  82% ~00s
|++++++++++++++++++++++++++++++++++++++++++       |  83% ~00s
|++++++++++++++++++++++++++++++++++++++++++       |  84% ~00s
|+++++++++++++++++++++++++++++++++++++++++++      |  85% ~00s
|+++++++++++++++++++++++++++++++++++++++++++      |  86% ~00s
|++++++++++++++++++++++++++++++++++++++++++++     |  87% ~00s
|++++++++++++++++++++++++++++++++++++++++++++     |  88% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++    |  89% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++    |  90% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++   |  91% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++   |  92% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++  |  93% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++  |  94% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++ |  95% ~00s
|++++++++++++++++++++++++++++++++++++++++++++++++ |  96% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 97% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 98% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 99% ~00s
|+++++++++++++++++++++++++++++++++++++++++++++++++| 100% elapsed=02s
```
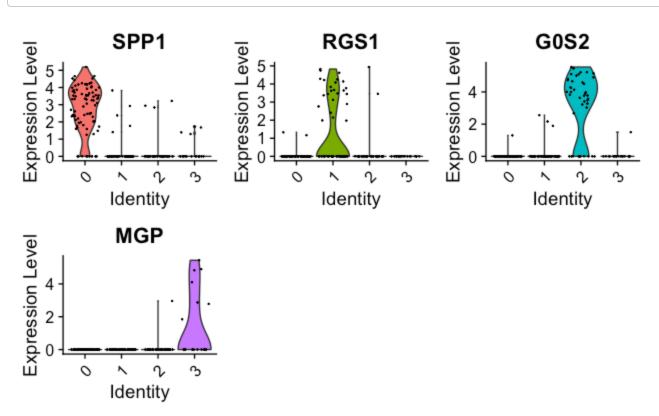
Hide

```
pdac1.markers %>%
    group_by(cluster) %>%
    slice_max(n = 1, order_by = avg_log2FC)
```

| p_val | avg_log2FC | pct.1 | pct.2 | p_val_adj | cluster | gene |
|---|---|---|---|---|---|---|
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <fctr> | <chr> |
| 1.260561e-30 | 3.830777 | 0.882 | 0.094 | 1.669991e-26 | 0 | SPP1 |
| 1.053404e-11 | 2.997340 | 0.369 | 0.027 | 1.395550e-07 | 1 | RGS1 |
| 1.967397e-29 | 5.863167 | 0.744 | 0.030 | 2.606407e-25 | 2 | G0S2 |
| 2.631740e-15 | 4.856750 | 0.368 | 0.005 | 3.486529e-11 | 3 | MGP |

4 rows

step 16: violin plot from 1 feature from each cluster

I used the features obtained above to create this violon plot to see the distribution

Hide

```
VlnPlot(object = pdac1, features = c("SPP1", "RGS1", "G0S2", "MGP"))
```



step 17: feature plot with the same features as step 16

Feature plot has been constructed with the same features as used in the creation of the violon plot above.

Hide

```
FeaturePlot(object = pdac1, features = c("SPP1", "RGS1", "G0S2", "MGP"))
```