

# scRNA-Seq Lab

*(solution)*

2017-03-02

---

## Introduction/Quality control

For this exercise, we will be analyzing the a single cell RNA-Seq dataset of Peripheral Blood Mononuclear Cells (PBMC) from the 10X Genomics platform. We will primarily be using the `seurat` package from the Satija Lab, which includes a vignette here. <http://satijalab.org/seurat/pbmc-tutorial.html>. The purpose of single cell RNA-Seq analysis is to uncover interesting biology that occurs at a granularity—the single cell—that isn't appreciated when these features become averaged in bulk. The goal of this analysis is to uncover heterogeneity in PBMCs and understanding the analysis workflows for single cell technologies.

**First, load the packages and the data object**

```
library(Seurat)
library(Matrix)
library(dplyr)

pbmc <- readRDS("pbmc.rds")
dim(pbmc@raw.data)
```

```
## [1] 32643 2001
```

*Note:* to achieve this object, the counts matrix had to be determined using a standard alignment protocol similar to bulk RNA-Seq analyses. The `.rds` object contains a **seurat** object with 2001 samples and over 32,000 genes. This sample set includes roughly 1,000 PBMC samples from two different batches.

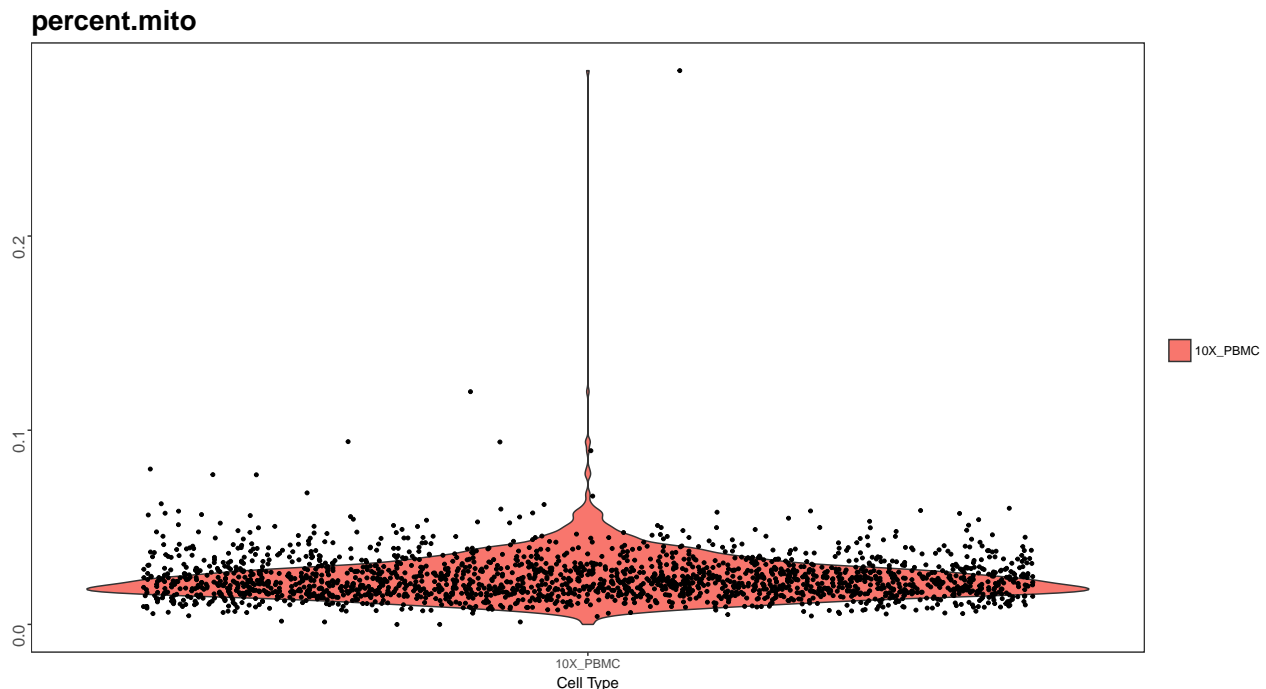
## Analysis

(1) The substantial sparsity associated with scRNA-Seq data makes analysis a unique challenge. Use the 'Setup' command to filter lowly expressed genes and weakly detected cells in this raw dataset. How many samples and genes are filtered afterwards?

```
pbmc <- Setup(pbmc, min.cells = 3, min.genes = 200, do.logNormalize = TRUE,  
             total.expr = 1e4, project = "10X_PBMC")  
# use 'str' command to determine still 2,000 samples and 12,800 genes remain after filter
```

(2) When trying to discover rare cell types, one has to be weary of technical confounders that imply heterogeneity that are actually false. Two measures of technical confounders are the number of mitochondrial reads mapped as well as the number of unique genes mapped. In this dataset, how many mitochondrial genes are there? What is the distribution of the proportion of reads for these mitochondrial genes? How many samples express a number of genes that significantly deviates from the rest?

```
mito.genes <- grep("^MT-", rownames(pbmc@data), value = TRUE)  
percent.mito <- colSums(expm1(pbmc@data[mito.genes, ]))/colSums(expm1(pbmc@data))  
  
#AddMetaData adds columns to object@data.info, and is a great place to stash QC stats  
pbmc <- AddMetaData(pbmc, percent.mito, "percent.mito")  
VlnPlot(pbmc, c("percent.mito"), nCol = 1)
```



```
pbmc <- SubsetData(pbmc, subset.name = "nGene", accept.high = 2500)  
pbmc <- SubsetData(pbmc, subset.name = "percent.mito", accept.high = 0.05)
```

## Linear/Non-linear dimensional reduction

(3) Rather than focusing on specific differentially expressed genes, a staple in scRNA-Seq analyses involves dimension reduction. Compute the top principal components using the

variable genes and determine which genes contribute most to these PCs. Find some other ways to display and interpret the results of the dimensionality reduction.

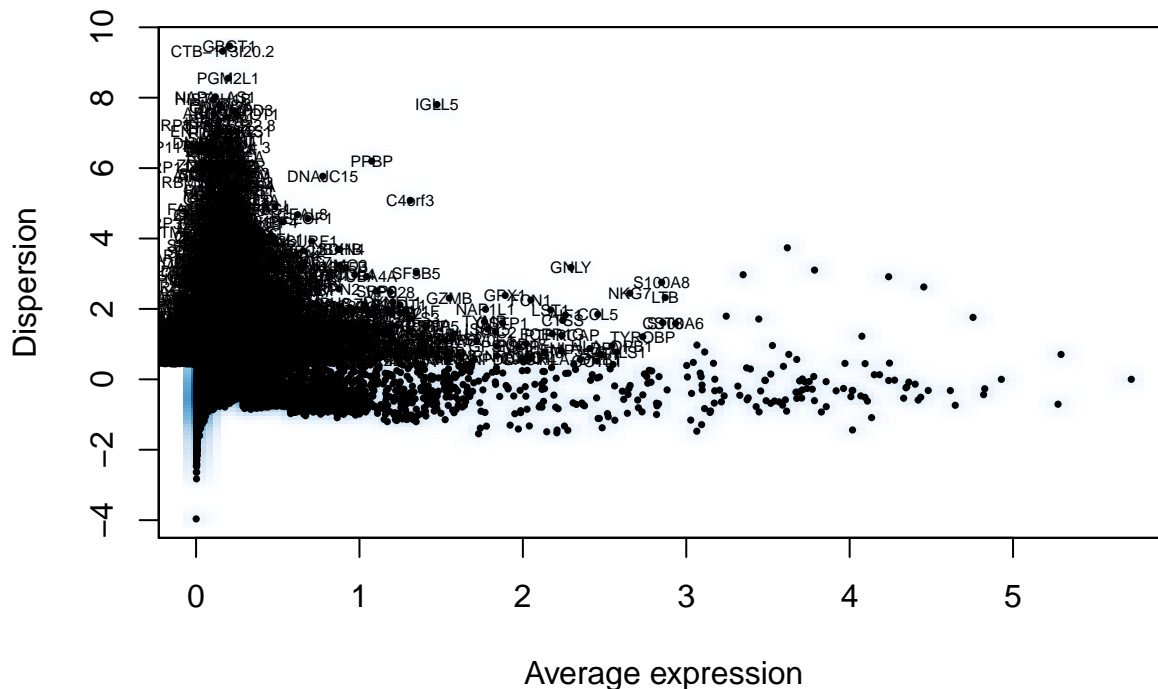
```
# Add this to compute variable genes to actually get this to work
```

```
pbmc <- MeanVarPlot(pbmc ,fxn.x = expMean, fxn.y = logVarDivMean, x.low.cutoff = 0.0125, x.high.cutoff =
```

```
## [1] "Calculating gene dispersion"
```

```
##
```

```
|
|                                     | 0%
|
|=====| 8%
|
|=====| 15%
|
|=====| 23%
|
|=====| 31%
|
|=====| 38%
|
|=====| 46%
|
|=====| 54%
|
|=====| 62%
|
|=====| 69%
|
|=====| 77%
|
|=====| 85%
|
|=====| 92%
|
|=====| 100%
```



```
pbmc <- PCA(pbmc, pc.genes = pbmc@var.genes, do.print = FALSE, pcs.print = 5, genes.print = 5)
```

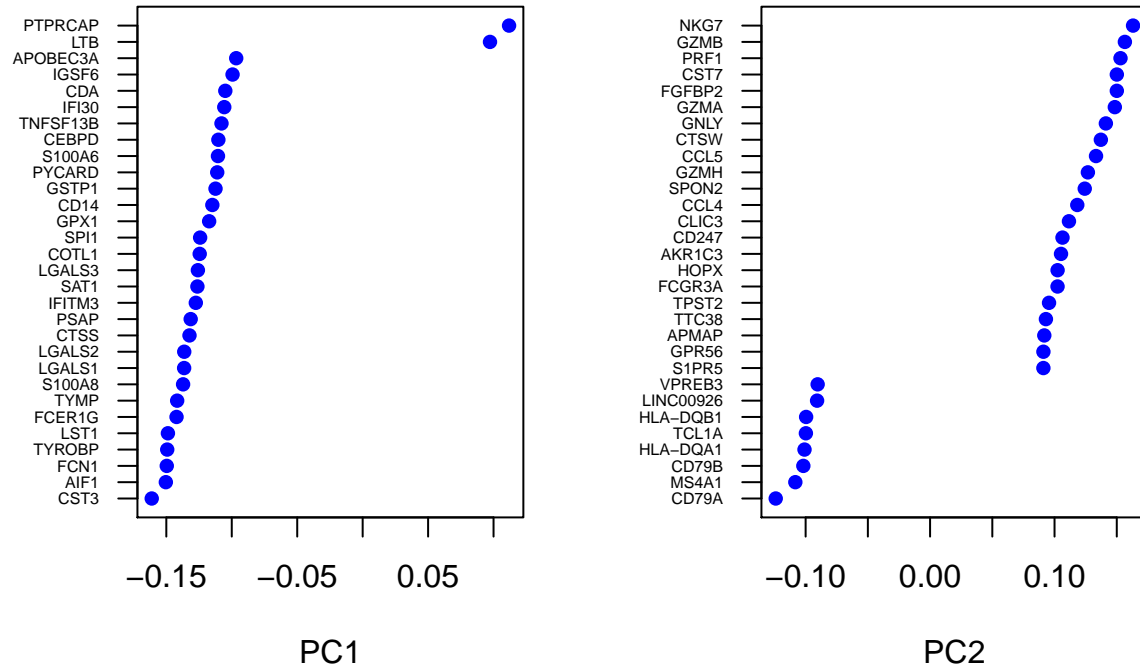
Seurat provides several useful ways of visualizing both cells and genes that define the PCA, including **PrintPCA()**, **VizPCA()**, **PCAPlot()**, and **PCHeatmap()**

```
# Examine and visualize PCA results a few different ways
PrintPCA(pbmc, pcs.print = 1:5, genes.print = 5, use.full = FALSE)
```

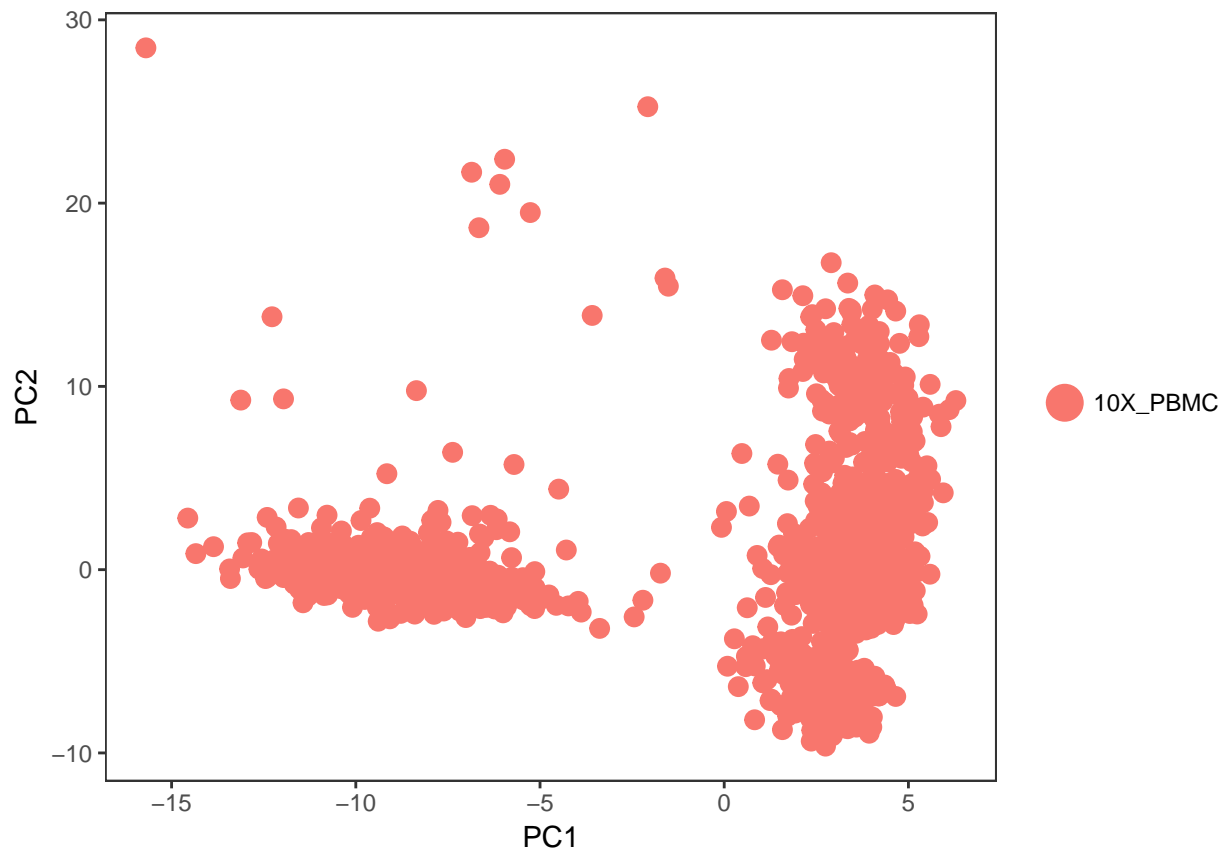
```
## [1] "PC1"
## [1] "CST3"   "AIF1"   "FCN1"   "TYROBP" "LST1"
## [1] ""
## [1] "PTPRCAP" "LTB"    "CD3D"   "CXCR4"   "AES"
## [1] ""
## [1] ""
## [1] "PC2"
## [1] "CD79A"   "MS4A1"   "CD79B"   "HLA-DQA1" "TCL1A"
## [1] ""
## [1] "NKG7"    "GZMB"    "PRF1"    "CST7"    "FGFBP2"
## [1] ""
## [1] ""
## [1] "PC3"
## [1] "NKG7"    "PRF1"    "CST7"    "FGFBP2" "GZMB"
## [1] ""
## [1] "PPBP"    "PF4"     "SDPR"    "GNG11"    "HIST1H2AC"
## [1] ""
## [1] ""
## [1] "PC4"
## [1] "CD3D"    "FYB"     "RGCC"    "CD27"    "NDFIP1"
## [1] ""
## [1] "CD79A"    "HLA-DQA1" "CD79B"    "HLA-DQB1" "HLA-DPB1"
## [1] ""
## [1] ""
## [1] "PC5"
```

```
## [1] "NAP1L1"      "LTB"          "HNRNPA2B1" "ABI3"         "PTGES3"
## [1] ""
## [1] "S100A8"      "S100A12"     "CD14"        "CLEC4E"      "GZMB"
## [1] ""
## [1] ""
```

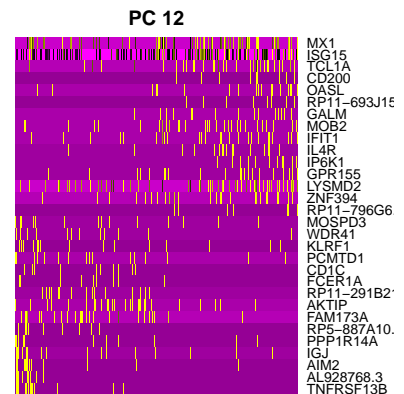
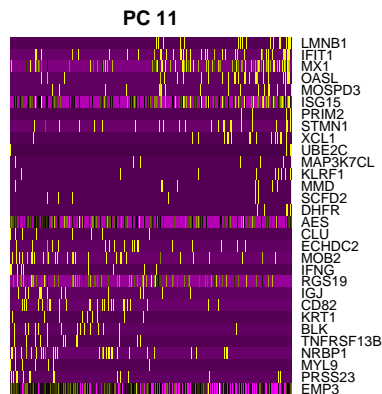
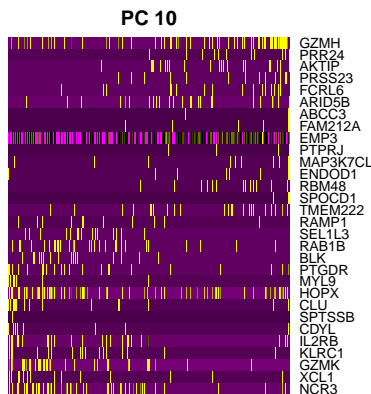
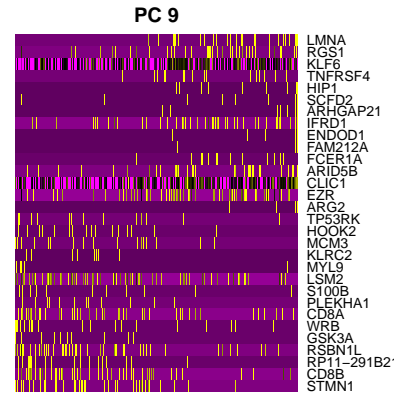
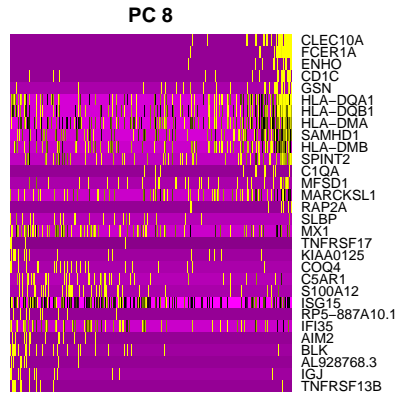
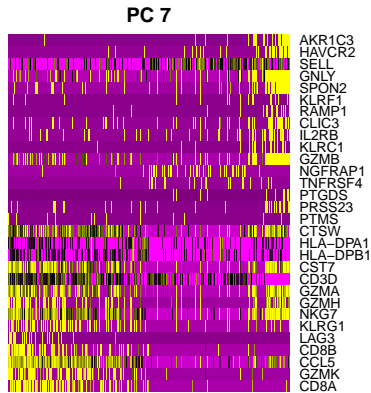
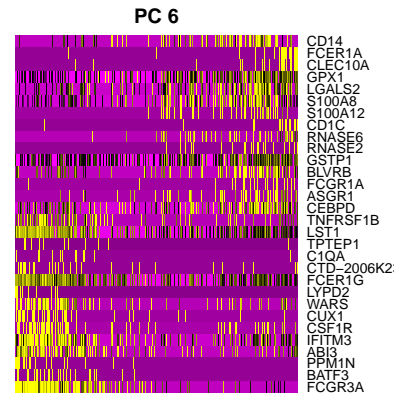
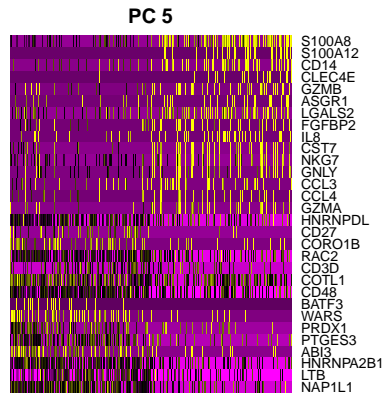
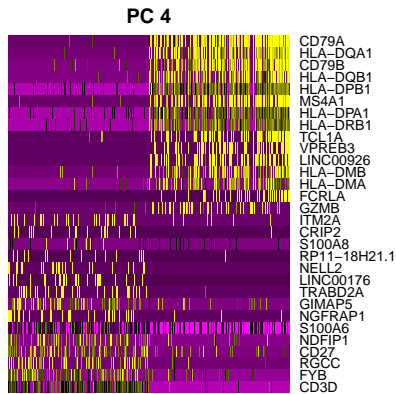
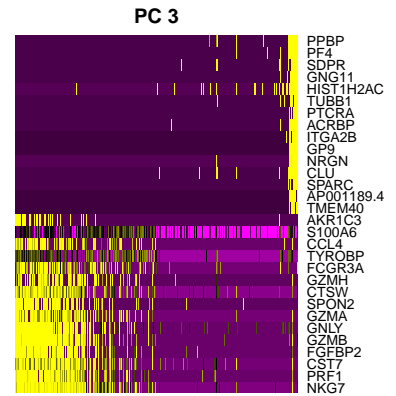
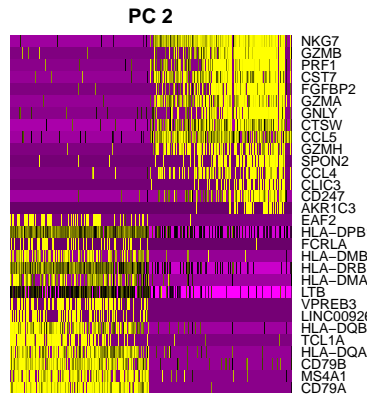
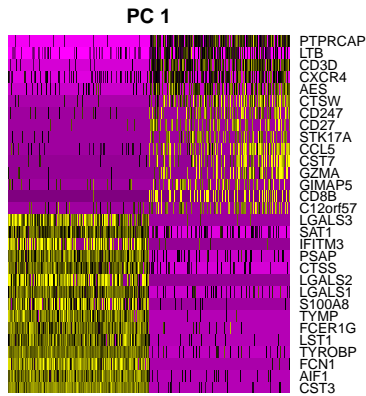
```
VizPCA(pbmc, 1:2)
```



```
PCAPlot(pbmc, 1, 2)
```

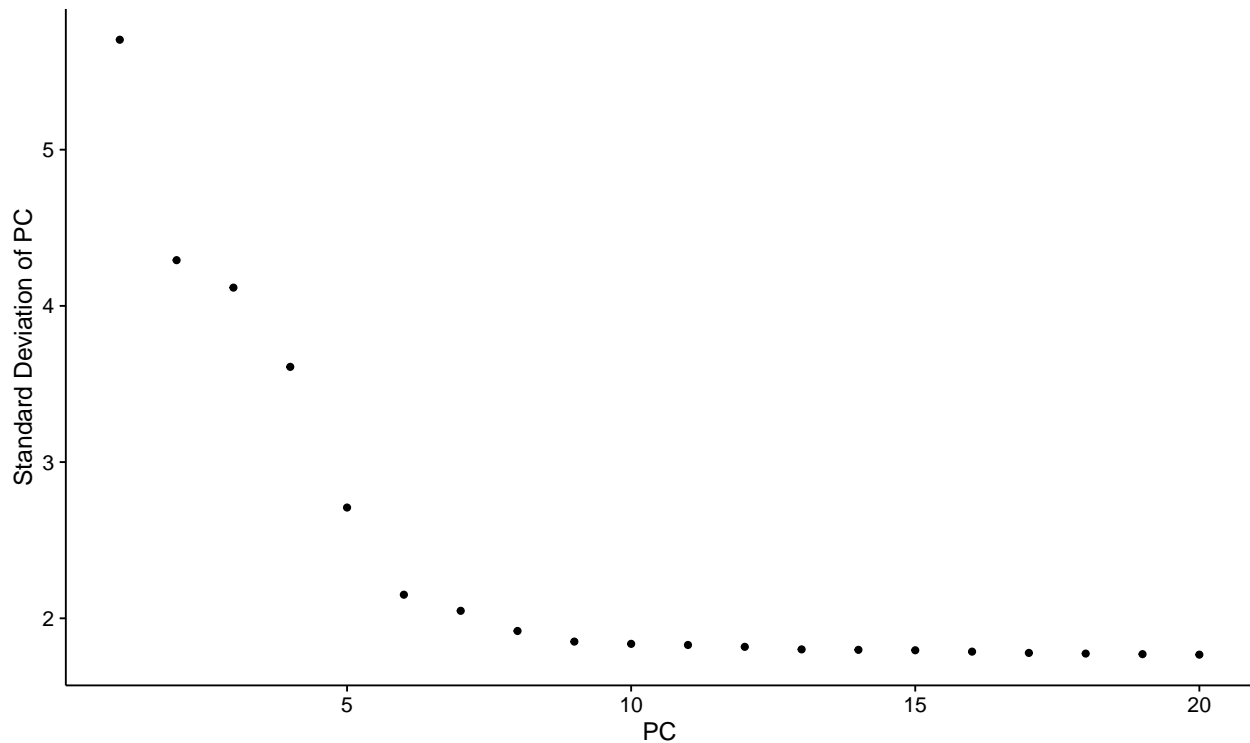


```
PCHeatmap(pbmc, pc.use = 1:12, cells.use = 500, do.balanced = TRUE,  
          label.columns = FALSE, use.full = FALSE)
```



(4) Which principal components are statistically significant? Comment on one or more approaches to determine this.

```
PCElbowPlot(pbmcc)
```



A more ad hoc method for determining which PCs to use is to look at a plot of the standard deviations of the principle components and draw your cutoff where there is a clear elbow in the graph. This can be done with `PCElbowPlot()`. In this example, it looks like the elbow would fall around PC 9.

(5) Use the `FindClusters` command to determine sample modules in the PBMC data. Comment on the type of clustering performed. Is it supervised or unsupervised?

```
pbmcc <- FindClusters(pbmcc, pc.use = 1:10, resolution = 0.6, print.output = 0, save.SNN = T)
```

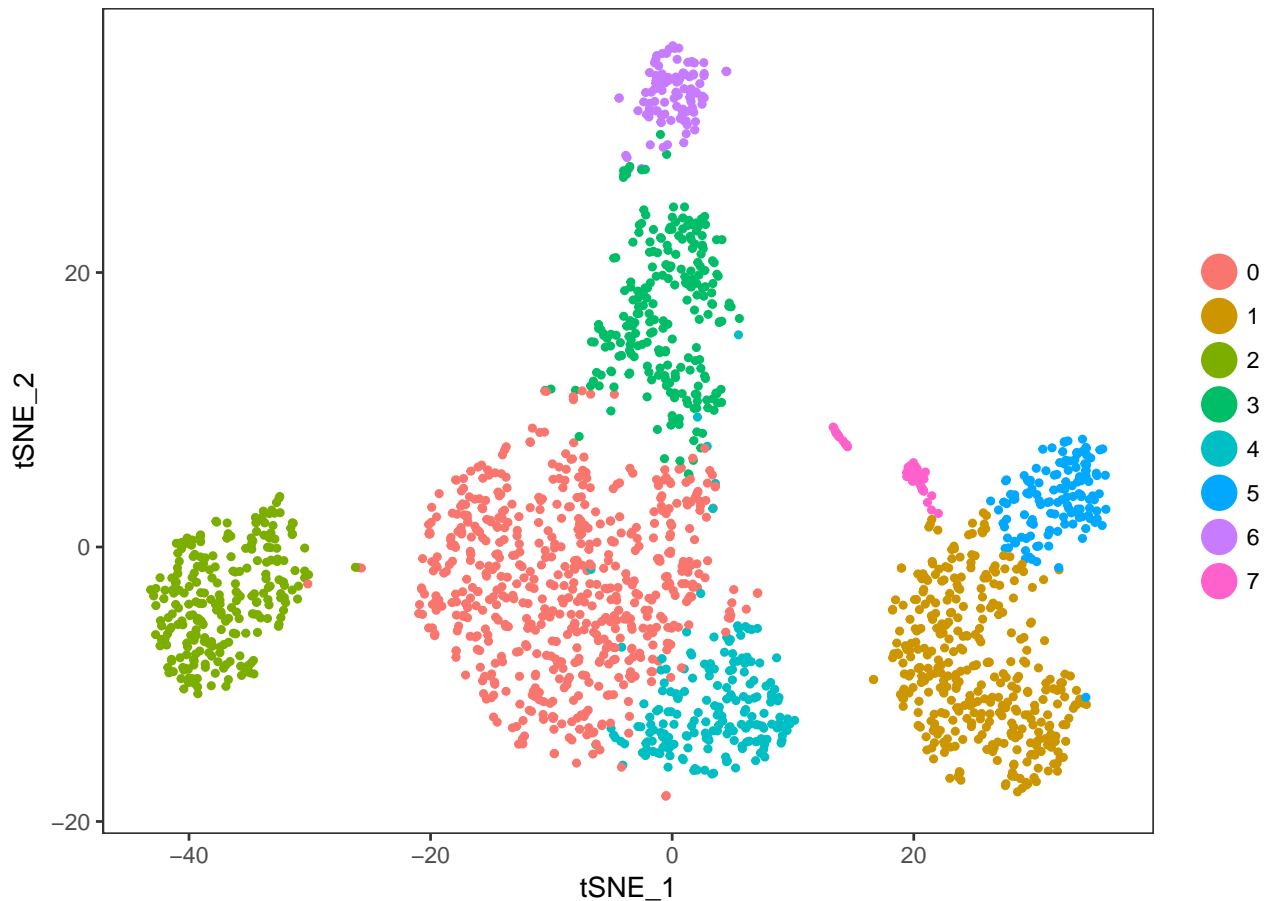
Details are in the Seurat source code as well as several paragraphs in the vignette. It is supervised.

(6) A popular method for displaying scRNA-Seq data is by creating two dimensions using tSNE. Run and visualize tSNE for this data. Comment on how this approach is different than PCA.

```
pbmcc <- RunTSNE(pbmcc, dims.use = 1:10, do.fast = TRUE)
```

```
TSNEPlot(pbmcc)
```





*Note:* tSNE is, by definition, a stochastic process so be sure to cache your data at this point or save the file image before re-running later steps! Main difference is linear/non linear effects; tSNE in this case is using the PCs as input

## Differentially expressed genes (cluster biomarkers)

(7) Now that we've defined data-driven clusters, we'd like to identify markers that define clusters via differential expression. What markers distinguish cluster 2? What markers distinguish cluster 2 from cluster 4? Every cluster from all others.

```
# find all markers of cluster 2
cluster2.markers <- FindMarkers(pbmc, ident.1 = 2, min.pct = 0.25)
print(head(cluster2.markers, 5))
```

```
##               p_val avg_diff pct.1 pct.2
## CD74      4.626575e-276 2.003223 1.000 0.754
## CD79A     3.774364e-266 2.989068 0.933 0.038
## CD79B     6.083228e-243 2.474387 0.922 0.131
## HLA-DRA   1.383353e-201 1.839170 0.996 0.464
## MS4A1     5.769042e-193 2.333413 0.816 0.045
```

```
cluster24.markers <- FindMarkers(pbmc, 2, 4, min.pct = 0.25)
print(head(cluster24.markers, 5))
```

```
##           p_val avg_diff pct.1 pct.2
## CD74      2.101969e-220 3.291178 1.000 0.520
## HLA-DRA   5.246923e-167 3.923024 0.996 0.133
## HLA-DPB1  7.946604e-123 3.046278 0.984 0.117
## HLA-DRB1  2.527228e-116 3.279697 0.973 0.026
## CD79A     1.155224e-103 3.064796 0.933 0.010

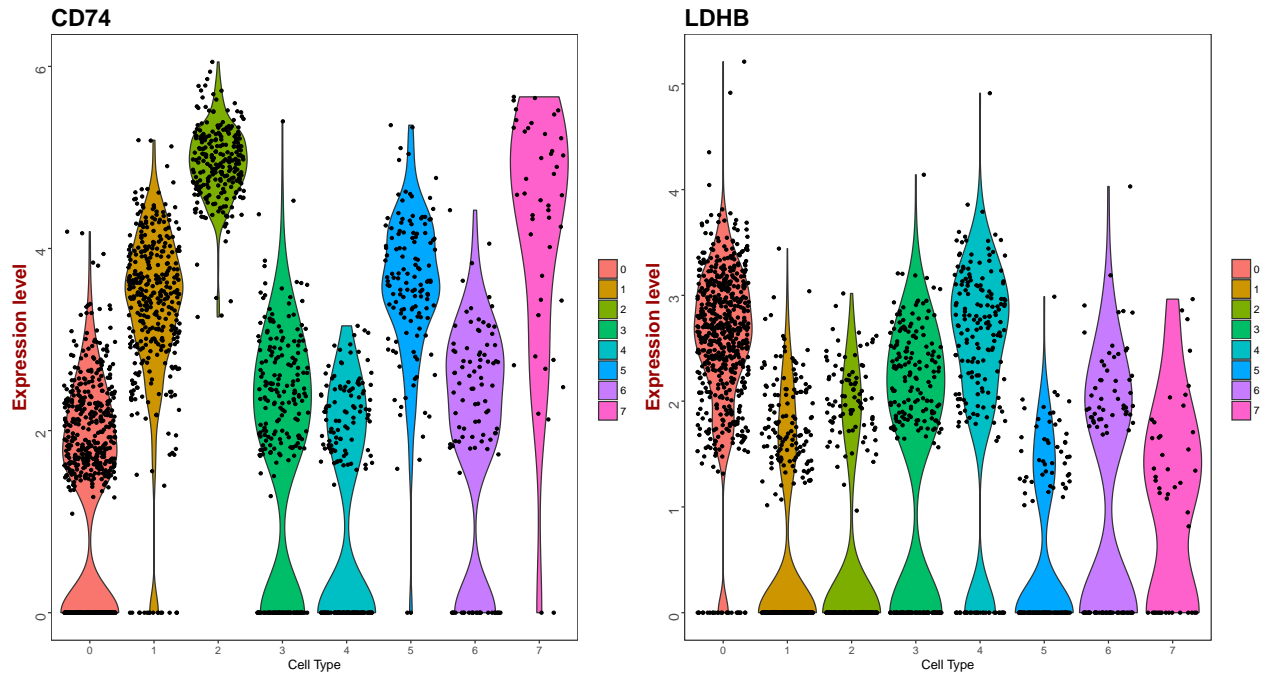
pbmc.markers <- FindAllMarkers(pbmc, only.pos = TRUE, min.pct = 0.25, thresh.use = 0.25)
pbmc.markers %>% group_by(cluster) %>% top_n(2, avg_diff)
```

```
## Source: local data frame [16 x 6]
## Groups: cluster [8]
```

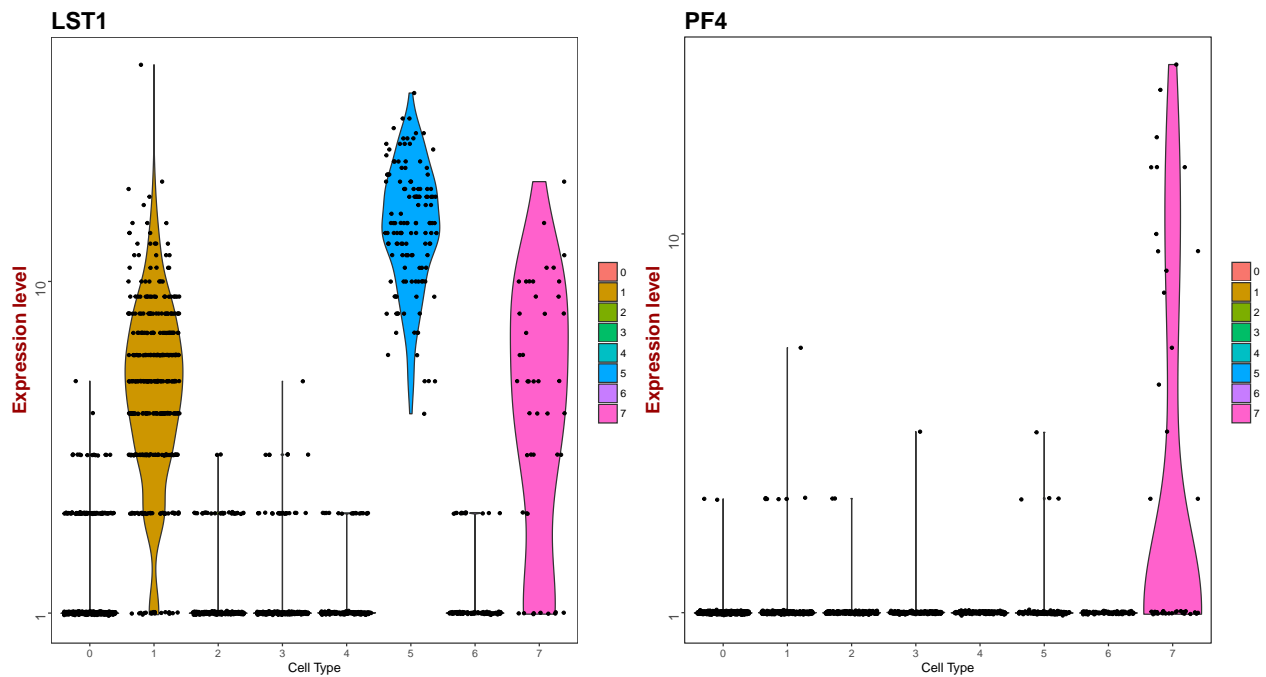
```
##           p_val  avg_diff pct.1 pct.2 cluster  gene
##           <dbl>    <dbl> <dbl> <dbl> <fctr>   <chr>
## 1  4.565200e-146 1.0133146 0.956 0.506      0  LDHB
## 2  8.547608e-137 0.8729406 0.884 0.304      0  CD3D
## 3   0.000000e+00 3.8247836 0.984 0.109      1 S100A8
## 4   0.000000e+00 3.7477009 0.978 0.205      1 S100A9
## 5  3.774364e-266 2.9890678 0.933 0.038      2  CD79A
## 6  3.053285e-139 2.5657555 0.588 0.024      2  TCL1A
## 7  1.142654e-192 2.3393136 0.964 0.192      3   CCL5
## 8   3.050984e-95 2.2246331 0.543 0.033      3  GZMK
## 9   6.368285e-14 0.7227104 0.357 0.166      4  CCR7
## 10  6.434473e-12 0.7038178 0.291 0.130      4  LEF1
## 11  2.098630e-125 1.9500161 1.000 0.309      5  LST1
## 12  1.043530e-104 2.1994271 0.940 0.121      5 FCGR3A
## 13  2.401586e-141 3.5404225 1.000 0.122      6   GNLY
## 14  2.330636e-138 3.2910876 1.000 0.076      6  GZMB
## 15  1.410915e-27 4.2288101 0.311 0.020      7  PPBP
## 16  8.260548e-26 3.3515661 0.356 0.008      7   PF4
```

(8) Using the biomarkers identified above, select a few markers to distinguish the various subgroups. Try plotting different measurements, including raw and normalized counts on/not on the log scale.

```
VlnPlot(pbmc, c("CD74", "LDHB"))
```



```
VlnPlot(pbm, c("LST1", "PF4"), use.raw = TRUE, y.log = TRUE)
```

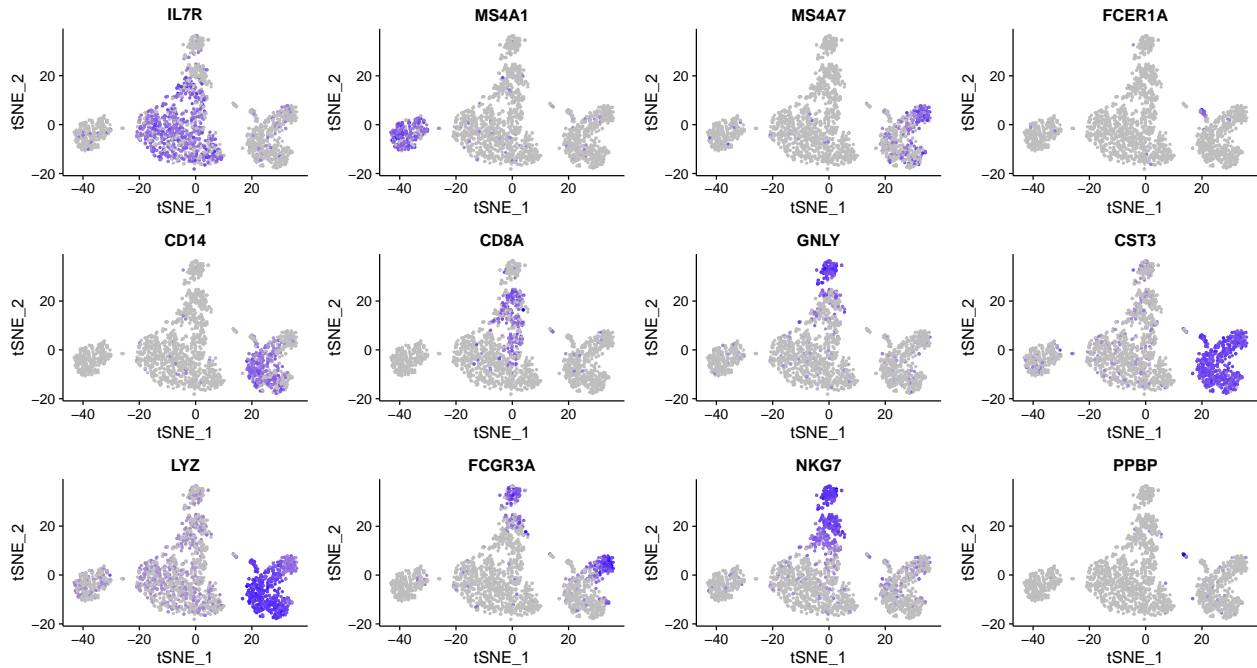


(9) Using the table below, identify which clusters correspond to which cell subtypes in your tSNE projection. Do you observe any rare populations or mixed populations? Explore some other markers to characterize the behavior of these populations.

Cluster ID	Markers	Cell Type
?	IL7R	CD4 T cells
?	CD14, LYZ	CD14+ Monocytes
?	MS4A1	B cells
?	CD8A	CD8 T cells

Cluster ID	Markers	Cell Type
?	FCGR3A, MS4A7	FCGR3A+ Monocytes
?	GNLY, NKG7	NK cells
?	FCER1A, CST3	Dendritic Cells
?	PPBP	Megakaryocytes

```
factors <- c("IL7R", "CD14", "LYZ", "MS4A1", "CD8A",
             "FCGR3A", "MS4A7", "GNLY", "NKG7", "FCER1A", "CST3", "PPBP")
FeaturePlot(pbm, factors, cols.use = c("grey", "blue"))
```



(10) Using the inference above, annotate your tSNE with the cell type names.

```
current.cluster.ids <- c(0, 1, 2, 3, 4, 5, 6, 7)
new.cluster.ids <- c("CD4 T cells", "CD14+ Monocytes", "B cells",
                    "CD8 T cells", "CD4 T cells", "FCGR3A+ Monocytes", "NK cells",
                    "Dendritic cells")
pbmc@ident <- plyr::mapvalues(pbm@ident, from = current.cluster.ids, to = new.cluster.ids)
TSNEPlot(pbm, do.label = TRUE, pt.size = 0.5)
```

