

Khan Inan
BI-GY 7633
HW 5 extra credit

*I apologize, I tried my best but I found this homework to be difficult. I hope the professor goes over the basic Seurat workflow, so I can understand it better. Thank you

Step 0: Please install and load the following packages

```
install.packages("dplyr")
install.packages("Seurat")
install.packages("patchwork")
library("dplyr")
library("Seurat")
library("patchwork")
```

Step 1: Load the TSV file

```
```{r}

pdac_data <- read.table("C:/Users/khani/Downloads/hw 5/GSM3036909.tsv",
header = TRUE, sep = "\t", row.names = 1, check.names = FALSE)
```
```

Step 2: Create a Seurat object.

```
1 library(Seurat)
2
3 pdac1 <- CreateSeuratObject(counts = pdac_data, project = "pdac1")
4 pdac1 <- FilterCells(pdac1, min.cells = 3)
5 pdac1 <- FilterCells(pdac1, min.features = 200)
```

Step 3: Label the Mitochondrial genes

```
1 pdac1 <- PercentageFeatureSet(pdac1, pattern = "^MT-")
```

Step 4: Visualize the distribution

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

Step 5: Filter data

```
pdac1 <- FilterCells(pdac1, subset.names = nFeature_RNA, low.thresholds =  
1 200, high.thresholds = 2500)  
pdac1 <- FilterCells(pdac1, subset.names = percent.mt, low.thresholds = 0,  
2 high.thresholds = 5)
```

Step 6: Normalize data

```
pdac1 <- NormalizeData(pdac1, normalization.method = "LogNormalize",  
1 scale.factor = 10000)
```

Step 6: Calculate gene variation

```
1 # Calculate gene variation  
pdac1 <- FindVariableFeatures(pdac1, selection.method = "vst", nfeatures =  
2 2000)  
3  
4 # Update Seurat object  
pdac1 <- ScaleData(pdac1, features = rownames(pdac1), display.progress =  
5 FALSE)  
6
```

Step 7: Scale data

```
1 pdac1 <- ScaleData(pdac1, vars.to.regress = "nGene")  
2
```

Step 8: PCA

```
pdac1 <- RunPCA(pdac1, features = pdac1@var.genes, ndims.print=1,  
1 verbose=FALSE)  
2
```

Step 9: Visualize data using VizDimLoadings and DimPlot functions. Can you tell from the PCA analysis, the number of cell types that are present?

```
1 library(RColorBrewer)
2 palette(brewer.pal(12, "Paired"))
3
4 DimPlot(pdac1, reduction = "pca", label = TRUE, pt.size = 0.5, group.by =
5 "ident")
```

It is difficult to determine the number of cell types by solely looking at the PCA analysis but it is possible to determine the variations between the cells and the amount of clustering that is present. It is also possible to find outliers within the data, but in order to determine the number of cell types we would need to do clustering analysis on the PCA data that is reduced

Step 10: PCA heatmaps

```
1 # Create heatmaps of the first 10 dimensions
2 DimHeatmap(object = pdac1, dims = 1:10, cells = 1:200, balanced = TRUE)
```

Step 11: Dimensionality

```
1 # Run jackstraw analysis
2 pdac1 <- JackStraw(pdac1, num.replicate = 100, dims = 1:20)
3
4 # Plot results
5 JackStrawPlot(pdac1, dims = 1:20)
6
```

Plot the results for the first 20 dimensions

```
1 # Run JackStraw analysis
2 set.seed(123)
3 pdac1 <- JackStraw(pdac1, num.replicate = 100, dims = 1:20)
4
5 # Plot JackStraw results
6 JackStrawPlot(pdac1, dims = 1:20)
7
```

Use the elbow plot

```
1 # Get the jackstraw p-values and q-values
2 pvals <- JackStraw(pdac1, num.replicate = 100)$pvals
3 qvals <- JackStraw(pdac1, num.replicate = 100)$qvals
4
5 # Create a data frame of the results
6 jack_df <- data.frame(Dimension = 1:20, Pvals = pvals[1:20], Qvals =
7   qvals[1:20])
8
9 # Plot the results
10 ggplot(jack_df, aes(x = Dimension, y = -log10(Qvals))) +
11   geom_line() +
12   geom_point() +
13   xlab("Principal Component") +
14   ylab("-log10(Q-value)") +
15   theme_classic()
```

Step 12: Clustering.

```
1 # Find neighbors and identify clusters
2 pdac1 <- FindNeighbors(pdac1, dims = 1:9)
3 pdac1 <- FindClusters(pdac1, resolution = 0.5)
```

Step 13: Perform a UMAP analysis using the first 9 dimensions using RunUMAP and then visualize it using DimPlot.

How many clusters do you get? How many possible mistakes do you see?

I see about 4 separate clusters and since there are some outliers I would say there are about 2 or 3 separate mistakes that are clearly identifiable

```
1 # Perform UMAP analysis
2 pdac1 <- RunUMAP(pdac1, dims = 1:9)
3
4 # Visualize UMAP projection
5 DimPlot(pdac1, reduction = "umap")
6
```

Step 14: Identify the markers that compare each cluster against all. Report only positively markers. Use the FindAllMarkers for this.

```
1 # Identify markers for each cluster
  pdac1 <- FindAllMarkers(pdac1, only.pos = TRUE, min.pct = 0.25,
2 logfc.threshold = 0.25)
3
```

Step 15: Create a violin plot using one feature from each cluster.

```
1 # Create a violin plot for one feature from each cluster
2 features <- c("CDK1", "G6PD", "RGS2", "DCN", "FGF13", "DCBLD2")
3 VlnPlot(pdac1, features = features, group.by = "ident")
4
```

Step 16: Create a feature plot using the same features as before.

```
1 # Create a feature plot for the same features as before
  FeaturePlot(pdac1, features = features, pt.size = 0.5, reduction.use =
2 "pca", cols = c("grey", "red", "blue", "green", "orange", "purple"))
3
```