# Homework 4

**Priyanka Emani - 1001981861**
**Aravind Kashyap - 1001956591**
**Shreyas Jagadeep - 1001888859**

UNIVERSITY OF
**TEXAS**
ARLINGTON

**CSE: 5370 Bioinformatics**

April 27, 2022

# Collaboration Statement

The assignment was carried out as a group. Following students were part of the group and made individual contributions in completing the assignment.

# 1. Bioinformatics

Single Cell RNA Analysis. The dataset used describes single-cell mouse tissues. The dataset contains 100,000 cells from the brain of a mouse. Each column in the expression matrix csv file corresponds to a a transcript (gene) wherelese each row corresponds to a single cell. The metadata csv file describes each cell.

# 2. Imports and Data Objects

```python
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
import scanpy as sc
import pandas as pd
import numpy as np
%matplotlib inline

count = pd.read_csv('/content/drive/MyDrive/data/brain_counts.csv', index_col=0)
metadata = pd.read_csv('/content/drive/MyDrive/data/brain_metadata.
,csv',index_col=0)
adata = sc.AnnData(X = count, obs = metadata)
```

# 3. Preprocessing

1.0 Spiked genes

```python
# record and keep count of spikes
is_spiked = {}
num_spikes = 0
for gene in adata.var_names:
if 'ERCC' in gene:
is_spiked[gene] = True
num_spikes += 1
else:
is_spiked[gene] = False
adata.var['ERCC'] = pd.Series(is_spiked)
# write
adata.write('/content/drive/MyDrive/data/brain_raw.h5ad')
# read
adata = sc.read('/content/drive/MyDrive/data/brain_raw.h5ad')
```
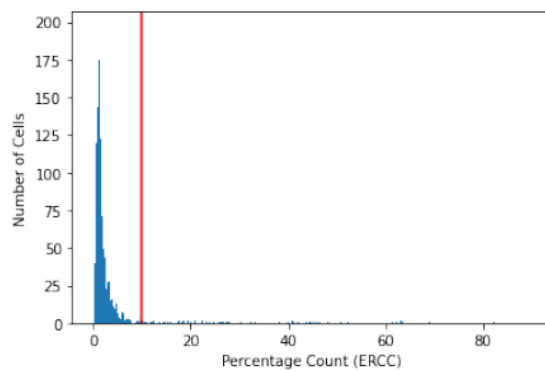
2.0 Quality control

```python
# Compute QC Metrics
qc = sc.pp.calculate_qc_metrics(adata, qc_vars = ['ERCC'])
cell_qc = qc[0]
gene_qc = qc[1]
```
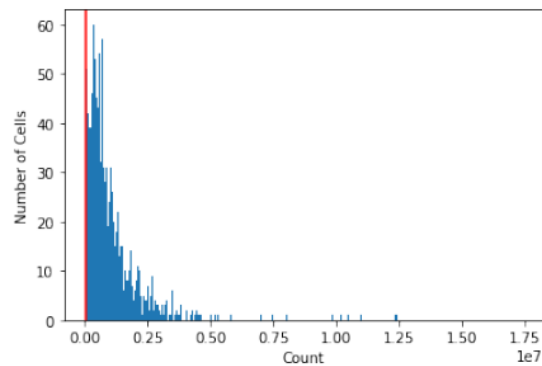
qc - cells

```python


# visualize spike-ins
# outliers to be filtered
plt.hist(cell_qc['pct_counts_ERCC'], bins=1000)
plt.xlabel('Percentage Count (ERCC)')
plt.ylabel('Number of Cells')
plt.axvline(10, color='red')
```
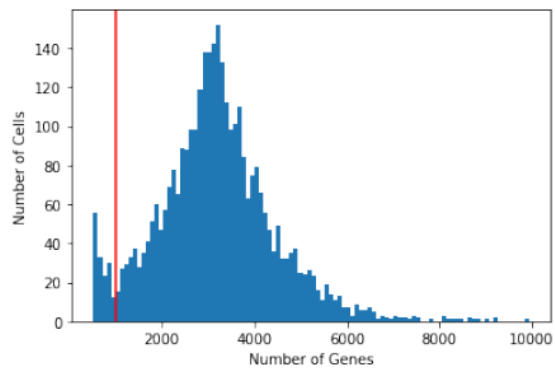
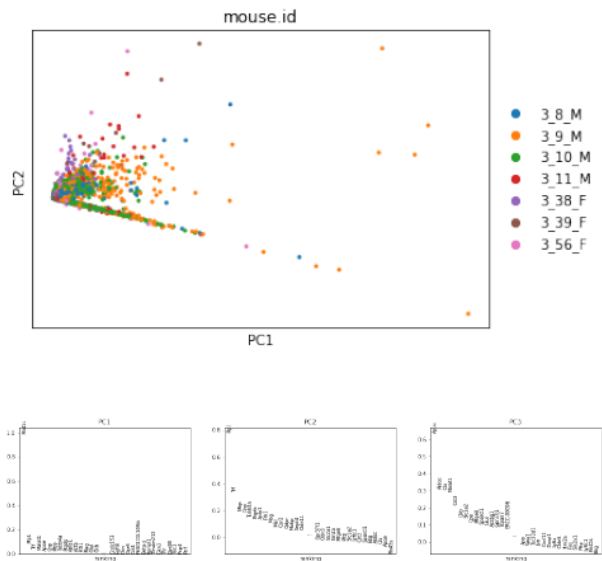<matplotlib.lines.Line2D at 0x7fdffb780c50>



```
1  # visualize lib size (reads per cell)
2  # cell with few reads to be filtered
3  plt.hist(cell_qc['total_counts'], bins=1000)
4  plt.xlabel('Count')
5  plt.ylabel('Number of Cells')
6  plt.axvline(50000, color='red')
```

<matplotlib.lines.Line2D at 0x7fdffad67710>



```
1  # visualize detected gene-count (in cells)
2  # outliers to be filtered
3  plt.hist(cell_qc['n_genes_by_counts'], bins=100)
4  plt.xlabel('Number of Genes')
5  plt.ylabel('Number of Cells')
6  plt.axvline(1000, color='red')
```

<matplotlib.lines.Line2D at 0x7fdffa2f8150>

3

```
1  # filtering
2  # by ERCC
3  low_ERCC = (cell_qc['pct_counts_ERCC'] < 10)
4  adata = adata[low_ERCC]
5  # by Gene count
6  sc.pp.filter_cells(adata, min_genes = 1000)
7  # write
8  adata.write('/content/drive/MyDrive/data/brain_qc.h5ad')
9  # read
10 adata = sc.read('/content/drive/MyDrive/data/brain_qc.h5ad')
```
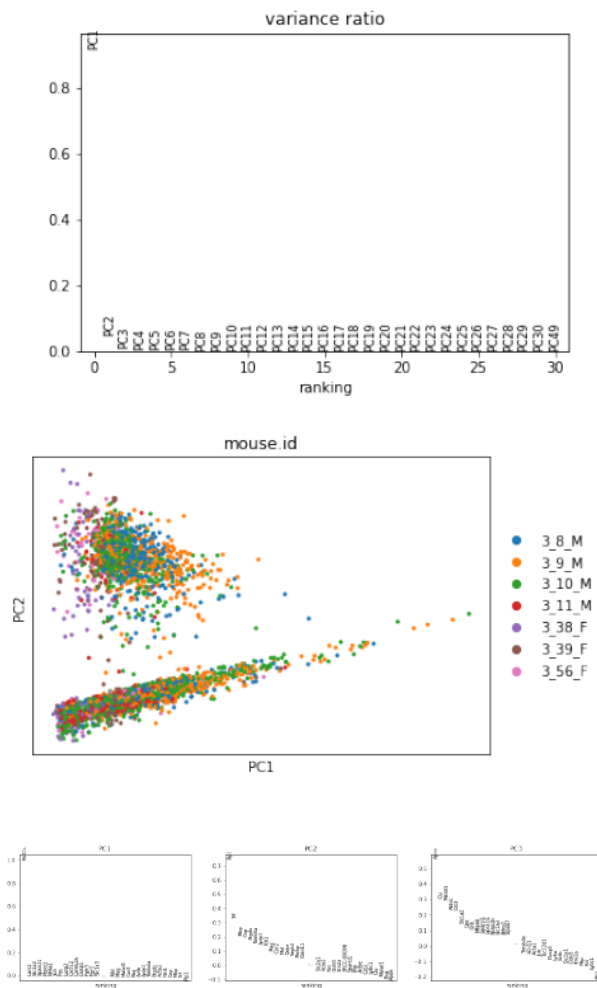
### 3.0 Normalization

Cell library normalization

```
1
2  # visualization of data before Normalization
3  sc.pp.pca(adata)
4  sc.pl.pca_overview(adata, color='mouse.id')
```
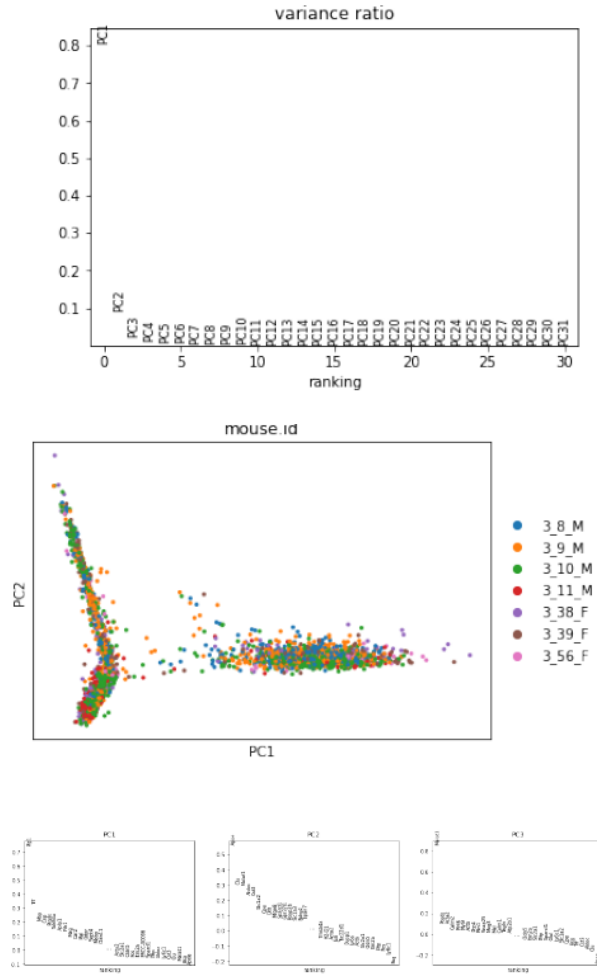
variance ratio



mouse.id

- 3_8_M
- 3_9_M
- 3_10_M
- 3_11_M
- 3_38_F
- 3_39_F
- 3_56_F



```
1    # copy data. copy and original to be used for comparison
2  adata_ = adata.copy()
3  # copy data before normalizing
4  adata_.raw = adata_
5  # CPM normalization
6  sc.pp.normalize_per_cell(adata_,counts_per_cell_after=1e6)
7  # visualization of data after Normalization
8  sc.pp.pca(adata_)
9  sc.pl.pca_overview(adata_, color='mouse.id')
10
```
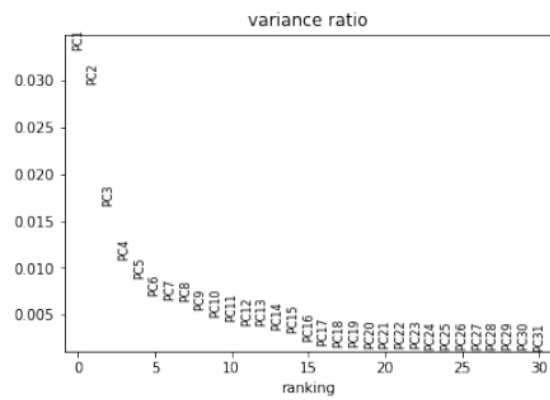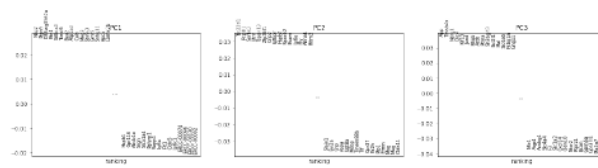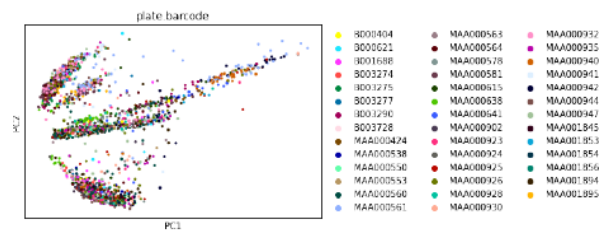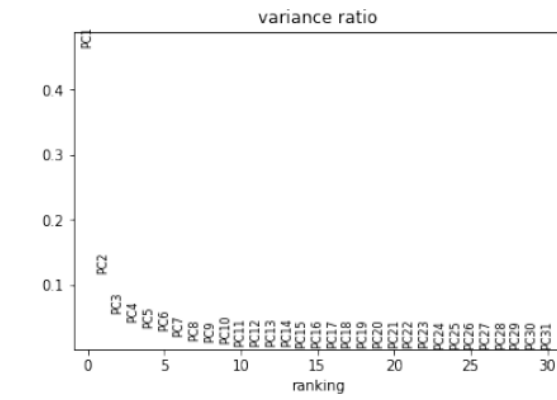
variance ratio



mouse.id

3_8_M
3_9_M
3_10_M
3_11_M
3_38_F
3_39_F
3_56_F



Gene expression normalization

```python
# Rn45s gene is dorminant in data.
# Use normalization to deal with imbalance
Rn45s_0 = adata_.var.index != 'Rn45s'
adata_Rn45s_0 = adata_[:, Rn45s_0]
# visualization of data after Normalization
sc.pp.pca(adata_Rn45s_0)
sc.pl.pca_overview(adata_Rn45s_0, color='mouse.id')

```

```python
        # centering and scaling gene expression data
sc.pp.log1p(adata_)
sc.pp.scale(adata_)
# visualization after Normalization
sc.pp.pca(adata_)
sc.pl.pca_overview(adata_, color='plate.barcode')
```

```python
# write
```

variance ratio



plate barcode



PC1          PC2          PC3

variance ratio

```
2 adata_.write('/content/drive/MyDrive/data/brain_normalized.h5ad')
3 # read
4 adata = sc.read('/content/drive/MyDrive/data/brain_normalized.h5ad')
```

# 4. Analysis

Dimensionality Reduction -UMAP

```
1 sc.pp.neighbors(adata)
2 sc.tl.umap(adata, min_dist=0.5, spread=1.0, n_components=2)
3 # UMAP visualization using Dimensionality Reduction
4 sc.pl.umap(adata, color='cell_ontology_class')
```

Clustering - UMAP

```
1
2 umap_= adata.obsm['X_umap']
3 kmeans = KMeans(n_clusters=4).fit(umap_)
4 adata.obs['K Means Clustering'] = kmeans.labels_
5 adata.obs['K Means Clustering'] = adata.obs['K Means Clustering'].astype(str)
6 # UMAP visualization using KMeans Clsutering
7 sc.pl.umap(adata, color='K Means Clustering')
```

# References

Below are the list of references we used for this project.

1. https://github.com/scverse/scanpy-tutorials/issues/28
2. https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM3109048
3. https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE49712
4. https://scanpy.readthedocs.io/en/stable/generated/scanpy.pp.calculate$_q c_m$etrics.html
5. https://github.com/scverse/scanpy/issues/978
6. https://github.com/scverse/scanpy/issues/1147
7. https://www.oreilly.com/library/view/python-data-science/9781491912126/ch04.html
8. https://matplotlib.org/stable/tutorials/introductory/pyplot.html
9. https://chanzuckerberg.github.io/scRNA-python-workshop/preprocessing/01-basic-qc.html
10. https://github.com/AmoDinho/datacamp-python-data-science-track/blob/master/Machine
11. https://scanpy.readthedocs.io/en/stable/generated/scanpy.pp.filter$_c$ells.html
12. https://scanpy.readthedocs.io/en/stable/generated/scanpy.pl.pca.html
13. https://github.com/scverse/scanpy/issues/324
14. https://nbisweden.github.io/workshop-scRNAseq/labs/compiled/scanpy/scanpy$_0$4$_c$lustering.html
15. https://scanpy.readthedocs.io/en/stable/generated/scanpy.pp.normalize$_p$er$_c$ell.html
16. https://scanpy.readthedocs.io/en/stable/generated/scanpy.pp.normalize$_t$otal.html
17. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4404308/
19. https://scanpy.readthedocs.io/en/stable/generated/scanpy.tl.umap.html
20. https://scanpy.readthedocs.io/en/stable/generated/scanpy.pp.neighbors.html
21. https://github.com/scverse/scanpy/blob/master/scanpy/tools/$_u$map.py
22. https://github.com/theislab/scvelo/issues/37
23. https://towardsdatascience.com/umap-and-k-means-to-classify-characters-league-of-legends-668a788cb3c1
24. https://umap-learn.readthedocs.io/en/latest/clustering.html
25. https://scanpy.readthedocs.io/en/stable/generated/scanpy.pp.log1p.html

Also we have taken help from Dr. Steven Fernandes who is post doctoral research at University of Central Florida