

## § Bioinfo-Lab-3

## § Mouse Retina Dataset

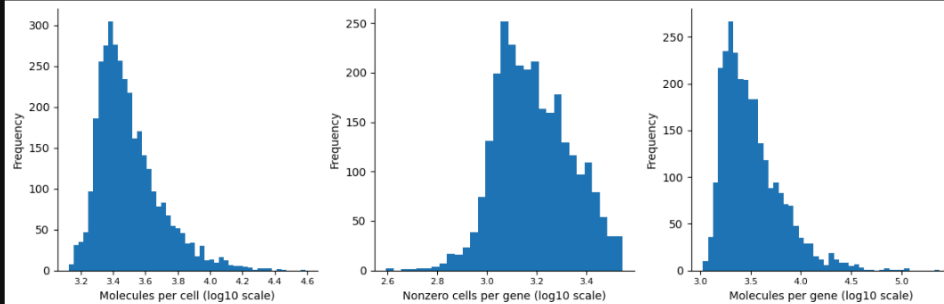
## Task 1: Filter and Normalization

```
[4]: # !pip install palantir
```

```
[5]: import palantir

fig, ax = palantir.plot.plot_molecules_per_cell_and_gene(data.T)

filtered_data = palantir.preprocess.filter_counts_data(data.T, cell_min_molecules=1000, genes_min_cells=10)
```



```
[6]: from sklearn.preprocessing import MinMaxScaler, StandardScaler

scaler_minmax = MinMaxScaler()

data_minmax = pd.DataFrame(scaler_minmax.fit_transform(filtered_data), columns=filtered_data.columns)

print("Min-Max Scaling:")
data_minmax
```

Min-Max Scaling:

```
[6]:
```

|      | CEP290   | CDC59    | PPFIA2   | LIN7A    | PPP1R12A | SYT1     | ZDHHC17 | OSBPL8 | NAP1L1   | KRR1   | ... | TBL1X    | GRIPAP1 | PQBP1    | PCSK1N   | DMD  | MAGED1   | ARHGEF9  | LAS1L    | DYNLT3   | PJA1     |
|------|----------|----------|----------|----------|----------|----------|---------|--------|----------|--------|-----|----------|---------|----------|----------|------|----------|----------|----------|----------|----------|
| 0    | 0.028571 | 0.157895 | 0.200000 | 0.214286 | 0.071429 | 0.723404 | 0.4     | 0.3    | 0.061728 | 0.4375 | ... | 0.583333 | 0.3     | 0.285714 | 0.531646 | 0.02 | 0.402985 | 0.692308 | 0.352941 | 0.882353 | 0.571429 |
| 1    | 0.085714 | 0.052632 | 0.266667 | 0.392857 | 0.357143 | 0.542853 | 0.4     | 0.3    | 0.061728 | 0.1250 | ... | 0.416667 | 0.4     | 0.714286 | 0.202532 | 0.08 | 0.298507 | 0.692308 | 0.176471 | 1.000000 | 0.642857 |
| 2    | 0.000000 | 0.052632 | 0.133333 | 0.214286 | 0.071429 | 0.531915 | 0.0     | 0.1    | 0.037037 | 0.0625 | ... | 0.333333 | 0.1     | 0.142857 | 0.215190 | 0.04 | 0.313433 | 0.153846 | 0.058824 | 0.470588 | 0.785714 |
| 3    | 0.057143 | 0.157895 | 0.200000 | 0.250000 | 0.000000 | 0.563830 | 0.4     | 0.5    | 0.024691 | 0.1250 | ... | 0.250000 | 0.3     | 0.142857 | 0.101266 | 0.00 | 0.149254 | 0.076923 | 0.058824 | 0.647059 | 0.142857 |
| 4    | 0.285714 | 0.157895 | 0.133333 | 0.000000 | 0.071429 | 0.265957 | 0.4     | 0.1    | 0.074074 | 0.0000 | ... | 0.083333 | 1.0     | 0.285714 | 0.025316 | 0.18 | 0.014925 | 0.000000 | 0.117647 | 0.294118 | 0.214286 |
| ...  | ...      | ...      | ...      | ...      | ...      | ...      | ...     | ...    | ...      | ...    | ... | ...      | ...     | ...      | ...      | ...  | ...      | ...      | ...      | ...      | ...      |
| 3512 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0     | 0.2    | 0.000000 | 0.0000 | ... | 0.000000 | 0.0     | 0.000000 | 0.000000 | 0.00 | 0.044776 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 3513 | 0.000000 | 0.052632 | 0.000000 | 0.000000 | 0.000000 | 0.010638 | 0.0     | 0.0    | 0.037037 | 0.0000 | ... | 0.000000 | 0.0     | 0.000000 | 0.000000 | 0.08 | 0.014925 | 0.000000 | 0.117647 | 0.000000 | 0.000000 |
| 3514 | 0.000000 | 0.263158 | 0.000000 | 0.000000 | 0.000000 | 0.021277 | 0.0     | 0.0    | 0.000000 | 0.0000 | ... | 0.000000 | 0.0     | 0.000000 | 0.012658 | 0.00 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| 3515 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.071429 | 0.053191 | 0.0     | 0.0    | 0.000000 | 0.0000 | ... | 0.083333 | 0.1     | 0.000000 | 0.000000 | 0.00 | 0.000000 | 0.153846 | 0.000000 | 0.000000 | 0.071429 |
| 3516 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.071429 | 0.000000 | 0.0     | 0.0    | 0.000000 | 0.0000 | ... | 0.000000 | 0.0     | 0.000000 | 0.000000 | 0.00 | 0.029851 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |

3517 rows x 2773 columns

```
[7]: scaler_zscore = StandardScaler()
data_zscore = pd.DataFrame(scaler_zscore.fit_transform(filtered_data), columns=filtered_data.columns)
print("\nZ-Score Normalization:")
data_zscore
```

Z-Score Normalization:

```
[7]:
```

|      | CEP290    | CDC59     | PPFIA2    | LIN7A     | PPP1R12A  | SYT1      | ZDHHC17   | OSBPL8    | NAP1L1    | KRR1      | ... | TBL1X     | GRIPAP1   | PQBP1     | PCSK1N    | DMD       | MAGED1    | ARHGEF9   | LAS1L     | DYNLT3    | PJA1      |
|------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 0    | -0.382141 | 1.237542  | 1.271030  | 1.015943  | 0.387928  | 6.149090  | 2.761013  | 2.454297  | 0.900040  | 6.980350  | ... | 4.820180  | 1.948172  | 1.371529  | 6.861591  | -0.272221 | 5.927894  | 5.784232  | 4.359605  | 11.830845 | 4.402182  |
| 1    | 0.283530  | -0.085186 | 1.910545  | 2.387668  | 4.172511  | 4.369979  | 2.761013  | 2.454297  | 0.900040  | 1.605089  | ... | 3.255785  | 2.817239  | 4.395771  | 2.355885  | 0.721278  | 4.272339  | 5.784232  | 1.838051  | 13.479709 | 5.039452  |
| 2    | -0.714976 | -0.085186 | 0.631514  | 1.015943  | 0.387928  | 4.265325  | -0.660581 | 0.457128  | 0.246023  | 0.530037  | ... | 2.473587  | 0.210039  | 0.363448  | 2.529181  | 0.058945  | 4.508847  | 0.869412  | 0.157015  | 6.059821  | 6.313992  |
| 3    | -0.049305 | 1.237542  | 1.271030  | 1.290228  | -0.558218 | 4.579286  | 2.761013  | 4.451465  | -0.080985 | 1.605089  | ... | 1.691389  | 1.948172  | 0.363448  | 0.969514  | -0.603388 | 1.907261  | 0.167294  | 0.157015  | 8.533117  | 0.578562  |
| 4    | 2.613376  | 1.237542  | 0.631514  | -0.630128 | 0.387928  | 1.648986  | 2.761013  | 0.457128  | 1.227049  | -0.545015 | ... | 0.126993  | 8.031640  | 1.371529  | -0.070265 | 2.377110  | -0.221310 | -0.534823 | 0.997533  | 3.586525  | 1.215832  |
| ...  | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ... | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       |
| 3512 | -0.714976 | -0.746550 | -0.647516 | -0.630128 | -0.558218 | -0.967354 | -0.660581 | 1.455713  | -0.735002 | -0.545015 | ... | -0.655205 | -0.659028 | -0.644633 | -0.416857 | -0.603388 | 0.251706  | -0.534823 | -0.683503 | -0.535635 | -0.695978 |
| 3513 | -0.714976 | -0.085186 | -0.647516 | -0.630128 | -0.558218 | -0.862700 | -0.660581 | -0.541456 | 0.246023  | -0.545015 | ... | -0.655205 | -0.659028 | -0.644633 | -0.416857 | 0.721278  | -0.221310 | -0.534823 | 0.997533  | -0.535635 | -0.695978 |
| 3514 | -0.714976 | 2.560270  | -0.647516 | -0.630128 | -0.558218 | -0.758047 | -0.660581 | -0.541456 | -0.735002 | -0.545015 | ... | -0.655205 | -0.659028 | -0.644633 | -0.243561 | -0.603388 | -0.457818 | -0.534823 | -0.683503 | -0.535635 | -0.695978 |
| 3515 | -0.714976 | -0.746550 | -0.647516 | -0.630128 | 0.387928  | -0.444086 | -0.660581 | -0.541456 | -0.735002 | -0.545015 | ... | 0.126993  | 0.210039  | -0.644633 | -0.416857 | -0.603388 | -0.457818 | 0.869412  | -0.683503 | -0.535635 | -0.058708 |
| 3516 | -0.714976 | -0.746550 | -0.647516 | -0.630128 | 0.387928  | -0.967354 | -0.660581 | -0.541456 | -0.735002 | -0.545015 | ... | -0.655205 | -0.659028 | -0.644633 | -0.416857 | -0.603388 | 0.015198  | -0.534823 | -0.683503 | -0.535635 | -0.695978 |

3517 rows x 2773 columns

```
[8]: import numpy as np
data_log = np.log(filtered_data+1)
print("Log-transformed Data:")
data_log
```

Log-transformed Data:

```
[8]:
```

|                   | CEP290   | CDC59    | PPFIA2   | LIN7A    | PPP1R12A | SYT1     | ZDHHC17  | OSBPL8   | NAP1L1   | KRR1     | ... | TBL1X    | GRIPAP1  | PQBP1    | PCSK1N   | DMD      | MAGED1   | ARHGEF9  | LAS1L    | DYNLT3   | PJA1     |
|-------------------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|-----|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| r1_GGCCOACGATCCG  | 0.693147 | 1.386294 | 1.386294 | 1.945910 | 0.693147 | 4.234107 | 1.609438 | 1.386294 | 1.791759 | 2.079442 | ... | 2.079442 | 1.386294 | 1.098612 | 3.761200 | 0.693147 | 3.332205 | 2.302585 | 1.945910 | 2.772589 | 2.197225 |
| r1_CTTGTGCGGGGAA  | 1.386294 | 0.693147 | 1.609438 | 2.484907 | 1.791759 | 3.951244 | 1.609438 | 1.386294 | 1.791759 | 1.098612 | ... | 1.791759 | 1.609438 | 1.791759 | 2.833213 | 1.609438 | 3.044522 | 2.302585 | 1.386294 | 2.890372 | 2.302585 |
| r1_GGCCAACGCTCTC  | 0.000000 | 0.693147 | 1.098612 | 1.945910 | 0.693147 | 3.931826 | 0.000000 | 0.693147 | 1.386294 | 0.693147 | ... | 1.609438 | 0.693147 | 0.693147 | 2.890372 | 1.098612 | 3.091042 | 1.098612 | 0.693147 | 2.197225 | 2.484907 |
| r1_GAATGTGGAGGACA | 1.098612 | 1.386294 | 1.386294 | 2.079442 | 0.000000 | 3.988984 | 1.609438 | 1.791759 | 1.098612 | 1.098612 | ... | 1.386294 | 1.386294 | 0.693147 | 2.197225 | 0.000000 | 2.397895 | 0.693147 | 0.693147 | 2.484907 | 1.098612 |
| r1_GTGCCGCTCTCTC  | 2.397895 | 1.386294 | 1.098612 | 0.000000 | 0.693147 | 3.258097 | 1.609438 | 0.693147 | 1.945910 | 0.000000 | ... | 0.693147 | 2.397895 | 1.098612 | 1.098612 | 2.302585 | 0.693147 | 0.000000 | 1.098612 | 1.791759 | 1.386294 |
| ...               | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ... | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      |
| p1_GTCTGATGCAAC   | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 1.098612 | 0.000000 | ... | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 1.386294 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |
| p1_AGTGCGACGCTCG  | 0.000000 | 0.693147 | 0.000000 | 0.000000 | 0.000000 | 0.693147 | 0.000000 | 0.000000 | 1.386294 | 0.000000 | ... | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 1.609438 | 0.693147 | 0.000000 | 1.098612 | 0.000000 | 0.000000 |
| p1_CCAATAGCCGAT   | 0.000000 | 1.791759 | 0.000000 | 0.000000 | 0.000000 | 1.098612 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | ... | 0.000000 | 0.000000 | 0.000000 | 0.693147 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |          |
| p1_TTAGTCTCTACCG  | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.693147 | 1.791759 | 0.000000 | 0.000000 | 0.000000 | ... | 0.693147 | 0.693147 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 1.098612 | 0.000000 | 0.000000 | 0.693147 |
| p1_GCATAAAGTATA   | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.693147 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | ... | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 1.098612 | 0.000000 | 0.000000 | 0.000000 | 0.000000 |

3517 rows x 2773 columns

## Task 2: MAGIC

```
[22]: # !pip install magic-impute
```

```
[10]: import magic
magic_operator=magic.MAGIC()
data_magic=magic_operator.fit_transform(data_log, genes='all_genes')

Calculating MAGIC...
Running MAGIC on 3517 cells and 2773 genes.
Calculating graph and diffusion operator...
Calculating PCA...
Calculated PCA in 0.54 seconds.
Calculating KNN search...
Calculated KNN search in 1.23 seconds.
Calculating affinities...
Calculated affinities in 1.14 seconds.
Calculated graph and diffusion operator in 2.94 seconds.
Calculating imputation...
Calculated imputation in 0.48 seconds.
Calculated MAGIC in 3.45 seconds.
```

```
[11]: data_magic
```

```
[11]:
```

|                  | CEP290   | CDC59    | PPFIA2   | LINTA    | PPP1R12A | SYT1     | ZDHHC17  | OSBPL8   | NAP1L1   | KRR1     | ... | TBL1X    | GRIPAP1  | PQBP1    | PCSKIN   | DMD      | MAGED1   | ARHGEF9  | LAS1L    | DYNLT3   | PJA1     |
|------------------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|-----|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| rl_GCCGCCAGTCCG  | 0.862227 | 1.083810 | 1.120356 | 2.028935 | 1.076615 | 3.900440 | 1.170198 | 0.980097 | 2.044446 | 1.054837 | ... | 1.724078 | 1.266112 | 1.446556 | 2.884915 | 0.683855 | 2.606584 | 1.556724 | 1.466776 | 1.977347 | 2.026418 |
| rl_CTTGTGCGGAGAA | 0.861889 | 1.080732 | 1.109982 | 2.023916 | 1.064033 | 3.894473 | 1.163940 | 0.975179 | 2.043852 | 1.043431 | ... | 1.710471 | 1.261051 | 1.442919 | 2.876661 | 0.689476 | 2.596321 | 1.540784 | 1.459034 | 1.964914 | 2.018262 |
| rl_GCGCAACTCCTC  | 0.844133 | 1.069069 | 1.068900 | 1.992562 | 1.020416 | 3.874745 | 1.136689 | 0.942444 | 2.040263 | 1.013713 | ... | 1.666048 | 1.236924 | 1.412667 | 2.857988 | 0.613030 | 2.572097 | 1.486537 | 1.424778 | 1.919669 | 2.001678 |
| rl_GATTGGGAGGCA  | 0.834013 | 1.000381 | 0.977214 | 1.901137 | 0.796713 | 3.686407 | 0.983522 | 0.867776 | 1.874945 | 0.794658 | ... | 1.376634 | 1.133459 | 1.165508 | 2.654247 | 0.422974 | 2.276723 | 1.361222 | 1.266365 | 1.611574 | 1.786427 |
| rl_GTGCCGCTCTCT  | 2.481328 | 1.651764 | 1.248324 | 0.475312 | 1.047405 | 2.948844 | 1.190809 | 1.141672 | 2.177910 | 0.964871 | ... | 1.008528 | 1.471604 | 1.023622 | 1.261460 | 2.576503 | 1.053146 | 0.272308 | 1.170301 | 0.933983 | 0.751174 |
| ...              | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ... | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      |
| p1_OTCTGATGCAAC  | 0.074809 | 0.250276 | 0.069107 | 0.263597 | 0.180710 | 0.127756 | 0.093311 | 0.183249 | 0.156478 | 0.099712 | ... | 0.395145 | 0.079952 | 0.150062 | 0.211371 | 0.221174 | 0.807290 | 0.112538 | 0.101429 | 0.396065 | 0.302892 |
| p1_AGTGCCAGCTCG  | 0.762621 | 0.420609 | 0.184430 | 0.033470 | 0.169211 | 1.284931 | 0.200687 | 0.238599 | 0.722147 | 0.182929 | ... | 0.183542 | 0.235960 | 0.152606 | 0.356440 | 0.834892 | 0.265389 | 0.048377 | 0.217577 | 0.196417 | 0.115361 |
| p1_CCAATAGCCGAT  | 0.743084 | 0.501508 | 0.146600 | 0.047412 | 0.153437 | 1.170968 | 0.148179 | 0.256126 | 0.727583 | 0.179621 | ... | 0.225725 | 0.234132 | 0.169391 | 0.394956 | 0.829194 | 0.235355 | 0.039011 | 0.202259 | 0.193969 | 0.119419 |
| p1_TAGTCTCTACCG  | 0.503462 | 0.339239 | 0.388419 | 0.629527 | 0.242635 | 2.016967 | 0.232638 | 0.230735 | 0.653998 | 0.207670 | ... | 0.382004 | 0.333536 | 0.339399 | 0.503003 | 0.406402 | 0.480049 | 0.343619 | 0.391367 | 0.175018 | 0.465066 |
| p1_GCAATAAGATATA | 0.096076 | 0.235807 | 0.057068 | 0.243076 | 0.173100 | 0.147513 | 0.099769 | 0.162345 | 0.162976 | 0.084243 | ... | 0.418160 | 0.087797 | 0.168241 | 0.197920 | 0.242122 | 0.764070 | 0.131225 | 0.121671 | 0.419430 | 0.326460 |

3517 rows x 2773 columns

```
[12]: # Task 3: t-SNE
```

```
[13]: from sklearn.manifold import TSNE
import matplotlib.pyplot as plt
```

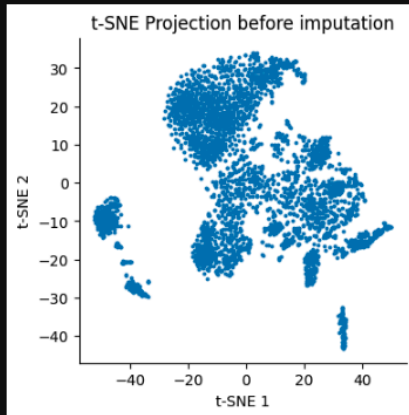
```
[14]: tsne = TSNE(n_components=2)
data_tsne=tsne.fit_transform(data_log)
data_magic_tsne = tsne.fit_transform(data_magic)
```

## Task 3: t-SNE

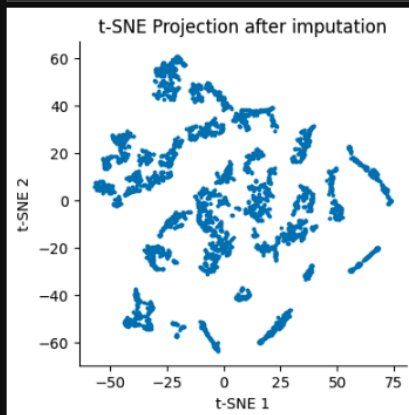
```
[15]: from sklearn.manifold import TSNE
import matplotlib.pyplot as plt

tsne = TSNE(n_components=2)
data_tsne=tsne.fit_transform(data_log)
data_magic_tsne = tsne.fit_transform(data_magic)
```

```
[16]: plt.scatter(data_tsne[:,0],data_tsne[:,1],s=3)
plt.title('t-SNE Projection before imputation')
plt.xlabel('t-SNE 1')
plt.ylabel('t-SNE 2')
plt.show()
```



```
[17]: plt.scatter(data_magic_tsne[:,0],data_magic_tsne[:,1],s=3)
plt.title('t-SNE Projection after imputation')
plt.xlabel('t-SNE 1')
plt.ylabel('t-SNE 2')
plt.show()
```



## Task 4: Dimension

```
[18]: from sklearn.decomposition import PCA

pca = PCA(n_components=2)
data_pca = pca.fit_transform(data_log)

import umap
reducer = umap.UMAP(n_neighbors=5)

data_umap = reducer.fit_transform(data_log)
```

## Task 5: Clustering

```
[19]: from sklearn.cluster import KMeans

kmeans = KMeans(n_clusters=10)
kmeans_labels = kmeans.fit_predict(data_umap)

from sklearn.cluster import AgglomerativeClustering

model = AgglomerativeClustering(n_clusters=10, metric='euclidean', linkage='ward')
hierarchical_labels = model.fit_predict(data_umap)
```

## Task 6: ARI

```
[20]: from sklearn.metrics import adjusted_rand_score
import pandas as pd

ground_truth = pd.read_csv("data/data1_mouse_retina/sample_cluster_ref_filtered.txt", sep='\t', header=None, index_col=0)

ground_truth.index = ground_truth.index.astype(str)
filtered_data.index = filtered_data.index.astype(str)

filter_truth = ground_truth[ground_truth.index.isin(filtered_data.index)]

truth_labels = filter_truth.iloc[:, 0].values.flatten()

ari_kmeans = adjusted_rand_score(truth_labels, kmeans_labels)
ari_hierarchical = adjusted_rand_score(truth_labels, hierarchical_labels)

ari_kmeans, ari_hierarchical
```

```
[20]: (0.3574350970964889, 0.3957469246365283)
```

## § Mesc Dataset

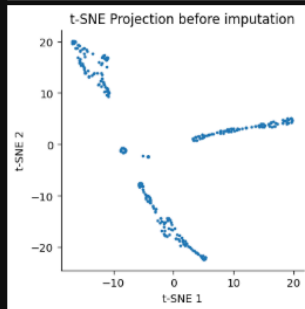


## Task 3: t-SNE

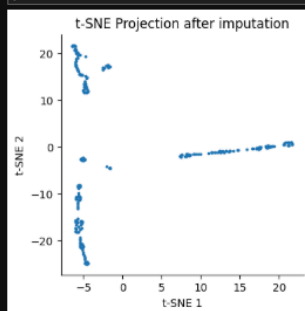
```
[15]: from sklearn.manifold import TSNE
import matplotlib.pyplot as plt

tsne = TSNE(n_components=2)
data_tsne = tsne.fit_transform(data_log)
data_magic_tsne = tsne.fit_transform(data_magic)

[16]: plt.scatter(data_tsne[:,0], data_tsne[:,1], s=3)
plt.title('t-SNE Projection before imputation')
plt.xlabel('t-SNE 1')
plt.ylabel('t-SNE 2')
plt.show()
```



```
[17]: plt.scatter(data_magic_tsne[:,0], data_magic_tsne[:,1], s=3)
plt.title('t-SNE Projection after imputation')
plt.xlabel('t-SNE 1')
plt.ylabel('t-SNE 2')
plt.show()
```



## Task 4: Dimension

```
[18]: from sklearn.decomposition import PCA

pca = PCA(n_components=2)
data_pca = pca.fit_transform(data_log)

import umap
reducer = umap.UMAP(n_neighbors=5)
data_umap = reducer.fit_transform(data_log)

/opt/conda/lib/python3.11/site-packages/sklearn/manifold/_spectral_embedding.py:329: UserWarning: Graph is not fully connected, spectral embedding may not work as expected.
warnings.warn()
```

## Task 4: Dimension

```
[18]: from sklearn.decomposition import PCA

pca = PCA(n_components=2)
data_pca = pca.fit_transform(data_log)

import umap
reducer = umap.UMAP(n_neighbors=5)
data_umap = reducer.fit_transform(data_log)

/opt/conda/lib/python3.11/site-packages/sklearn/manifold/_spectral_embedding.py:329: UserWarning: Graph is not fully connected, spectral embedding may not work as expected.
warnings.warn()
```

## Task 5: Clustering

```
[19]: from sklearn.cluster import KMeans

kmeans = KMeans(n_clusters=10)
kmeans_labels = kmeans.fit_predict(data_umap)

from sklearn.cluster import AgglomerativeClustering

model = AgglomerativeClustering(n_clusters=10, metric='euclidean', linkage='ward')
hierarchical_labels = model.fit_predict(data_umap)
```

## Task 6: ARI

```
[20]: from sklearn.metrics import adjusted_rand_score
import pandas as pd

ground_truth = pd.read_csv("data/data2_mesc/sample_cluster_ref_filtered.txt", sep='\t', header=None, index_col=0)

ground_truth.index = ground_truth.index.astype(int)
filtered_data.index = filtered_data.index.astype(int)

filter_truth = ground_truth[ground_truth.index.isin(filtered_data.index)]

truth_labels = filter_truth.iloc[:, 0].values.flatten()

ari_kmeans = adjusted_rand_score(truth_labels, kmeans_labels)
ari_hierarchical = adjusted_rand_score(truth_labels, hierarchical_labels)

ari_kmeans, ari_hierarchical
```

```
[20]: (0.24332549678853352, 0.3809128475299934)
```

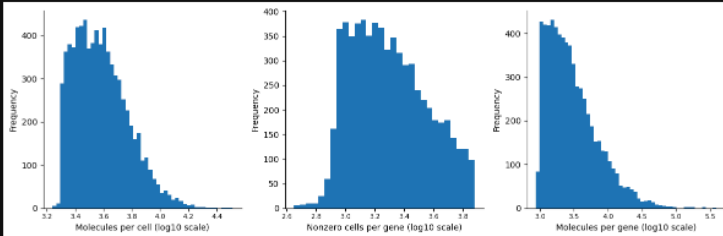
## Embryo cortex Dataset

## Task 1: Filter and Normalization

```
[4]: !pip install palantir
```

```
[5]: import palantir
```

```
fig, ax = palantir.plot.plot_molecules_per_cell_and_gene(data.T)
filtered_data = palantir.preprocess.filter_counts_data(data.T, cell_min_molecules=1000, genes_min_cells=10)
```



```
[6]: from sklearn.preprocessing import MinMaxScaler, StandardScaler
```

```
scaler_minmax = MinMaxScaler()
data_minmax = pd.DataFrame(scaler_minmax.fit_transform(filtered_data), columns=filtered_data.columns)
print("Min-Max Scaling:")
data_minmax
```

```
Min-Max Scaling:
```

|      | 0610007N19RK | 0610007P14RK | 0610008B22RK | 0610008D07RK | 0610009C020RK | 0610010F05RK | 0610011F06RK | 0610030E20RK | 0610037L13RK | 1110001A16RK | ... | l7Rn6 | mt-Ct1   | mt-Cytb  | mt-Nd1   | mt-Nd2   | mt-Nd4   | mt-Nd5   | mt-Nd6   | mt-Rnr1  | mt-Rnr2  |
|------|--------------|--------------|--------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|-----|-------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| 0    | 0.0          | 0.166667     | 0.0          | 0.000000     | 0.0           | 0.000000     | 0.000        | 0.0          | 0.000000     | 0.000000     | ... | 0.125 | 0.024390 | 0.072464 | 0.075949 | 0.000000 | 0.000000 | 0.014286 | 0.000000 | 0.000000 | 0.030973 |
| 1    | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.0           | 0.000000     | 0.000        | 0.0          | 0.000000     | 0.000000     | ... | 0.000 | 0.000000 | 0.004203 | 0.088608 | 0.036364 | 0.035714 | 0.000000 | 0.000000 | 0.000000 | 0.026549 |
| 2    | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.0           | 0.066667     | 0.000        | 0.0          | 0.000000     | 0.000000     | ... | 0.125 | 0.000000 | 0.069657 | 0.113924 | 0.072727 | 0.125000 | 0.100000 | 0.142857 | 0.063333 | 0.053097 |
| 3    | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.0           | 0.066667     | 0.000        | 0.0          | 0.000000     | 0.000000     | ... | 0.000 | 0.024390 | 0.130435 | 0.101266 | 0.090909 | 0.089286 | 0.014286 | 0.000000 | 0.000000 | 0.048673 |
| 4    | 0.0          | 0.000000     | 0.0          | 0.153846     | 0.0           | 0.000000     | 0.125        | 0.0          | 0.000000     | 0.000000     | ... | 0.000 | 0.024390 | 0.043478 | 0.063291 | 0.054545 | 0.053571 | 0.000000 | 0.142857 | 0.000000 | 0.044248 |
| ...  | ...          | ...          | ...          | ...          | ...           | ...          | ...          | ...          | ...          | ...          | ... | ...   | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      |
| 7596 | 0.0          | 0.000000     | 0.0          | 0.461538     | 0.0           | 0.066667     | 0.000        | 0.0          | 0.000000     | 0.285714     | ... | 0.000 | 0.048780 | 0.188406 | 0.253185 | 0.127273 | 0.107143 | 0.057143 | 0.142857 | 0.250000 | 0.132743 |
| 7597 | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.5           | 0.000000     | 0.000        | 0.0          | 0.000000     | 0.142857     | ... | 0.000 | 0.000000 | 0.007246 | 0.025316 | 0.018182 | 0.014286 | 0.000000 | 0.063333 | 0.181416 |          |
| 7598 | 0.0          | 0.000000     | 0.0          | 0.076923     | 0.0           | 0.066667     | 0.000        | 0.0          | 0.000000     | 0.142857     | ... | 0.125 | 0.024390 | 0.028986 | 0.063291 | 0.000000 | 0.000000 | 0.014286 | 0.142857 | 0.041667 | 0.088496 |
| 7599 | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.0           | 0.000000     | 0.000        | 0.0          | 0.000000     | 0.000000     | ... | 0.000 | 0.024390 | 0.072464 | 0.101266 | 0.018182 | 0.035714 | 0.028571 | 0.000000 | 0.041667 | 0.188142 |
| 7600 | 0.0          | 0.000000     | 0.0          | 0.153846     | 0.0           | 0.000000     | 0.000        | 0.0          | 0.142857     | 0.000000     | ... | 0.125 | 0.073171 | 0.043478 | 0.075949 | 0.036364 | 0.035714 | 0.042857 | 0.000000 | 0.000000 | 0.119469 |

```
7601 rows x 6470 columns
```

```
[7]: scaler_score = StandardScaler()
data_score = pd.DataFrame(scaler_score.fit_transform(filtered_data), columns=filtered_data.columns)
print("Z-Score Normalization:")
data_score
```

```
Z-Score Normalization:
```

|      | 0610007N19RK | 0610007P14RK | 0610008B22RK | 0610008D07RK | 0610009C020RK | 0610010F05RK | 0610011F06RK | 0610030E20RK | 0610037L13RK | 1110001A16RK | ... | l7Rn6     | mt-Ct1    | mt-Cytb   | mt-Nd1    | mt-Nd2    | mt-Nd4    | mt-Nd5    | mt-Nd6    | mt-Rnr1   | mt-Rnr2   |
|------|--------------|--------------|--------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|-----|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 0    | -0.282364    | 0.968239     | -0.535158    | -1.032482    | -0.355989     | -0.663598    | -0.503545    | -0.371784    | -0.583722    | -0.537763    | ... | 0.405438  | -0.626387 | -0.756392 | -0.696406 | -1.182838 | -1.37843  | -0.741729 | -0.449034 | -0.713510 | -0.919172 |
| 1    | -0.282364    | -0.527295    | -0.535158    | -1.032482    | -0.355989     | -0.663598    | -0.503545    | -0.371784    | -0.583722    | -0.537763    | ... | -0.703499 | -0.193408 | -0.507933 | -0.570999 | -0.649768 | -0.883994 | -1.022178 | -0.449034 | -0.713510 | -0.986133 |
| 2    | -0.282364    | -0.527295    | -0.535158    | -1.032482    | -0.355989     | 0.176650     | -0.503545    | -0.371784    | -0.583722    | -0.537763    | ... | 0.405438  | -0.193408 | -0.507933 | -0.320007 | -0.196999 | 0.200629  | 0.940969  | 1.27460   | 0.156562  | -0.584368 |
| 3    | -0.282364    | -0.527295    | -0.535158    | -1.032482    | -0.355989     | 0.176650     | -0.503545    | -0.371784    | -0.583722    | -0.537763    | ... | -0.703499 | -0.626387 | -0.093835 | -0.445503 | 0.149836  | -0.233220 | -0.741729 | -0.449034 | -0.713510 | -0.651329 |
| 4    | -0.282364    | -0.527295    | -0.535158    | 0.237918     | -0.355989     | -0.663598    | 1.202027     | -0.371784    | -0.583722    | -0.537763    | ... | -0.703499 | -0.626387 | -1.087671 | -0.821992 | -0.383234 | -0.667069 | -1.022178 | 1.27460   | -0.713510 | -0.718290 |
| ...  | ...          | ...          | ...          | ...          | ...           | ...          | ...          | ...          | ...          | ...          | ... | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       | ...       |
| 7596 | -0.282364    | -0.527295    | -0.535158    | 2.778719     | -0.355989     | 0.176650     | -0.503545    | -0.371784    | -0.583722    | 2.181170     | ... | -0.703499 | -0.218366 | 0.568722  | 1.060454  | 0.682905  | -0.016296 | 0.099620  | 1.27460   | 2.976707  | 0.820927  |
| 7597 | -0.282364    | -0.527295    | -0.535158    | -1.032482    | 4.528259      | -0.663598    | -0.503545    | -0.371784    | -0.583722    | 0.811704     | ... | -0.703499 | -0.193408 | -0.507933 | -1.198481 | -0.963031 | -1.37843  | -0.741729 | -0.449034 | 0.156562  | 1.357497  |
| 7598 | -0.282364    | -0.527295    | -0.535158    | -0.307282    | -0.355989     | 0.176650     | -0.503545    | -0.371784    | -0.583722    | 0.811704     | ... | 0.405438  | -0.626387 | -1.253310 | -0.821992 | -1.182838 | -1.37843  | -0.741729 | 1.27460   | -0.098474 | -0.046861 |
| 7599 | -0.282364    | -0.527295    | -0.535158    | -1.032482    | -0.355989     | -0.663598    | -0.503545    | -0.371784    | -0.583722    | -0.537763    | ... | -0.703499 | -0.626387 | -0.756392 | -0.445503 | -0.963031 | -0.883994 | -0.461729 | -0.449034 | -0.098474 | 1.566614  |
| 7600 | -0.282364    | -0.527295    | -0.535158    | 0.237918     | -0.355989     | -0.663598    | -0.503545    | -0.371784    | 0.698241     | -0.537763    | ... | 0.405438  | 0.193656  | -1.087671 | -0.696406 | -0.649768 | -0.883994 | -0.180629 | -0.449034 | -0.713510 | 0.420045  |

```
7601 rows x 6470 columns
```

```
[8]: import numpy as np
```

```
data_log = np.log(filtered_data)
```

```
print("Log-transformed Data:")
data_log
```

```
Log-transformed Data:
```

|                         | 0610007N19RK | 0610007P14RK | 0610008B22RK | 0610008D07RK | 0610009C020RK | 0610010F05RK | 0610011F06RK | 0610030E20RK | 0610037L13RK | 1110001A16RK | ... | l7Rn6    | mt-Ct1   | mt-Cytb  | mt-Nd1   | mt-Nd2   | mt-Nd4   | mt-Nd5   | mt-Nd6   | mt-Rnr1  | mt-Rnr2  |          |
|-------------------------|--------------|--------------|--------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|-----|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| #t4.WT10_AAAAATCTCTCC   | 0.0          | 0.693147     | 0.0          | 0.000000     | 0.000000      | 0.000000     | 0.000000     | 0.0          | 0.000000     | 0.000000     | ... | 0.693147 | 0.693147 | 2.397895 | 1.945910 | 0.000000 | 0.000000 | 0.693147 | 0.000000 | 0.000000 | 0.000000 | 2.079442 |
| #t4.WT10_AAAACACATCTCC  | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.000000      | 0.000000     | 0.000000     | 0.0          | 0.000000     | 0.000000     | ... | 0.000000 | 0.000000 | 2.639057 | 2.079442 | 1.098612 | 1.098612 | 0.000000 | 0.000000 | 0.000000 | 1.945910 |          |
| #t4.WT10_AAAACGCTGGAT   | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.000000      | 0.693147     | 0.000000     | 0.0          | 0.000000     | 0.000000     | ... | 0.693147 | 0.000000 | 2.584949 | 2.305585 | 1.609438 | 2.079442 | 2.079442 | 0.693147 | 1.098612 | 2.584949 |          |
| #t4.WT10_AAAACGCGCCGA   | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.000000      | 0.000000     | 0.000000     | 0.0          | 0.000000     | 0.000000     | ... | 0.000000 | 0.693147 | 2.944439 | 2.197225 | 1.917559 | 1.917559 | 0.693147 | 0.000000 | 0.000000 | 2.484907 |          |
| #t4.WT10_AAAATGACCTCA   | 0.0          | 0.000000     | 0.0          | 1.098612     | 0.000000      | 0.000000     | 0.693147     | 0.0          | 0.000000     | 0.000000     | ... | 0.000000 | 0.693147 | 1.945910 | 1.917559 | 1.386294 | 1.386294 | 0.000000 | 0.693147 | 0.000000 | 2.397895 |          |
| ...                     | ...          | ...          | ...          | ...          | ...           | ...          | ...          | ...          | ...          | ...          | ... | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      | ...      |          |
| #t4.WT10_2_TTTTGTCCGCG  | 0.0          | 0.000000     | 0.0          | 1.945910     | 0.000000      | 0.693147     | 0.000000     | 0.0          | 0.000000     | 1.098612     | ... | 0.000000 | 1.098612 | 3.295837 | 3.044522 | 2.079442 | 1.945910 | 1.609438 | 0.693147 | 1.945910 | 3.433987 |          |
| #t4.WT10_2_TTTTGTCACTCT | 0.0          | 0.000000     | 0.0          | 0.000000     | 1.098612      | 0.000000     | 0.000000     | 0.0          | 0.000000     | 0.693147     | ... | 0.000000 | 0.693147 | 0.000000 | 0.693147 | 1.098612 | 0.693147 | 0.000000 | 0.693147 | 0.000000 | 1.098612 |          |
| #t4.WT10_2_TTTTATGATGTN | 0.0          | 0.000000     | 0.0          | 0.693147     | 0.000000      | 0.693147     | 0.000000     | 0.0          | 0.000000     | 0.693147     | ... | 0.693147 | 0.693147 | 1.609438 | 1.917559 | 0.000000 | 0.000000 | 0.693147 | 0.693147 | 0.693147 | 3.044522 |          |
| #t4.WT10_2_TTTTATCATCGG | 0.0          | 0.000000     | 0.0          | 0.000000     | 0.000000      | 0.000000     | 0.000000     | 0.0          | 0.000000     | 0.000000     | ... | 0.000000 | 0.693147 | 2.397895 | 2.197225 | 0.693147 | 1.098612 | 1.098612 | 0.000000 | 0.693147 | 3.663562 |          |
| #t4.WT10_2_TTTTATCACTTG | 0.0          | 0.000000     | 0.0          | 1.098612     | 0.000000      | 0.000000     | 0.000000     | 0.0          | 0.693147     | 0.000000     | ... | 0.693147 | 1.386294 | 1.945910 | 1.945910 | 1.098612 | 1.098612 | 1.386294 | 0.000000 | 0.000000 | 3.332205 |          |

```
7601 rows x 6470 columns
```

## Task 2: MAGIC

```
[9]: !pip install magic-impute
```

```
[10]: import magic
```

```
magic_operator = magic.MAGIC()
```

```
data_magic = magic_operator.fit_transform(data_log, genes='all_genes')
```

```
Calculating MAGIC...
```

```
Running MAGIC on 7601 cells and 6470 genes.
```

```
Calculating graph and diffusion operator...
```

```
Calculating PCA...
```

```
Calculated PCA in 1.25 seconds.
```

```
Calculating NN search...
```

```
Calculated NN search in 5.99 seconds.
```

```
Calculating affinities...
```

```
Calculated affinities in 5.38 seconds.
```

```
Calculating graph and diffusion operator in 12.66 seconds.
```

```
Running MAGIC with 'solver=exact' on 6470-dimensional data may take a long time. Consider downsampling specific genes with 'genes=list-like' or using 'solver=approximate'.
```

```
Calculating imputation...
```

```
Calculated imputation in 2.97 seconds.
```

```
Calculated MAGIC in 35.73 seconds.
```

```
[11]: data_magic
```

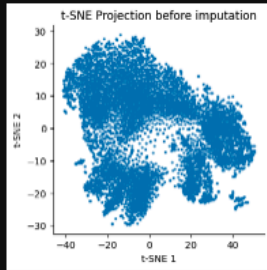
|  | 0610007N19RK | 0610007P14RK | 0610008B22RK | 0610008D07RK | 0610009C020RK | 0610010F05RK | 0610011F06RK | 0610030E20RK | 0610037L13RK | 1110001A16RK | ... | l7Rn6 | mt-Ct1 | mt-Cytb | mt-Nd1 | mt-Nd2 | mt-Nd4 |
|--|--------------|--------------|--------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|-----|-------|--------|---------|--------|--------|--------|
|--|--------------|--------------|--------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|-----|-------|--------|---------|--------|--------|--------|

## Task 3: t-SNE

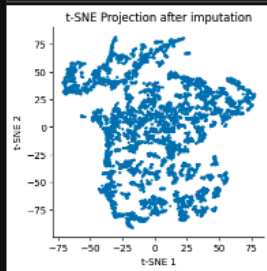
```
[15]: from sklearn.manifold import TSNE
import matplotlib.pyplot as plt

tsne = TSNE(n_components=2)
data_tsne = tsne.fit_transform(data_log)
data_magic_tsne = tsne.fit_transform(data_magic)
```

```
[16]: plt.scatter(data_tsne[:,0],data_tsne[:,1],s=1)
plt.title('t-SNE Projection before imputation')
plt.xlabel('t-SNE 1')
plt.ylabel('t-SNE 2')
plt.show()
```



```
[17]: plt.scatter(data_magic_tsne[:,0],data_magic_tsne[:,1],s=1)
plt.title('t-SNE Projection after imputation')
plt.xlabel('t-SNE 1')
plt.ylabel('t-SNE 2')
plt.show()
```



## Task 4: Dimension

```
[18]: from sklearn.decomposition import PCA

pca = PCA(n_components=2)
data_pca = pca.fit_transform(data_log)

import umap
reducer = umap.UMAP(n_neighbors=5)

data_umap = reducer.fit_transform(data_log)
```

## Task 5: Clustering

```
[19]: from sklearn.cluster import KMeans

kmeans = KMeans(n_clusters=10)
kmeans_labels = kmeans.fit_predict(data_umap)

from sklearn.cluster import AgglomerativeClustering

model = AgglomerativeClustering(n_clusters=10, metric='euclidean', linkage='ward')
hierarchical_labels = model.fit_predict(data_umap)
```

## Task 6: ARI

```
•[25]: from sklearn.metrics import adjusted_rand_score
import pandas as pd

ground_truth = pd.read_csv("data/data3_embryo_cortex/sample_cluster_ref_filtered.txt", sep=' ', header=None, index_col=0)

ground_truth.index = ground_truth.index.astype(str)
filtered_data.index = filtered_data.index.astype(str)
print(ground_truth.head())
filter_truth = ground_truth[ground_truth.index.isin(filtered_data.index)]

truth_labels = filter_truth.iloc[:, 0].values.flatten()

ari_kmeans = adjusted_rand_score(truth_labels, kmeans_labels)
ari_hierarchical = adjusted_rand_score(truth_labels, hierarchical_labels)

ari_kmeans, ari_hierarchical

ground_truth 的形状: (7601, 1)
ground_truth 的前几行:

```

|                       | 1             |
|-----------------------|---------------|
| 0                     |               |
| e14.WT10_AAAAATCTCTCC | RG1           |
| e14.WT10_AAAACACATTCC | LayerV-VI     |
| e14.WT10_AAAACCGTGGAT | LayerV-VI     |
| e14.WT10_AAAACGCGCCGA | LayerV-VI     |
| e14.WT10_AAAATGGACTCA | Striatal_inh2 |

```
[25]: (0.3384487409938317, 0.3289003480647635)
```