

```

1. #=====
2. #===== SLDM hw6, PROBLEM 1 =====
3. #=====
4.
5.
6.
7. # Need to pip install: "phate", "scprep" on the command line before running
8. # the rest of this code in an IDE.
9.
10.
11.
12. #=====
13. #==== Loading the 10X data =====
14. #=====
15. # PART A from the homework
16.
17. # Import the data (code given on the tutorial website at:
18. # https://github.com/KrishnaswamyLab/PHATE/blob/master/Python
19. # /tutorial/EmbryoidBody.ipynb)
20. import os
21. import zipfile
22. from urllib.request import urlopen
23. download_path = os.path.expanduser("~")
24. print(download_path)
25.
26. if not os.path.isdir(os.path.join(download_path, "scRNAseq", "T0_1A")):
27.     if not os.path.isdir(download_path):
28.         os.mkdir(download_path)
29.     zip_data = os.path.join(download_path, "scRNAseq.zip")
30.     if not os.path.isfile(zip_data):
31.         with urlopen("https://data.mendeley.com/datasets/v6n743h5ng/1/files/7489a88f-9ef6-4dff-a8f8-1381d046afe3/scRNAseq.zip?dl=1") as url:
32.             print("Downloading data file...")
33.             # Open our local file for writing
34.             with open(zip_data, "wb") as handle:
35.                 handle.write(url.read())
36.     print("Unzipping...")
37.     with zipfile.ZipFile(zip_data, 'r') as handle:
38.         handle.extractall(download_path)
39.     print("Done.")
40.
41.
42.
43. # Need these libraries
44. import pandas as pd
45. import numpy as np
46. import phate
47. import scprep
48.
49. # Now use scprep to import data into a Pandas dataframe, using the
50. # scprep.io.load_10x function.
51. sparse=True
52. T1 = scprep.io.load_10X(os.path.join(download_path, "scRNAseq", "T0_1A"),
53.                         sparse=sparse,
54.                         gene_labels='both')
55. T2 = scprep.io.load_10X(os.path.join(download_path, "scRNAseq", "T2_3B"),
56.                         sparse=sparse,
57.                         gene_labels='both')
58. T3 = scprep.io.load_10X(os.path.join(download_path, "scRNAseq", "T4_5C"),
59.                         sparse=sparse,
60.                         gene_labels='both')
61. T4 = scprep.io.load_10X(os.path.join(download_path, "scRNAseq", "T6_7D"),

```

```

62.                 sparse=sparse,
63.                 gene_labels='both')
64. T5 = scprep.io.load_10X(os.path.join(download_path, "scRNAseq", "T8_9E"),
65.                 sparse=sparse,
66.                 gene_labels='both')
67. T1.head()
68.
69. # Now, merge all datasets and create a vector representing the time point of
70. # each sample
71. EBT_counts, sample_labels = scprep.utils.combine_batches(
72.     [T1, T2, T3, T4, T5],
73.     ["Day 0-3", "Day 6-9", "Day 12-15", "Day 18-21", "Day 24-27"],
74.     append_to_cell_names=True
75. )
76. del T1, T2, T3, T4, T5 # removes objects from memory
77. EBT_counts.head()
78.
79.
80. #=====
81. #===== Preprocessing: filtering, normalizing, and transforming =====
82. #=====
83. # This is PART B from the homework (technically most of this isn't asked for in
84. # the homework, but I figured it made more sense to follow the tutorial exactly
85.
86. # Remove (suspected) dead cells
87. mito_genes = scprep.utils.get_gene_set(EBT_counts, starts_with="MT-")
88. # Get all mitochondrial genes. There are 14, FYI.
89. scprep.plot.plot_gene_set_expression(EBT_counts, mito_genes, percentile=90)
90. # Plot number of cells that have a certain amount of mitochondrial RNA,
91. # remove cells that are above the 90th percentile. (Line below)
92. EBT_counts, sample_labels = scprep.filter.filter_gene_set_expression(
93.     EBT_counts, mito_genes,
94.     percentile=90,
95.     keep_cells='below',
96.     sample_labels=sample_labels)
97.
98. # Now filter out cells that have either very large or very small library sizes.
99. # Library size is somewhat analogous to sample size. We'll eliminate the
100. # bottom 20% of cells for each sample.
101. scprep.plot.plot_library_size(EBT_counts, percentile=20)
102. EBT_counts, sample_labels = scprep.filter.filter_library_size(
103.     EBT_counts, percentile=20,
104.     keep_cells='above',
105.     sample_labels=sample_labels,
106.     filter_per_sample=True)
107.
108. # Now remove rare genes (genes expressed in 10 or fewer cells)
109. EBT_counts = scprep.filter.remove_rare_genes(EBT_counts, min_cells=10)
110.
111. # Normalization: accounting for differences in library sizes, divide each cell
112. # by its library size and then rescale by the median library size.
113. EBT_counts = scprep.normalize.library_size_normalize(EBT_counts)
114.
115. # Transformation: use square root transform (similar to using log transform
116. # but has the added benefit of dealing with 0's automatically).
117. EBT_counts = scprep.transform.sqrt(EBT_counts)
118.
119.
120. #=====
121. #===== Applying PHATE to data =====
122. #=====

```

```

123.# (In the tutorial, this section is "Embedding Data Using PHATE")
124.
125.# Default parameters for the PHATE function are:
126.# k: number of nearest neighbors, default is 5
127.# a: alpha decay, default is 40
128.# t: number of times to power the operator, default "auto", 21 for these data
129.# gamma: informational distance constant, default is 1.
130.phate.PHATE()
131.phate_operator = phate.PHATE(n_jobs=-2)
132.Y_phate = phate_operator.fit_transform(EBT_counts)
133.
134.# Now plot using phate.plot.scatter2d
135.phate.plot.scatter2d(Y_phate,
136.                      c=sample_labels,
137.                      s=3,
138.                      figsize=(12,8),
139.                      cmap="Spectral")
140.
141.# PART D (from homework): run PHATE on the data using a different value of t.
142.# Plot the PHATE coordinates colored by time point and include the plot.
143.# Based on the results, do you think your chosen value of t is better than the
144.# parameter chosen using the "knee point" of the VNE plot (the default value)?
145.# Will the VNE be higher or lower for your chosen value of t than that
146.# selected (by default) in part(b)?
147.phate_op_2 = phate.PHATE(n_jobs = -2,
148.                          t = 30)
149.Y_phate_2 = phate_op_2.fit_transform(EBT_counts)
150.
151.# Now plot using phate.plot.scatter2d having used t = 30
152.phate.plot.scatter2d(Y_phate_2,
153.                      c=sample_labels,
154.                      s=3,
155.                      figsize=(12,8),
156.                      cmap="Spectral")
157.
158.# PART E (from homework): run PHATE on the data using default parameters to
159.# obtain 3d coordinates. Plot the 3d coordinates. Rotate the plot such that
160.# it's different from what the tutorial has.
161.phate.plot.scatter3d(phate_operator, c=sample_labels, s=3, figsize=(8,6),
162.                    cmap="Spectral")
163.# This saves the 3D plot as a gif
164.phate.plot.rotate_scatter3d(phate_operator, c=sample_labels,
165.                            s=3, figsize=(8,6), cmap="Spectral",
166.                            filename="phate.gif")
167.# This saves the 3D plot as an MP4 (which is also a cool gun in my opinion)
168.phate.plot.rotate_scatter3d(phate_operator, c=sample_labels,
169.                            s=3, figsize=(8,6), cmap="Spectral",
170.                            filename="phate.mp4")

```