

- This homework is due at 11:59pm on February 24, 2023. **Please submit as a PDF** by email to natalies@cs.unc.edu+comp790.
- Homework not submitted as a PDF will not be graded. If you use a Jupyter notebook, please convert to a PDF.
- There are 3 data files provided for the following questions, including, `Levine_matrix.csv`, `cell_graph.edgelist`, and `population_names_Levine_13dim.txt`. Instructions for how to use these data will be provided in each homework problem.
- You are welcome to any available libraries, like numpy and networkx.
- Please feel free consult with other colleagues, but please write up your own independent solution.
- You are welcome to use Python, Julia, or R here. All hints and sparse code are given for Python.
- You are welcome to write up your assignment using the `HW1_790-166.tex` template, or write up the solutions in the method of your choice.
- This homework is worth 70 points total.

## Problem 1

**Warm Up - Counting and Adjacency Matrix Math** (9 Points Total) Consider an undirected, unweighted graph with  $N$  nodes and its corresponding adjacency matrix,  $\mathbf{A}$ .

1. (5 points) Consider a graph with  $N$  nodes and its corresponding adjacency matrix,  $\mathbf{A} \in \mathbb{R}^{N \times N}$ . Let  $\mathbf{1} \in \mathbb{R}^N$  be the column vector of only 1s (e.g.  $N$  1s.) Using  $\mathbf{A}$  and  $\mathbf{1}$ , write an expression for a vector of node degrees,  $\mathbf{k}$ , such that the  $i$ th entry of  $\mathbf{k}$ ,  $k_i$ , represents the number of neighbors of node  $i$ . (Hint, you should be writing this expression in terms of  $\mathbf{A}$  and  $\mathbf{1}$ .)
2. (3 points) Again, using  $\mathbf{A}$  and  $\mathbf{1}$  write an expression for the total number of edges in the graph.
3. (1 point) Verify the two expressions that you just defined by plugging the  $\mathbf{A}$  defined below. Show that you get  $\mathbf{k} = [1, 2, 1]$  and that the number of edges in the graph is 2. You can show this by drawing the graph. You can either verify this by writing out the algebra or by writing code to compute the expressions and showing the output.

$$\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$$

## Problem 2

**Nice Properties of the Graph Laplacian** (14 Points Total)

- Here we will walk through computing the Graph Laplacian and exploring nice properties about it that relate to graph structure. You are welcome to paste the code that you write for each of these sub-problems.

Consider the following graph,  $\mathcal{G}$  of 6 nodes. This graph has two components. We will see how what happens with our graph Laplacian in this situation.

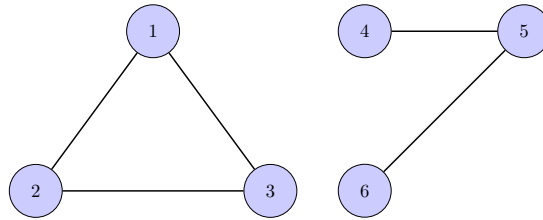


Figure 1: Graph,  $\mathcal{G}$  with six nodes and two components.

1. (3 points) Create the adjacency matrix,  $\mathbf{A}$ , corresponding to  $\mathcal{G}$ . (Note that node  $i$  should correspond to the  $i$ th row of  $\mathbf{A}$ .)
2. (2 points) We will go through some steps to write the code to compute the Graph Laplacian. The first matrix we need to define is the diagonal degree matrix,  $\mathbf{D}$ . Create the degree matrix ( $\mathbf{D}$ ) for  $\mathcal{G}$ .
3. (2 points) Define the Laplacian matrix,  $\mathbf{L}$  as  $\mathbf{L} = \mathbf{D} - \mathbf{A}$ .
4. (2 points) An eigenvector ( $\mathbf{v} \in \mathbb{R}^6$ ) and its corresponding eigenvalue ( $\lambda$ ) of the graph Laplacian,  $\mathbf{L}$ , must satisfy  $\mathbf{L}\mathbf{v} = \lambda\mathbf{v}$ . As this graph has two components, we will be able to find two unique eigenvectors ( $\mathbf{v}_1$  and  $\mathbf{v}_2$ ) that correspond to the zero eigenvalue, such that  $\mathbf{L} \times \mathbf{v}_1 = 0 \times \mathbf{v}_1$  and  $\mathbf{L} \times \mathbf{v}_2 = 0 \times \mathbf{v}_2$ . Consider  $\mathbf{v}_1 = [1, 1, 1, 0, 0, 0]^T$ . You will notice that there are 1s in the first three entries, corresponding to the nodes that are together in a component. Check that  $\mathbf{v}_1$  does indeed correspond to eigenvalue 0 and satisfies  $\mathbf{L} \times \mathbf{v}_1 = 0 \times \mathbf{v}_1$ .

If you are using numpy, you can check by computing,

```
L @ np.matrix([[1],[1],[1],[0],[0],[0]])
```

Discuss your findings.

5. (3 points) Given this pattern you just observed with  $\mathbf{v}_1$ , find  $\mathbf{v}_2$ . Show that your  $\mathbf{v}_2$  satisfies  $\mathbf{L} \times \mathbf{v}_2 = 0 \times \mathbf{v}_2$ .

You can fill in the following with your entries of  $\mathbf{v}_2$ .

```
L @ np.matrix([[ ],[ ],[ ],[ ],[ ],[ ]])
```

6. (2 points) As a further numerical experiment, find the eigenvalues of  $\mathbf{L}$  (you can use <https://numpy.org/doc/stable/reference/generated/numpy.linalg.eig.html#numpy.linalg.eig>). What do you notice about the smallest two eigenvalues? What about the third smallest eigenvalue?

## Problem 3

**(A Single-Cell Graph in the Wild)** (3 points) A graph,  $\mathcal{G}_{cell}$  was constructed between 5000 cells contained in `cell_graph.edgelist`. Make sure for the subsequent problems you interpret this as an undirected graph and the adjacency matrix is symmetric.

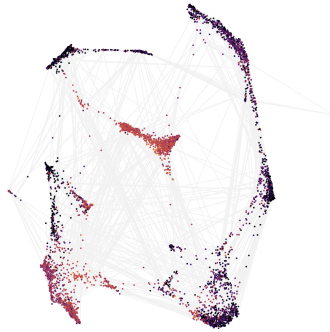


Figure 2: A visualization of the graph in homework problem 3.

- (3 points) Use all of your hard work from Problem 2 to compute the graph Laplacian,  $\mathbf{L}$  for  $\mathcal{G}_{cell}$ . Make a histogram to visualize the distribution of eigenvalues of  $\mathbf{L}$ . What is the smallest eigenvalue of  $\mathbf{L}$ ? What are the smallest and second smallest eigenvalues of  $\mathbf{L}$  and how does the observation relate to the structure of the graph?

## Problem 4

### (14 Points Total) Playing with Single Cell Data

We will consider data from a CyTOF experiment obtained from <http://flowrepository.org/id/FR-FCM-ZZPH>. Here, we are considering the expression of 13 different protein markers across a set of cells. It has already been pre-processed for you. From the entire dataset, 5,000 cells were sampled for further analysis. You can use the following accompanying data as follows.

- `Levine_matrix.csv` is the cell  $\times$  marker matrix. Note that the last column is labels for the cells. Let's call this matrix,  $\mathbf{X}$ . **You should not use the last column (named label) for any kind of clustering.** Some cells are not labeled (hence are called NaN).
- `population_names_Levine_13dim.txt` maps the cell labels from the last column of  $\mathbf{X}$  (number labels) to biologically-interpretable cell-type names.
- `cell_graph.edgelist` is an edgelist for a between-cell graph. We will call this,  $\mathbf{G}$ . Note that the nodes in  $\mathbf{G}$  correspond to the rows in  $\mathbf{X}$ . So, node  $i$  maps to row  $i$  of  $\mathbf{X}$ , etc.

1) **Clustering on cell  $\times$  marker data** (7 points): Use a clustering algorithm of your choice to generate a cell-to-cluster partition for the cells, using the matrix,  $\mathbf{X}$ . Use normalized mutual information (NMI) to compute overlap between the true and predicted cell labels. **Note that because not all cells are labeled, you can compute this only based on the labeled cells.** Feel free to use an available implementation, such as, [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.normalized\\_mutual\\_info\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.normalized_mutual_info_score.html).

2) **Graph Partitioning** (7 points): Use a graph clustering algorithm to partition  $\mathbf{G}$  into clusters. Similar to part 1 (1) compute NMI between the labels obtained in graph partitioning and the true cell labels.

3) (3 points) Comment on any observations you observe between the quality of the partitions obtained clustering on  $\mathbf{X}$  in comparison to partitioning  $\mathbf{G}$ . Which approach do you think works better, using the original data, or the graph?

4) **Rare Cell-types** (5 points) `Plasmacytoid_DC_cells`, or pDCs (label 21) are a popular rare cell type, meaning many clustering algorithms will not be able to reliably find them. Report the number of distinct clusters where you found pDCs in both the clustering of  $\mathbf{X}$  and in the partitioning of  $\mathbf{G}$ .

5) **Cell Classification** (10 points) Select cells from  $\mathbf{X}$  with the following labels,  $\{11, 12, 17, 18\}$  and  $\{1, 2, 3\}$ . In general, cells with labels  $\{11, 12, 17, 18\}$  are T-cells and cells with labels  $\{1, 2, 3\}$  are monocytes. Convert this to a binary classification problem by labeling T-cells with 0 and monocytes with 1. Use your favorite classifier to predict the labels of these cells. Use an ROC curve to visualize the performance. Please be sure to use an appropriate training/testing or cross validation procedure. If the performance was not good, explain what could have gone wrong. If your performance is very good, can you identify features from  $\mathbf{X}$  that were helpful in predicting labels?

## Problem 5

(30 points total) **node2vec**

We will use the implementation of node2vec available in github, <https://github.com/eliorc/node2vec>.

1) **(Clustering on Node2Vec Features** (10 points)) First, use default parameters and follow the instructions in the README on the graph in `cell_graph.edgelist`. This will create a 128-dimensional vector for each node. Cluster the nodes based on these vectors and compare to the ground truth labels in the last column of `Levine_matrix.csv` using NMI. Compare your results to Problem 4, question 3. Does an embedding of the graph offer any apparent advantages in classifying cells?

2) **(Parameters, part 1** (5 points)) Try a few different values for the number of dimensions `--dimensions`, such that some of them are less than 128, and some of them are more than 128. Cluster cells again with the embeddings obtained in different dimensions. Again, you can compute the NMI between the cluster assignments and the ground truth labels. Comment on some observations, and show a plot of NMI plotted against the number of dimensions used.

3) **(Parameters, part 2** (5 points)) Recall that the parameters  $p$  and  $q$  control the ‘breadth’ vs ‘depth of the walk’. Choose one of these parameters to vary, and repeat the previous question using the default 128 dimensions, but varying values for either  $p$  or  $q$ . Comment on some observations, and show a plot of NMI against  $p$  or  $q$  (whichever one you chose).

4) **(Cell Classification, Part II** (10 points)) Repeat Problem 4, question (5). However, instead of using only  $\mathbf{X}$  as the feature matrix, we are going to combine the marker expressions with node2vec features. Let  $\mathbf{N}$  be your matrix generated through node2Vec. Create a new matrix called  $\mathbf{X} = [\mathbf{X}|\mathbf{N}]$ . That is, you will simply concatenate  $\mathbf{X}$  and  $\mathbf{N}$ . Formulate the same classification problem from Problem 4, question (5) to classify T-cells from monocytes. Again, report your ROC curve. Please use an appropriate training/testing or cross-validation procedure. Comment on the performance, especially in comparison to the results obtained in Problem 4, question (5).