

ADVANCED LEARNING FOR TEXT AND GRAPH DATA

Lab session 5: Graph Mining

Lecture: Prof. Michalis Vazirgiannis

Lab: Giannis Nikolentzos

Monday, December 16, 2019

This handout includes theoretical introductions, [coding tasks](#) and [questions](#). Before the deadline, you should submit here a **.zip** file (max 10MB in size) containing a `/code/` folder (itself containing your scripts with the gaps filled) and an answer sheet named `firstname_lastname.pdf`, following the template available [here](#), and containing your answers to the questions. Your answers should be well constructed and well justified. They should not repeat the question or generalities in the handout. When relevant, you are welcome to include figures, equations and tables derived from your own computations, theoretical proofs or qualitative explanations. **One submission is required for each student. The deadline for this lab is December 31, 2019 11:59 PM.** No extension will be granted. Late policy is as follows: $]0, 24]$ hours late \rightarrow -5 pts; $]24, 48]$ hours late \rightarrow -10 pts; > 48 hours late \rightarrow not graded (zero).

1 Learning objective

In this lab, you will implement some basic techniques for dealing with different graph mining problems. Specifically, the lab is divided into three parts. In the first part, you will study the dynamics of a real-world graph. Then, you will use some clustering algorithms to reveal its community structure. Finally, you will use graph kernels to measure the similarity between graphs and to perform graph classification. We will use Python 3.6, and the NetworkX library (<http://networkx.github.io/>).

2 Analyzing a Real-World Graph

In this part of the lab, we will analyze the CA-HepTh collaboration network, examining several structural properties. The Arxiv HEP-TH (High Energy Physics - Theory) collaboration network comes from the e-print arXiv and covers scientific collaborations between authors of papers submitted to the High Energy Physics - Theory category. If an author i co-authored a paper with author j , the graph contains an undirected edge from i to j .

2.1 Load Graph and Simple Statistics

The graph is stored in the `CA-HepTh.txt` file¹, as an edge list:

¹The graph can be downloaded from the following link: <https://snap.stanford.edu/data/ca-HepTh.txt.gz>.

```
# Directed graph (each unordered pair of nodes is saved once): CA-HepTh.txt
# Collaboration network of Arxiv High Energy Physics Theory category (there is an edge i
# Nodes: 9877 Edges: 51971
# FromNodeId      ToNodeId
24325      24394
24325      40517
24325      58507
24394      3737
24394      3905
24394      7237
...
```

Let's first create an undirected NetworkX graph based on the data contained in this file.

Task 1

Load the network data into an undirected graph G , using the `read_edgelist()` function of NetworkX. Note that, the delimiter used to separate values is the tab character `\t` and additionally, that lines that start with the `#` character are comments. Furthermore, compute and print the following network characteristics: (1) number of nodes, (2) number of edges.

Question 1 (3 points)

What is the maximum number of edges and the maximum number of triangles an undirected graph of n nodes without self-loops can have?

2.2 Connected Components

A connected component of an undirected graph is a subgraph in which any two vertices are connected to each other by one or more paths.

Task 2

Print the number of connected components. If the graph is not connected, retrieve the largest connected component subgraph (also known as *giant connected component*) (Hint: you can use the `connected_components()` function of NetworkX). Find the number of nodes and edges of the largest connected component and examine in what fraction of the whole graph they correspond.

2.3 Analysis of the Degree Distribution of the Graph

Extract the degree sequence of the graph using the following code:

```
degree_sequence = [G.degree(node) for node in G.nodes()]
```

The above code returns a list that contains the degree of all the nodes of the graph.

Task 3

Find and print the minimum, maximum, and mean degree of the nodes of the graph (Hint: you can use the built-in functions `min()`, `max()`, `mean()` of the NumPy library).

Let's now compute and plot the degree distribution of the graph.

Task 4

Plot the degree histogram using the `matplotlib` library of Python (Hint: use the `degree_histogram()` function that returns a list of the frequency of each degree value). Produce again the plot using log-log axis.

Question 2 (3 points)

What do you observe? How this type of distribution is called?

3 Community Detection

In the second part of the lab, we will focus on the community detection (or clustering) problem in graphs. Typically, a community corresponds to a set of nodes that highly interact among each other, compared to the intensity of interactions (as expressed by the number of edges) with the rest nodes of the graph. The experiments for this part will also be performed in the CA-HepTh collaboration network.

3.1 Spectral Clustering

We will first implement and apply a very popular graph clustering algorithm, called *Spectral Clustering* [4]. The basic idea of the algorithm is to utilize information associated to the spectrum of the graph, in order to identify well-separated clusters. Algorithm 1 illustrates the pseudocode of Spectral Clustering.

Algorithm 1 Spectral Clustering

Input: Graph $G = (V, E)$ and parameter k

Output: Clusters C_1, C_2, \dots, C_k (i.e., cluster assignments of each node of the graph)

- 1: Let \mathbf{A} be the adjacency matrix of the graph
 - 2: Compute the Laplacian matrix $\mathbf{L} = \mathbf{D} - \mathbf{A}$. Matrix \mathbf{D} corresponds to the diagonal degree matrix of graph G (i.e., degree of each node v (= number of neighbors) in the main diagonal)
 - 3: Apply eigenvalue decomposition to the Laplacian matrix \mathbf{L} and compute the eigenvectors that correspond to d smallest eigenvalues. Let $\mathbf{U} = [\mathbf{u}_1 | \mathbf{u}_2 | \dots | \mathbf{u}_d] \in \mathbb{R}^{m \times d}$ be the matrix containing these eigenvectors as columns
 - 4: For $i = 1, \dots, m$, let $y_i \in \mathbb{R}^d$ be the vector corresponding to the i -th row of \mathbf{U} . Apply k -means to the points $(y_i)_{i=1, \dots, m}$ (i.e., the rows of \mathbf{U}) and find clusters C_1, C_2, \dots, C_k
-

We will now implement the Spectral Clustering algorithm.

Task 5

Fill in the body of the `spectral_clustering()` function which implements the Spectral Clustering algorithm. The algorithm must return a dictionary keyed by node to the cluster to which the node belongs (Hint: to perform k -means, you can use `scikit-learn`'s implementation of the algorithm).

Question 3 (4 points)

Why spectral clustering focuses on the smallest eigenvalues of the Laplacian matrix \mathbf{L} ? What is the problem that is optimized by the eigenvalue decomposition?

Task 6

Apply the Spectral Clustering algorithm to the giant connected component of the CA-HepTh dataset, trying to identify 50 clusters.

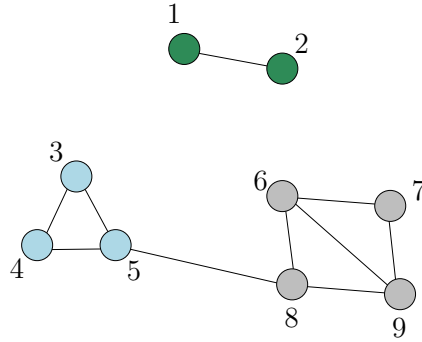


Figure 1: A graph where nodes have been assigned to 3 clusters.

3.2 Modularity

To assess the quality of a clustering algorithm, several metrics have been proposed. *Modularity* is one of the most popular and widely used metrics to evaluate the quality of a network's partition into communities [2]. Considering a specific partition of the network into clusters, **modularity measures the number of edges that lie within a cluster compared to the expected number of edges of a null graph (or configuration model), i.e., a random graph with the same degree distribution.** In other words, the measure of modularity is built upon the idea that random graphs are not expected to present inherent community structure; thus, comparing the observed density of a subgraph with the expected density of the same subgraph in case where edges are placed randomly, leads to a community evaluation metric. Modularity is given by the following formula:

$$Q = \sum_{c=1}^{n_c} \left[\frac{l_c}{m} - \left(\frac{d_c}{2m} \right)^2 \right]$$

where, $m = |E|$ is the total number of edges in the graph, n_c is the number of communities in the graph, l_c is the number of edges within the community c and d_c is the sum of the degrees of the nodes that belong to community c . Modularity takes values in the range $[-1, 1]$, with higher values indicating better community structure.

Question 4 (3 points)

Compute (showing your calculations) the modularity of the clustering result shown in Figure 1. Note that different colors correspond to different clusters.

Task 7

Fill in the body of the `modularity()` function that computes the modularity of a clustering result.

Next, we will use modularity to evaluate two clustering results of the nodes of the giant connected component of the CA-HepTh dataset.

Task 8

Compute the modularity of the following two clustering results: (i) the one obtained by the Spectral Clustering algorithm using $k = 50$, and (ii) the one obtained if we randomly partition the nodes into 50 clusters (Hint: to assign each node to a cluster, use the `randint(a, b)` function which returns a random integer n such that $a \leq n \leq b$).

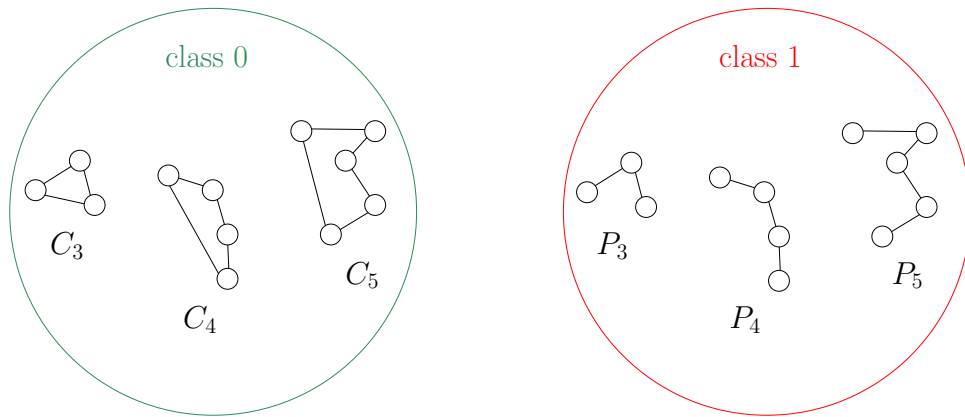


Figure 2: Dataset consisting of two sets of graphs: cycle graphs (left) and path graphs (right).

4 Graph Classification using Graph Kernels

In the last part of the lab, we will focus on the problem of graph classification. Graph classification arises in the context of a number of classical domains such as chemical data, biological data, and the web. In order to perform graph classification, we will employ **graph kernels, a powerful framework for graph comparison.**

Kernels can be intuitively understood as functions measuring the similarity of pairs of objects. More formally, for a function $k(x, x')$ to be a kernel, it has to be (1) symmetric: $k(x, x') = k(x', x)$, and (2) positive semi-definite. If a function satisfies the above two conditions on a set \mathcal{X} , it is known that there exists a map $\phi : \mathcal{X} \rightarrow \mathcal{H}$ into a Hilbert space \mathcal{H} , such that $k(x, x') = \langle \phi(x), \phi(x') \rangle$ for all $(x, x') \in \mathcal{X}^2$ where $\langle \cdot, \cdot \rangle$ is the inner product in \mathcal{H} . Kernel functions thus compute the inner product between examples that are mapped in a higher-dimensional feature space. However, they do not necessarily explicitly compute the feature map ϕ for each example. One advantage of kernel methods is that they can operate on very general types of data such as images and graphs. Kernels defined on graphs are known as *graph kernels*. **Most graph kernels decompose graphs into their substructures** and then to measure their similarity, they count the number of common substructures. Graph kernels typically focus on some structural aspect of graphs such as random walks, shortest paths, subtrees, cycles, and graphlets.

4.1 Dataset Generation

We will first create a very simple graph classification dataset. The dataset will contain two types of graphs: (1) cycle graphs, and (2) path graphs. A cycle graph C_n is a graph on n nodes containing a single cycle through all nodes, while a path graph P_n is a tree with two nodes of degree 1, and all the remaining $n - 2$ nodes of degree 2. Each graph is assigned a class label: label 0 if it is a cycle or label 1 if it is a path. Figure 2 illustrates such a dataset consisting of three cycle graphs and three path graphs. Use the `cycle_graph()` and `path_graph()` functions of NetworkX to generate 100 cycle graphs and 100 path graphs of size $n = 3, \dots, 102$, respectively. Store the 200 graphs in a list and their class labels in another list.

Task 9

Fill in the body of the `create_dataset()` function to generate the dataset as described above.

Before computing the kernels, it is necessary to split the dataset into a training and a test set. We can use the `train_test_split()` function of scikit-learn as follows:

```
from sklearn.model_selection import train_test_split

G_train, G_test, y_train, y_test = train_test_split(Gs, y, test_size=0.1)
```

4.2 Implementation of Graphlet Kernel

We will next investigate if graph kernels can distinguish cycle graphs from path graphs. We will use the following two graph kernels: (1) **shortest path kernel**, and (2) **graphlet kernel**. This kernel has already been implemented for you (i.e., `shortest_path_kernel()` function). The shortest path kernel counts the number of **shortest paths of equal length** in two graphs [1]. It can be shown that in the case of unlabeled graphs, the kernel maps the graphs into a feature space where each feature corresponds to a shortest path distance and the value is equal to the frequency of that distance in the graph (see Figure 3 for an illustration).

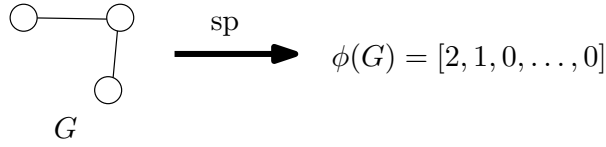


Figure 3: Example of feature map of the shortest path kernel. There are 2 shortest paths of distance 1 and 1 shortest path of distance 2 in this graph.

Question 5 (3 points)

Give an example of two non-isomorphic graphs which are mapped to the same representation by the shortest path kernel?

The graphlet kernel decomposes graphs into graphlets (i.e., small subgraphs with k nodes where $k \in \{3, 4, 5\}$) and counts matching graphlets in the input graphs [3]. For example, the set of graphlets of size 3 is shown in Figure 4 below.

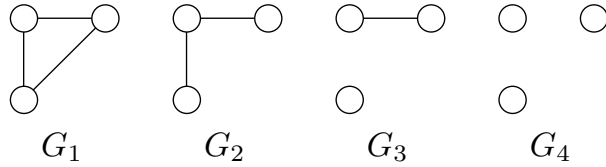


Figure 4: Set of the graphlets of size 3.

The graphlet kernel samples a number of small subgraphs from a graph, and computes their distribution. Here, we will focus on graphlets of size 3. Let $\{\text{graphlet}_1, \text{graphlet}_2, \text{graphlet}_3, \text{graphlet}_4\}$ be the set of size-3 graphlets (i.e., those shown in Figure 4). The graphlet kernel maps each graph G into a vector $f_G \in \mathbb{N}^4$ whose i -th entry is equal to the number of sampled subgraphs from G that are isomorphic to graphlet_i . Then, the graphlet kernel is defined as follows:

$$k(G, G') = f_G^\top f_{G'} \quad (1)$$

Given a set of training graphs (with cardinality N_1), a set of test graphs (with cardinality N_2) and a graph kernel, we are interested in generating two matrices. A symmetric matrix $\mathbf{K}_{train} \in \mathbb{R}^{N_1 \times N_1}$ which contains the kernel values for all pairs of training graphs, and a second matrix $\mathbf{K}_{test} \in \mathbb{R}^{N_2 \times N_1}$ which stores the kernel values between the graphs of the test set and those of the training set. For the shortest path kernel, we can produce these two matrices as follows:

```
K_train_sp, K_test_sp = shortest_path_kernel(G_train, G_test)
```

We will next implement the graphlet kernel.

Task 10

Fill in the body of the `graphlet_kernel()` function. The function generates the feature maps of equation 1 by sampling `n_samples` size-3 graphlets from each graph. Then, it generates the K_{train} and K_{test} matrices by computing the inner products between the feature maps (Hint: you can use the `random.choice()` function of NumPy to sample 3 nodes from the set of nodes of a graph. Given a set of nodes `s`, use the `G.subgraph(s)` function of NetworkX to obtain the subgraph induced by set `s`. To test if a subgraph is isomorphic to a graphlet, use the `is_isomorphic()` function of NetworkX).

Task 11

Use the `graphlet_kernel()` function that you implemented to compute the kernel matrices associated with the graphlet kernel.

4.3 Graph Classification using SVM

After generating the K_{train} and K_{test} matrices, we can use the SVM classifier to perform graph classification. More specifically, as shown below, we can directly feed the kernel matrices to the classifier to perform training and make predictions:

```
from sklearn.svm import SVC

# Initialize SVM and train
clf = SVC(kernel='precomputed')
clf.fit(K_train, y_train)

# Predict
y_pred = clf.predict(K_test)
```

Task 12

Train two SVM classifiers (i.e., one using the kernel matrix generated by the shortest path kernel, and the other using the kernel matrix generated by the graphlet kernel). Then, use the two classifiers to make predictions. Evaluate the two kernels (i.e., shortest path and graphlet) by computing the classification accuracies of the corresponding models (Hint: use the `accuracy_score()` function of scikit-learn).

Question 6 (4 points)

Compare the performance of the two kernels (i.e., shortest path and graphlet). What do you observe? Did the graphlet kernel achieve a high accuracy? Why?

References

- [1] Karsten M Borgwardt and Hans-Peter Kriegel. Shortest-path kernels on graphs. In *Proceedings of the 5th IEEE International Conference on Data Mining*, pages 74–81, 2005.
- [2] Mark EJ Newman. Modularity and community structure in networks. *Proceedings of the national academy of sciences*, 103(23):8577–8582, 2006.

- [3] Nino Shervashidze, SVN Vishwanathan, Tobias Petri, Kurt Mehlhorn, and Karsten M Borgwardt. Efficient Graphlet Kernels for Large Graph Comparison. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, pages 488–495, 2009.
- [4] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.