

Lab activity: analysis of music networks.

1 Session 2: Data preprocessing

1.1 Introduction

In this second session of the practice, we will work on the preprocessing of the data downloaded in the first session. In other words, we will prepare the data to make it easier to extract information from it later on using analysis algorithms.

Firstly, we will implement three functions that work with NetworkX graphs, which will allow us to convert directed graphs into undirected graphs and prune graphs according to different criteria. Then, we will create two additional functions that will generate artist similarity graphs based on the audio features of their songs. Finally, we will perform an initial exploration of the original data and the data obtained with these preprocessing functions.

1.2 Data preprocessing

In this second session of the practice, we will use the artist graphs and the song dataset that we obtained in the previous session. Along with the instructions for this second session, a file is attached specifying the functions that need to be implemented for each activity (the file `Lab.AGX.202223.S2.skeleton.py`).

1. (1 point) Implement the function `retrieve_bidirectional_edges` with the following header:

```
1 def retrieve_bidirectional_edges(g: nx.DiGraph, out_filename: str) -> nx.  
  Graph:  
2     """  
3     Convert a directed graph into an undirected graph by considering  
  bidirectional edges only.  
4  
5     :param g: a networkx digraph.  
6     :param out_filename: name of the file that will be saved.  
7     :return: a networkx undirected graph.  
8     """
```

The function will receive a directed graph as input and return an undirected graph. The resulting undirected graph will have as edges those edges from the original graph that existed in both directions. In other words, the edge $e = (v_i, v_j)$ will exist in the new graph if and only if both edges $e = (v_i, v_j)$ and $e' = (v_j, v_i)$ existed in the input graph. The nodes of the new graph will be defined by the edges contained in the new graph (only nodes that have incident edges in the new graph should appear).

The function `retrieve_bidirectional_edges` will return the resulting graph and, additionally, it will save it to disk in a file with the specified name as a parameter (in `graphml` format).

2. (1 point) Implement the function `prune_low_degree_nodes` with the following header:

```
1 def prune_low_degree_nodes(g: nx.Graph, min_degree: int, out_filename: str)
  -> nx.Graph:
2     """
3     Prune a graph by removing nodes with degree < min_degree.
4
5     :param g: a networkx graph.
6     :param min_degree: lower bound value for the degree.
7     :param out_filename: name of the file that will be saved.
8     :return: a pruned networkx graph.
9     """
```

The function will receive an undirected input graph and generate a graph from which all nodes with a degree less than `min_degree` have been removed. The function will have three input parameters: an undirected `networkx` graph, a variable with the minimum degree value, and the output file name where the generated graph will be saved.

This removal of nodes with a degree less than `min_degree` should be done in a single pass. In other words, at the end of the traversal of all nodes, there may still be nodes with a degree less than `min_degree`. In this case, we will keep these nodes. After this process, you should remove the zero-degree nodes from the resulting graph.

3. (1.5 points) Implement the function `prune_low_weight_edges` with the following header:

```
1 def prune_low_weight_edges(g: nx.Graph, min_weight: float = None,
  min_percentile: int = None,
2                               out_filename: str = None) -> nx.Graph:
3     """
4     Prune a graph by removing edges with weight < threshold. Threshold can
  be specified as a value or as a percentile.
5
6     :param g: a weighted networkx graph.
7     :param min_weight: lower bound value for the weight.
8     :param min_percentile: lower bound percentile for the weight.
9     :param out_filename: name of the file that will be saved.
10    :return: a pruned networkx graph.
11    """
```

The function will receive an undirected input graph with weighted edges and generate a graph from which all edges with a weight less than the value specified as a parameter will be removed. This weight can be specified in two different ways: either directly with the threshold value (parameter `min_weight`) or with the percentile value (parameter `min_percentile`). The function should raise an exception if neither of the two parameters is specified or if both parameters are specified (i.e., the function call should only specify one of the two parameters).

After the edge removal process, you should remove the zero-degree nodes from the resulting graph.

In addition to returning the resulting graph, the function will also save it to disk in the file specified by the parameter name (in `graphml` format).

4. (1 point) Implement the function `compute_mean_audio_features` with the following header:

```

1 def compute_mean_audio_features(tracks_df: pd.DataFrame) -> pd.DataFrame:
2     """
3     Compute the mean audio features for tracks of the same artist.
4
5     :param tracks_df: tracks dataframe (with audio features per each track)
6     :return: artist dataframe (with mean audio features per each artist).
7     """

```

The function will receive a **dataframe** of songs (the result of the `get_track_data` function implemented in session 1) and will return another **dataframe** with the average audio characteristics of each artist. In other words, the resulting **dataframe** will have one row for each artist, which will contain both the identification data of that artist (at least, the identifier and name) and the average values of each audio characteristic (*danceability*, *energy*, *loudness*, ...) of all the songs by that artist.

5. (1.5 points) Implement the function `create_similarity_graph` with the following header:

```

1 def create_similarity_graph(artist_audio_features_df: pd.DataFrame,
2                             similarity: str, out_filename: str = None) -> nx.Graph:
3     """
4     Create a similarity graph from a dataframe with mean audio features per
5     artist.
6
7     :param artist_audio_features_df: dataframe with mean audio features per
8     artist.
9     :param similarity: the name of the similarity metric to use (e.g. "
10    cosine" or "euclidean").
11    :param out_filename: name of the file that will be saved.
12    :return: a networkx graph with the similarity between artists as edge
13    weights.
14    """

```

The function will receive a **dataframe** with the average audio characteristics for different artists and a similarity measure, and will construct a complete graph (undirected and with weighted edges) that stores the similarity between each pair of artists as the weight of the connecting edge.

The function will return the resulting graph and also save it to the specified file (in `graphml` format).

6. (1 point) Using the previous functions, obtain:

- Two undirected graphs ($g'B$ and $g'D$) of artists obtained by applying the programmed function in exercise 1, `retrieve_bidirectional_edges`, to the graphs obtained by the crawler in session 1, g_B and g_D .
- Two undirected graphs with weights (g_B^w and g_D^w) obtained from the similarity between the artists. To obtain them, it will be necessary to calculate the vector of average audio features for each artist (`compute_mean_audio_features`), create a similarity graph with these features (`create_similarity_graph`), and prune the resulting graph (`prune_low_weight_edges`) to achieve the desired size. Specifically, the size of the graph should be as similar as possible to the size of graphs $g'B$ and $g'D$, respectively.¹

¹As a reference, the difference should not exceed 50 edges.

1.3 Questions to be addressed in the report

1. **(0.5 points)** Provide the order and size of the four obtained undirected graphs ($g'B$, $g'D$, g_B^w , and g_D^w).
2. **(1 point)** Justify the strategy used to obtain g_B^w and g_D^w .
3. **(0.5 points)** Justify whether the directed graphs obtained from the initial exploration of the crawler (g_B and g_D) can have more than one weakly connected component and strongly connected component, and explain why. Indicate the relationship with the selection of a single seed.
4. **(0.5 points)** Also justify the relationship between the previous results and the number of connected components in the undirected graphs (g'_B and g'_D).
5. **(0.5 points)** Compute the size of the largest connected component from g'_B and g'_D . Which one is bigger? Justify the result.

1.4 Evaluation and expected results for this part of the practice

This second session of the practice will account for 20% of the total grade for the practice.

Remember that it is important that **the Python code includes sufficient comments** to understand its functionality and **respect the headers** of the provided functions.

The files to be obtained in this first part of the practice for the final submission are as follows:

- The file `Lab_AGX_202223_S2_skeleton.py` where each of the functions defined in this document has been implemented.
- A **pdf** file where all the questions posed in this statement have been answered in detail. The answers to the questions should be ordered.
- The data files obtained in exercise 6 (four **graphml** files, containing the graphs $g'B$, $g'D$, g_B^w , and g_D^w).