Popular Erasmus destinations

Alessandro Lotta, Youssef Ben Khalifa

January 30, 2023

Contents

1	Dataset creation and pre-processing	2
2	Graph creation	2
3	Analysis performed	2
4	Analysis results and comparisons	4
5	Further analysis and improvements	6

1 Dataset creation and pre-processing

To manage the large amount of data we had to work with, we decided to make use of a custom created class called **Dataset**, in which we define all the methods and objects we used.

The dataset class makes use of the **pandas** library, which is a Python library for data manipulation and analysis, with which we perform the import, the preprocessing and the manipulation of the datasets.

Dataset cleaning

The first thing we need to do is to clear out the entries that are either incomplete or not useful for our analysis. To do that, we start by applying a simple filter on the dataset to select only the columns we are interested in, to then remove all the rows that have at least one missing value. From the resulting dataset we can start to filter out the type of entries we need for each of our analysis: for example if we want to perform our analysis only on the students that are currently studying in a university, we can filter out all the entries

2 Graph creation

The entire project is based on the usage of the python library **NetworkX** which is a Python library for the creation, manipulation and study of the structure, dynamics, and functions of complex networks. through the usage of this library we were able to create the graphs we needed to perform our analysis.

Through the project m

3 Analysis performed

In this section we will go over all the analysis we have performed on the graph we generated from the data we had. As we said in the project proposal we will use two graphs $C = (V_c, E_c)$ and $U = (V_u, E_u)$: one having countries as nodes and the other one having universities as nodes, in which the edges and their weights represent the amount of students moving/received from one node to the other.

PageRank coefficient

The **PageRank** coefficient of a node v expresses the "importance" of a node in the graph, this is done by considering the number of incoming edges and the PageRank coefficient of the nodes that

are connected to it. Analytically, the PageRank coefficient of a node v is defined as:

$$Pr(v) = (1 - d) + d \sum_{u \in V} \frac{Pr(u)}{deg(u)}$$
 (3.1)

where d is the **damping factor** given by the user, V is the set of nodes in the graph and deg(u) is the degree of node u.

This particular feature is very useful when we compute it on the graph in which we define as the set of nodes and edges, respectively the countries/universities and as edges both the students that are sent and received. We can find a measure of how important a country or a university can be in terms of students flow.

PageRank coefficient implementation

We implemented the PageRank algorithm using the **NetworkX** library, which is a Python library for the creation, manipulation and study of the structure, dynamics, and functions of complex networks. The implementation of the algorithm can be found in the appendix.

Closeness centrality

The Closeness centrality of a node v is mathematically defined as:

$$C(v) = \frac{n-1}{\sum_{u \in V} d(v, u)}$$
 (3.2)

where V is the set of nodes in the graph and d(v, u) is the length of the shortest path between nodes v and u. The Closeness centrality of a node v is yet another measure through which we can obtain useful information from our graphs. In particular, computing the Closeness centrality on the graph which is defined using

- as nodes the set of countries and/or universities;
- as edges the amount the students sent from the country/university to another country/university;

we can determine the "well" the students from that specific country/university are distributed over Europe.

Closeness centrality implementation

Once again, we used the **NetworkX** library to implement the Closeness centrality algorithm. The implementation can be found in the appendix.

4 Analysis results and comparisons

Each analysis was done on two different machines in order to try and extrapolate a measure of the efficiency of the algorithms we used. However, the comparisons between the different time results on the algorithm execution are not very reliable, since the machines used are affected by many other factors other that the hardware itself.

The main focus of the analysis is on how the algorithms implementations we adopted perform w.r.t. the size of the graph we created.

To analyze and compare the results we decided to concentrate only on the Students that participate on the Erasums program while being subscribed to a Master Degree, to do this we applied a filter on the "Education Level" column selecting the value "ISCED-7".

Countries Results

The results obtained are presented by showing the top countries based on the values obtained from PageRank and HITS.

Country	PageRank	Country	Authorities	-	Country	Hubs
France	5.719243	France	9.997428		Tunisia	11.644793
Italy	4.760962	Italy	5.675233		Vietnam	3.165710
Germany	4.507946	Germany	3.920258		Algeria	2.065967
Tunisia	4.419083	Spain	1.839538		Egypt	1.954555
Spain	4.150749	Poland	1.600657		Montenegro	1.025746

For the country graphs we also decided to implement the code to generate the geographical heatmaps by using the library 'geopandas', so that we can obtain more representative results. The geographical maps were created by only considering Europe and excluding the country Russia. An example of these is the following map obtained with the PageRank results



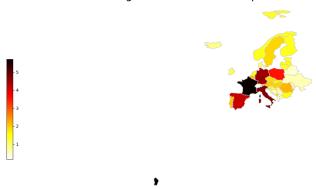


Figure 1: PageRank Heatmap

From the results we obtained it is possible to see that there are some differences between the two methods, more in detail in the top results for PageRank and Authorities values there are almost the same countries(France,Italy,Germany,Spain) instead the Hubs values are completely different.

Universities Results

In the case of the results obtained with the UniNodes graph, we can observe that the PageRank and Hits values have different rankings for top Universities, except for "UNIVERSIDADE DE LISBOA" which is the most important for both PageRank and Authorities. From the table containing the Hubs scores it is clear that the most relevant Universities in Italy are also the ones that have students traveling to the universities that have high Authorities values. In particular "UNIVERSITA DEGLI STUDI DI ROMA" and "POLITECNICO DI MILANO" have also the top PageRank values.

Organization	PageRank
UNIVERSIDADE DE LISBOA	49.156703
Stichting ArtEZ	45.381767
UNIVERSITA DEGLI STUDI DI ROMA LA SAPIENZA	44.403827
Eesti Kunstiakadeemia	42.335885
POLITECNICO DI MILANO	41.404278

Organization	Authorities
UNIVERSIDADE DE LISBOA	8.715438e+01
NORGES TEKNISK-NATURVITENSKAPELIGE UNIVERSITET	$6.628848e{+01}$
UNIVERSITAT POLITECNICA DE CATALUNYA	6.617643e + 01
UNIVERSITAT POLITECNICA DE VALENCIA	6.022172e+01
KATHOLIEKE UNIVERSITEIT LEUVEN	5.544156e + 01

Organization	Hubs
ALMA MATER STUDIORUM - UNIVERSITA DI BOLOGNA	9.454380e+01
UNIVERSITA DEGLI STUDI DI PADOVA	7.960611e+01
UNIVERSITA DEGLI STUDI DI ROMA LA SAPIENZA	7.632965e+01
POLITECNICO DI MILANO	7.479510e + 01
UNIVERSIDADE DE LISBOA	6.127300e+01

5 Further analysis and improvements

- 1. Use of the **Random graphs** method to verify if the features we extracted actually give out interesting information;
- 2.