

# Computing top-k Closeness Centrality Faster in Unweighted Graphs

Luca Lombardo

**Abstract**

TO DO!

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	The Problem . . . . .	3
<b>2</b>	<b>The algorithm</b>	<b>5</b>
2.1	The lower bound technique . . . . .	5
<b>3</b>	<b>The IMDB Case Study</b>	<b>7</b>
3.1	Data Structure . . . . .	7
3.2	Filtering . . . . .	9
3.2.1	name.basics.tsv . . . . .	9
3.2.2	title.basics.tsv . . . . .	9
3.2.3	title.principals.tsv . . . . .	10

# 1 Introduction

A graph  $G = (V, E)$  is a pair of a sets. Where  $V = \{v_1, \dots, v_n\}$  is the set *nodes*, and  $E \subseteq V \times V$ ,  $E = \{(v_i, v_j), \dots\}$  is the set of *edges* (with  $|E| = m \leq n^2$ ).

In this paper we discuss the problem of identifying the most central nodes in a network using the measure of *closeness centrality*. Given a connected graph, the closeness centrality of a node  $v \in V$  is defined as the reciprocal of the sum of the length of the shortest paths between the node and all other nodes in the graph. Normalizing we obtain the following formula:

$$c(v) = \frac{n - 1}{\sum_{w \in V} d(v, w)} \quad (1)$$

where  $n$  is the cardinality of  $V$  and  $d(v, w)$  is the distance between  $v, w \in V$ . This is a very powerful tool for the analysis of a network: it ranks each node telling us the most efficient ones in spreading information through all the other nodes in the graph. As mentioned before, the denominator of this definition give us the length of the shortest path between two nodes. This means that for a node to be central, the average number of links needed to reach another node has to be low. The goal of this paper is to computer the  $k$  vertices with the higher closeness centrality.

As case study we are using the collaboration graph of the actors in the *Internet Movie Database* (IMDB). On this data we define an undirected graph  $G = (V, E)$  where

- the vertex  $V$  are the actor and the actress
- the non oriented edges in  $E$  links the actors and the actresses if they played together in a movie.

## 1.1 The Problem

We are dealing with a web-scale network: any brute force algorithm would require years to end. The main difficulty is caused by the computation of distance  $d(v, w)$ . This is a well know problem: *All Pairs Shortest Paths or APSP problem*.

We can solve the APSP problem either using the fast matrix multiplication or, as I did, implementing a breath-first-search (BFS) method. There are several reason to prefer this second approach over the first one in this type of problems.

A graph is a data structure and we can describe it in different ways. Choosing one over another can have an enormous impact on performance. In this case, we need to remember the type of graph that we are dealing with: a very big and

sparse one. The fast matrix multiplication requires to consider our graph as an  $n \times n$  matrix where the position  $(i, j)$  is zero if the nodes  $i, j$  are not linked, 1 (or a generic number if weighted) otherwise. This method requires  $O(n^2)$  space in memory, that is an enormous quantity on a web-scale graph. Furthermore the time complexity is  $O(n^{2.373} \log n)$  [Zwick 2002; Williams 2012]

Using the BFS method the space complexity is  $O(n + m)$ , which is a very lower value compared to the previous method. In terms of time, the complexity is  $O(nm)$ . Unfortunately, this is not enough to compute all the distances in a reasonable time. It is also been proven that this method can not be improved. In this paper I will propose an exact algorithm to compute the top- $k$  nodes with the higher closeness centrality. I will also discuss an interesting and original relation between the physics of the visualized graph and the nodes with different centrality values.

## 2 The algorithm

In a connected graph, given a node  $v \in V$ , we can define the its farness as

$$f(v) = \frac{1}{c(v)} = \frac{1}{n-1} \sum_{w \in V} d(v, w) \quad (2)$$

where  $c(v)$  is the closeness centrality defined in (1). Since we are working with a disconnected graph, a natural generalization of this formula is

$$f(v) = \frac{1}{c(v)} = \frac{1}{r(v)-1} \sum_{w \in V} d(v, w) \quad (3)$$

where  $r(v) = |R(v)|$  is the cardinality of the set of reachable nodes from  $v$ . To avoid any problem during the computation, this formula still needs to be modified. Let's assume the nodes  $v$  that we are considering has just a link at distance 1 with another node  $w$  with *out-degree* 0. If we consider the formula (3) we will get a false result:  $v$  would appear to be very central, even if it's obviously very peripheral. To avoid this problem, we can generalize the formula (3) normalizing as suggested in [Lin 1976; Wasserman and Faust 1994; Boldi and Vigna 2013; 2014; Olsen et al. 2014]

$$f(v) = \frac{n-1}{(r(v)-1)^2} \sum_{w \in R(v)} d(v, w) \quad (4)$$

With the convention that is a case of  $\frac{0}{0}$  we set the closeness of  $v$  to 0

### 2.1 The lower bound technique

During the computation of the farness, for each node, we have to compute the distance from that node and all the other one reachable from it. Since we are dealing with millions of nodes, it's not possible in a reasonable time. In order to compute only the top- $k$  most central node we need to find a way to avoid computing BFS for nodes that won't be in the top- $k$ .

The idea is to keep track of a lower bound on the farness for each node that we will compute. This will allow us to kill the BFS operation before reaches the end if the lower bound tell us that the node will not be in the top- $k$ . More precisely:

- The algorithm will compute the farness of the first  $k$  nodes, saving them in a vector `top-actors`. From now on, this vector will be full
- Then, for all the next vertices, it defines a lower bound

$$\frac{n-1}{(n-1)^2} (\sigma_{d-1} + n_d \cdot d) \quad (5)$$

where  $\sigma_d$  is the partial sum in (4) at the level of exploration  $d$ . The lower bound (5) is updated every time that we change level of exploration during the BFS. In this way, if at a change of level the lower bound of the vertex that we are considering is bigger than the  $k - th$  element of **top-actors**, we can kill the BFS. The reason behind that is very simple: the vector **top-actors** is populated with the top-k nodes in order and the farness is inversely proportional to the closeness centrality. So if at that level the lower bound is already bigger than the last element of the vector, there is no need to compute the other level of the BFS since it will not be added in **top-actors** anyway.

The (5) it's a worst case scenario, and makes it perfect for a lower bound. If we are at the level  $d$  of exploration, we have already computed the sum in (4) up to the level  $d - 1$ . Then we need consider in our computation of the sum the current level of exploration: the worst case gives us that it's linked to all the nodes at distance  $d$ . We also put  $r(v) = n$ , in the case that our graph is strongly connected and all vertices are reachable from  $v$

SCRIVERE PSEUDOCODICE

## 3 The IMDB Case Study

The algorithm shown before can be applied to any dataset on which is possible to build a graph on. In this case we are considering the data taken from the *Internet Movie Database* (IMDB).

### 3.1 Data Structure

All the data used can be downloaded here: <https://datasets.imdbws.com/>

In particular we're interested in 3 files

- `title.basics.tsv`
- `title.principals.tsv`
- `name.basics.tsv`

Let's have a closer look to this 3 files:

#### **title.basics.tsv.gz**

*Contains the following information for titles:*

- `tconst` (string) - alphanumeric unique identifier of the title
- `titleType` (string) - the type/format of the title (e.g. movie, short, tvseries, tvepisode, video, etc)
- `primaryTitle` (string) - the more popular title / the title used by the filmmakers on promotional materials at the point of release
- `originalTitle` (string) - original title, in the original language
- `isAdult` (boolean) - 0: non-adult title; 1: adult title
- `startYear` (YYYY) - represents the release year of a title. In the case of TV Series, it is the series start year
- `endYear` (YYYY) - TV Series end year.
- `runtimeMinutes` - primary runtime of the title, in minutes
- `genres` (string array) - includes up to three genres associated with the title

### **title.principals.tsv.gz**

*Contains the principal cast/crew for titles:*

- **tconst** (string) - alphanumeric unique identifier of the title
- **ordering** (integer) – a number to uniquely identify rows for a given titleId
- **nconst** (string) - alphanumeric unique identifier of the name/person
- **category** (string) - the category of job that person was in
- **job** (string) - the specific job title if applicable
- **characters** (string) - the name of the character played if applicable

### **name.basics.tsv.gz**

*Contains the following information for names:*

- **nconst** (string) - alphanumeric unique identifier of the name/person
- **primaryName** (string)– name by which the person is most often credited
- **birthYear** – in YYYY format
- **deathYear** – in YYYY format if applicable
- **primaryProfession** (array of strings)– the top-3 professions of the person
- **knownForTitles** (array of tconsts) – titles the person is known for



## 3.2 Filtering

This is a crucial section for the algorithm in this particular case study. This raw data contains a huge amount of un-useful information that will just have a negative impact on the performance during the computation. We are going to see in detail all the modification made for each file. All this operation have been implemented using `python` and the `pandas` library.

### 3.2.1 name.basics.tsv

For this file we only need the following columns

- `nconst`
- `primaryTitle`
- `primaryProfession`

Since all the actors starts with the string `nm0` we can remove it to clean the output. Furthermore a lot of actors/actresses do more than one job (director etc..). To avoid excluding important actors we consider all the ones that have the string `actor/actress` in their profession. In this way, both someone who is classified as `actor` or as `actor, director` is taken into consideration.

Then we can generate the final filtered file `Attori.txt` that has only two columns: `nconst` and `primaryName`

### 3.2.2 title.basics.tsv

For this file we only need the following columns

- `tconst`
- `primaryTitle`
- `isAdult`
- `titleType`

Since all the movies starts with the string `t0` we can remove it to clean the output. In this case, we also want to remove all the movies for adults. This part can be optional if we are interest only in the closeness and harmonic centrality. Even if the actors and actresses of the adult industry use to make a lot of movies together, this won't alter the centrality result. As we know, an higher closeness centrality can be seen as the ability of a node to spread efficiently information in the network. Including the adult industry would lead to the creation of a very dense and isolated neighborhood. But none of those nodes will have an higher closeness centrality because they only spread information in their community. This phenomenon will be discussed more deeply in the analysis of the graph

visualized.

We can also notice that there is a lot of *junk* in IMDb. To avoid dealing with un-useful data, we are considering all the non-adult movies in this whitelist

- `movie`
- `tvSeries`
- `tvMovie`
- `tvMiniSeries`

The reason to only consider this categories is purely to optimize the performance during the computation. On IMDb each episode is listed as a single element: to remove them without losing the most important relations, we only consider the category `tvSeries`. This category list a TV-Series as a single element, not divided in multiple episodes. In this way we will lose some of the relations with minor actors that may appear in just a few episodes. But we will have preserved the relations between the protagonists of the show.

Then we can generate the final filtered file `FilmFiltrati.txt` that has only two columns: `tconst` and `primaryTitle`

### 3.2.3 `title.principals.tsv`

For this file we only need the following columns

- `tconst`
- `nconst`
- `category`

As before, we clean the output removing unnecessary strings. Then we create an array of unique actor ids (`nconst`) and an array of how many times they appear (`counts`). This will give us the number of movies they appear in. And here it comes the core of this filtering.

Let's define a constant `MINMOVIES`. This integer is the minimum number of movies that an actor needs to have made in his carrier to be considered in this graph. The reason to do that it's purely computational. If an actor/actress has less than a reasonable number of movies made in his carrier, there is a high probability that he/she has an important role in our graph during the computation of the centralities.