

# A Social Network Analysis based approach to deriving knowledge about research scenarios in a set of countries<sup>\*</sup>

Paolo Lo Giudice<sup>1</sup>, Paolo Russo<sup>1</sup>, and Domenico Ursino<sup>2</sup>

<sup>1</sup> DIIES, University “Mediterranea” of Reggio Calabria

<sup>2</sup> DICEAM, University “Mediterranea” of Reggio Calabria

(DISCUSSION PAPER)

**Abstract.** In this paper, we propose a new Social Network Analysis based approach to providing a multi-dimensional picture of the research scenarios of a set of countries of interest and to detecting possible hubs operating therein. This knowledge allows the understanding of the impact of different socio-economic conditions on research. Furthermore, it may support the design of policies for sustaining the accumulation of scientific and technological capabilities. We apply our approach to four North African countries (i.e., Algeria, Egypt, Morocco and Tunisia) in such a way as to show its potential.

## 1 Introduction

In the last years, scientometrics and bibliometrics received a growing interest both in research literature and as objective ways for evaluating the performances of researchers, universities, institutions, etc. Data available for scientometrics and bibliometrics investigations are growing at a very rapid rate. As a matter of fact, currently, the problem of extracting useful knowledge from the large amount of available scientometrics and bibliometrics data can be seen as a Data Mining problem, and in the very next future, big data approaches for solving it will be unavoidable. The obvious consequence of this reasoning is that more and more innovative approaches to facing this issue are necessary.

Social Network Analysis [12, 8] and, more in general, graph theory, have been a prominent family of approaches adopted in the past in this context (see, for instance [10, 4, 5, 11, 2, 9, 7, 6, 3]). Furthermore, it is possible to foresee that they will be even more exploited in the future, due to the more and more increasing number of proposals somehow involving them.

This paper aims at providing a contribution in this setting. Indeed, it proposes a new Social Network Analysis based approach to deriving knowledge about research scenarios and hubs in a set of countries of interest. As for this paper, a hub is a research institution that operates as a guide or stimulus to the research in its country and, at the same time, is capable of stimulating cooperations with institutions of other countries. Our hub definition is strongly fitted to

---

<sup>\*</sup> This work was partially supported by Aubay Italia S.p.A.

our scenario of interest. It does not claim to have a mathematical foundation, but it strongly benefits from the observations, suggestions and experience of innovation management researchers, who guided us in its formulation. Our approach is general and can be directly applied to any set of countries. The only requirement is to have at disposal the set of the publications of all the research institutions of the countries to investigate. In this paper, we applied it to four North African countries (e.g., Algeria, Egypt, Morocco and Tunisia) and we exploited all the publications of all the research institutions of the four countries of interest in the time interval [2003, 2013], as stored in the Web of Science repository [1]. The most important support data structure is a social network whose nodes represent institutions and whose edges denote collaborations among institutions. Starting from it, other important support data structures and accompanying parameters (some of which were never defined in the literature) are introduced.

This paper is organized as follows. In Section 2, we illustrate our approach. In Section 3, we apply it to the four North African countries mentioned above. Finally, in Section 4, we draw our conclusions and overview some possible future developments.

## 2 Description of our approach

Before starting the description of our approach, we must define some sets that formalize available data and, therefore, will be extensively used below. The first set regards the set  $RA$  of research areas. It consists of the following elements:  $RA = \{\text{'NS'}, \text{'AS'}, \text{'MH'}, \text{'SS'}, \text{'HU'}, \text{'ET'}\}$ , where ‘NS’ (resp., ‘AS’, ‘MH’, ‘SS’, ‘HU’, ‘ET’) stands for ‘Natural Science’ (resp., ‘Agricultural Science’, ‘Medical and Health Science’, ‘Social Science’, ‘Humanities’, ‘Engineering and Technology’). The second set concerns the overall set  $Pub$  of publications at our disposal. Given a publication  $p \in Pub$ . The third basic set regards the set  $C$  of the countries to investigate.

### 2.1 Hub characterization and detection

In this section, we aim at detecting a method for detecting both hubs and their features in a set of countries. For this purpose, we preliminarily introduce a first support data structure. It is a social network:  $G = \langle N, E \rangle$ .  $N$  is the set of the nodes of  $G$ . A node  $n_i \in N$  corresponds to exactly one institution registered in our database. Since there is a biunivocal correspondence between a node of  $N$  and the corresponding institution, in the following, we will use the symbol  $n_i$  to indicate both of them. Each node of  $N$  is labeled with an element of  $C$  depending on the country of the corresponding institution. We indicate by  $l_i$  the label of  $n_i$ .  $E$  is the set of the edges of  $G$ . There exists an edge  $e_{ij} = (n_i, n_j, w_{ij}) \in E$  if there exists at least one publication involving one author of  $n_i$  and one author of  $n_j$ .  $w_{ij}$  is the weight of  $e_{ij}$ ; it denotes the number of publications having at least one researcher of  $n_i$  and one researcher of  $n_j$  among their authors.

Now, we are able to introduce the concept of hub. With regard to this fact, we point out that we do not aim at proposing a new concept characterized by a mathematical foundation supporting it. Instead, we would like to introduce an informal and empirical, yet reasonable, concept, which can support innovation managers to make their decisions. In carrying out this activity, we strongly benefited from the observations, suggestions and experience of innovation management researchers, who guided us in its formulation. Taking this purpose into account, we can say that, in order to be a hub, an institution must satisfy the following conditions: *(i)*  $C_1$ : it should have published a very high number of papers; *(ii)*  $C_2$ : it should have published a high number of papers in cooperation with institutions different from the ones of its country; *(iii)*  $C_3$ : it should have published many papers in cooperation with institutions of its country.

The reasons underlying these three conditions are the following: *(i)* If an institution published very few papers, even if all in co-authorship with foreign institutions, it cannot have a weight such as to influence the research scenario of its country. *(ii)* If an institution published a high number of papers, but all in co-authorship with other institutions of its own country only, it would be certainly an important research center in the context of its country, but it would not have the capability, required to hubs, of stimulating contacts with foreign countries. *(iii)* If an institution published even a lot of papers, but all in co-authorship with foreign institutions only, it would not be able to strongly influence the research of its own country.

To “quantify” conditions  $C_1$ ,  $C_2$  and  $C_3$ , we use three metrics, namely  $M_1$ ,  $M_2$  and  $M_3$ , respectively.  $M_1$  coincides with the classical weighted degree centrality,  $M_2$  coincides with the normalized weighted degree centrality and  $M_3$  is analogous to the E-I index [8].

As theoretically conjectured in the past literature, and as verified for the countries composing our case study,  $M_1$ ,  $M_2$  and  $M_3$  follow a power law distribution. Taking all these considerations into account, the set  $\mathcal{H}^X$  of hubs for the countries into consideration can be defined as the set of those institutions simultaneously belonging to the top  $X\%$  of institutions with the highest values of  $M_1$ ,  $M_2$  and  $M_3$  (we call  $I_1^X$ ,  $I_2^X$  and  $I_3^X$  these three sets, when considered separately). In this definition,  $X$  is a threshold allowing the selection of the institutions having the highest values of  $M_1$ ,  $M_2$  and  $M_3$ . The choice to use  $X$  as a threshold parameter derives from the power law distributions characterizing all the three metrics. Reasonable values of  $X$  could be 10, 15 and 20. After several experiments, we decided to consider a default value of  $X$  equal to 20. As a consequence, in the following, when  $X$  is not specified, we intend that it is equal to 20. In the following, we use the symbol  $\mathcal{H}_k^X$  to indicate the hubs of a given country  $k$ .

## 2.2 Investigation of the research scenarios for the countries of interest

In this section, we aim at analyzing the research scenarios of the countries into examination. Initially, we can introduce three indicators that could give us some

knowledge about the research scenarios of the countries into consideration. The first one,  $RQ$ , is an indicator of the overall research quality in the countries of interest. In fact, it measures how many institutions of  $I_1$  belong to these countries. The second one,  $FC$ , indicates how many institutions, among the top ones of the countries of interest, publish many papers with foreign institutions. The third one,  $TP$ , indicates how many institutions that publish very much with foreign institutions belong to the top institutions of the countries of interest.

In the investigation of the research scenario of a country  $k$  and of the role of its hubs, it appears very interesting to study its paper distribution. For this purpose, we introduce the average number  $AvgPub_k^{\mathcal{H}}$  of the publications of its hubs. Another interesting issue to investigate is to verify if a hub of  $k$  publishes more with institutions of  $k$  (we call “internal” the corresponding publications) than with foreign ones (we call “external” the corresponding publications) or alone. To carry out this investigation, we introduce: (i) the average number  $AvgHubPub_k^I$  of publications performed by the hubs of  $k$  with institutions of  $k$ ; (ii) the average number  $AvgHubPub_k^F$  of publications performed by the hubs of  $k$  with foreign institutions; (iii) the average number  $AvgHubPub_k^A$  of publications performed alone by the hubs of  $k$  (we call them “alone publications” in the following).

A further interesting analysis is devoted to understand if, in their cooperation with foreign institutions, the hubs of a given country  $k$  privilege one or few countries. For this purpose, we specialize to our research context the Herfindahl Index. This index is very used in economics is defined as the sum of the squares of the market shares of the firms within the industry, where market shares are expressed as fractions. It can range from 0.0 to 1.0, moving from a huge number of very small firms to a single monopolistic producer. In our case, we extend the Herfindahl index to our context and define the Herfindahl Index  $HI_k$  associated with the papers published by the hubs of  $k$  to verify if these hubs published in cooperation with institutions of few (implying high values of  $HI_k$ ) or many (implying low values of  $HI_k$ ) countries.

### 2.3 Cooperation among hubs of the same country

In this section, we aim at investigating the cooperation levels of the hubs of a given country  $k$ . For this purpose, we preliminarily define a support data structure called *clique social network*. In particular, let  $G$  be the social network defined in Section 2.1 and let  $G_k$  be its “projection” on the country  $k$ . Let  $\mathcal{C}_k$  be the set of cliques of  $G_k$  and let  $\mathcal{H}_k$  be the set of the hubs of  $k$ . A *clique social network*  $CG_k$  has a node for each hub of  $\mathcal{H}_k$  belonging to at least one clique of  $\mathcal{C}_k$ . Each node  $n_i$  of  $CG_k$  has associated a weight  $w_i$  denoting the number of cliques of  $\mathcal{C}_k$  which it belongs to. An edge  $(n_i, n_j)$  of  $CG_k$  denotes that  $n_i$  and  $n_j$  together belong to at least one clique of  $\mathcal{C}_k$ .

Some measures capable of quantitatively representing the differences that characterize the cooperation among hubs are the following: (i) the number of cliques  $|\mathcal{C}_k|$ ; (ii) the absolute dimension  $d_{\mathcal{C}_k}$  of the largest clique in  $\mathcal{C}_k$ ; (iii) the relative dimension  $\frac{d_{\mathcal{C}_k}}{|\mathcal{H}_k|}$  of the largest clique in  $\mathcal{C}_k$ ; (iv) the fraction  $f_{\mathcal{C}_k}^{\mathcal{H}}$  of hubs

belonging to at least one clique of  $\mathcal{C}_k$ . In order to avoid that results are biased by the number of publications (which can be very different in the different countries of interest), we define a normalized version  $\widehat{CG}_k$  of  $CG_k$ . Finally, we searched for some measures to compare clique social networks. After several experiments, we found that the most significant ones were: (i) the number of nodes; (ii) the number of edges; (iii) density<sup>3</sup>.

## 2.4 Investigation about the quality of publications

All indicators introduced above are based only on the number of publications. Actually, it would be important to take also their quality into account. One way to do this consists in taking their impact factor into consideration; another way consists in considering the number of citations received by papers. Impact factors are measured only for journal papers. As a consequence, if we want to exploit this measure, we must define a new support data structure. This structure, that we indicate by  $G'$ , is, once again, a social network. It is defined as  $G' = \langle N', E' \rangle$ . There is a node  $n_i \in N'$  for each institution having at least one author that published at least one journal paper. An edge  $e'_{ij} = (n'_i, n'_j, w'_{ij})$  has a semantics similar to the one of  $e_{ij}$  except that the weight  $w'_{ij} = \sum_{p \in (Pub_{ij} \cap JPub)} IF_p$  considers both the number of publications simultaneously performed by  $n_i$  and  $n_j$  and the corresponding impact factors. Paper citations are valid both for conference proceedings and for journal papers. However, in order to make our analyses about the quality of publications homogeneous, we chose to investigate only journal papers. In this case, we used the same support social network as the one exploited for impact factors but the edge weights  $w'_{ij}$  was computed as:  $w'_{ij} = \sum_{p \in (Pub_{ij} \cap JPub)} CitN_p$ , where  $CitN_p$  is the number of citations of  $p$ .

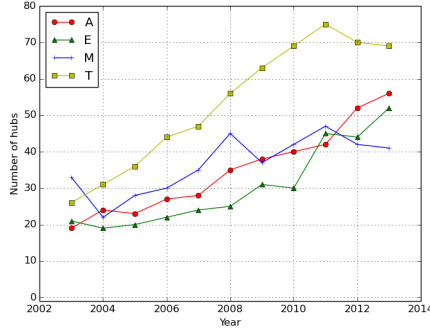
## 2.5 Characterization of hub neighborhoods

A first parameter useful to characterize hub neighbors is the average number  $AvgPub$  of publications of the hub neighborhoods. A second parameter regards their average dimension  $AvgDim$ . Even in this case, we disaggregate data per country and we call  $AvgDim_k$  the corresponding parameter for the country  $k$ .

A next analysis regards the cooperation level among the institutions belonging to hub neighborhoods. To perform this task, we define a new support social network. We call it *nbh social network* and we represent it by means of the symbol  $NbhG_i$ . Given a neighborhood  $nbh_i$ , the corresponding *nbh* social network is defined as follows:  $nbhG_i = \langle nbh_i, nbhE_i \rangle$ . There is a node in  $NbhG_i$  for each node of  $nbh_i$ ; there is an edge  $(n_i, n_j) \in nbhE_i$  if there exists at least one publication between an author of  $n_i$  and an author of  $n_j$ .

After having introduced this social network, we define a first parameter on it. This parameter is called  $AvgCFrac$  and corresponds to the average fraction of the real number of cliques existing in hub neighborhoods against the possible

<sup>3</sup> Actually, this last measure can be derived from the two other ones, but it is very expressing and, consequently, we decided to explicitly consider it.



**Fig. 1.** Number of hubs for each country in the year interval [2003,2013]

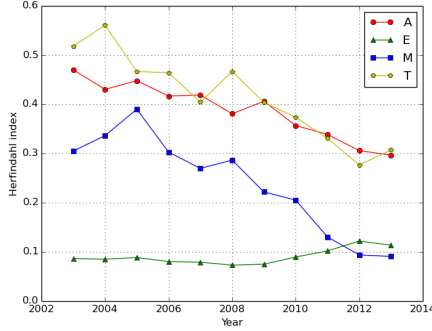
number of them. It is an indicator of the cooperation level among hubs. As usual, we call  $AvgCFrac_k$  the “projection” of  $AvgCFrac$  on the country  $k$ . A second parameter about intra-neighborhood cooperation regards the average fraction  $AvgCNbh$  of the number of cliques existing in hub neighborhoods against the number of neighborhood nodes. Again, we call  $AvgCNbh_k$  the “projection” of  $AvgCNbh$  on the country  $k$ . A final parameter measuring the cooperation level between hub neighbors is the average density  $AvgDens$  of the  $nbh$  social network. As usual, we call  $AvgDens_k$  the “projection” of  $AvgDens$  on the country  $k$ .

### 3 Application of our approach to four North African countries

As pointed out in the introduction, we applied our approach to four North African countries, namely Algeria, Egypt, Morocco and Tunisia. As a consequence, in our case study, the set  $C$  introduced in Section 2, consisted of the following elements:  $C = \{‘A’, ‘E’, ‘M’, ‘T’, ‘O’\}$  where ‘A’ (resp., ‘E’, ‘M’, ‘T’, ‘O’) stands for ‘Algeria’ (resp., ‘Egypt’, ‘Morocco’, ‘Tunisia’, ‘Others’). Clearly, ‘O’ indicates all the countries different from the four into examination. The reasons for adding ‘O’ will be clear below. Due to space limitations we can present only a very limited number of the results that we have obtained.

In Figure 1, we report the variation of the number of hubs for each country. From the analysis of this figure, we can see that the country with the highest number of hubs is Tunisia. This result was unexpected also because both the extension and the number of citizens of Tunisia were smaller than the ones of the other three countries.

In Figure 2, we report the Herfindahl index  $HI_k$  for the four countries. From the analysis of this figure we can observe that Tunisia and Algeria have a high Herfindahl index, which implies that their hubs cooperate mostly with one or few countries. By contrast, Egypt has a very low Herfindahl index, i.e., its hubs cooperate with many countries. An interesting trend is the one of Morocco; in



**Fig. 2.** Herfindahl index over time for the four countries

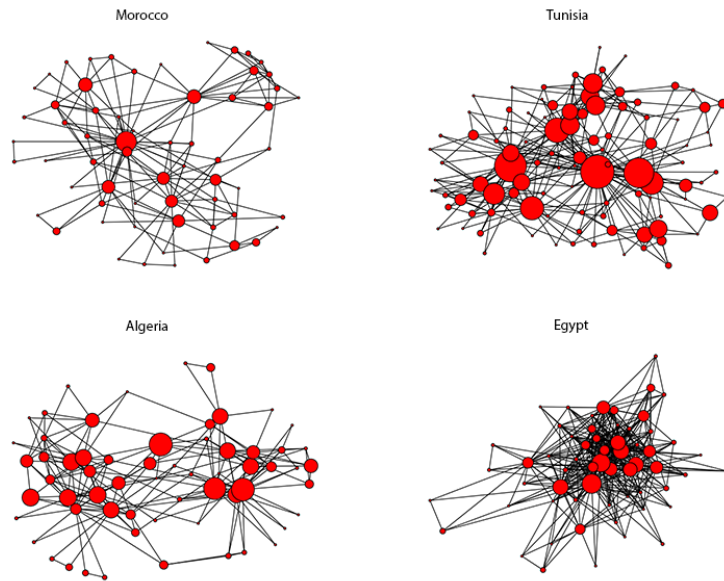
Country	$ C1_k $	$d1_{C_k}$	$\frac{d1_{C_k}}{ H_k }$	$f1_{C_k}^H$
Algeria	292	7	0.152	0.913
Egypt	38	13	0.351	0.973
Tunisia	130	8	0.116	0.942
Morocco	82	7	0.127	0.818

**Table 1.** Quantitative differences characterizing the cooperation behaviors of hubs in the four countries in the time interval [2003, 2009]

fact, it initially has a behavior like the ones of Tunisia and Algeria, whereas, in the last years, it shows a behavior like the one of Egypt.

To determine the cooperation levels among hubs for the four North African countries into consideration, for each country  $k$ , we performed the following tasks: (i) we considered the two time intervals [2003, 2009] and [2007, 2013]; (ii) we computed the clique social networks  $CG1_k$  (resp.,  $CG2_k$ ), corresponding to the first and the second time intervals, respectively; (iii) we measured the four parameters introduced in Section 2.2 for quantitatively evaluating clique social networks. Obtained results for the first time interval are reported in Table 1. From the analysis of these tables we can draw the following conclusions: (i) Egypt has the largest clique in both periods; the clique is much larger than the maximum cliques of the other countries; (ii) in Egypt almost all hubs belong to at least one clique. These results indicate that Egyptian hubs are more prone to cooperation than the hubs of the other countries.

In Figure 3, we report the graphs  $CG2_k$  for all the four countries; in these graphs the dimension of nodes is proportional to the corresponding weight, i.e., to the number of cliques they belong to. The analysis of this figure confirms the previous conjecture; in fact, the number of edges in the Egyptian graph is much higher than in the other graphs. This fact, along with the presence of many not very large nodes, allows us to derive another important knowledge pattern, i.e., that research cooperation in Egypt is more advanced than in the other countries.



**Fig. 3.** Graphs  $CG2_k$  for all the four countries

## 4 Conclusion

In this paper, we have proposed a new SNA-based approach to investigating the research scenarios of a set of countries of interest and to detecting possible hubs operating in these countries. Extracted knowledge allows the evaluation of the impact of different socio-economic conditions on research and favors the design of policies for supporting innovation in the countries of interest. We applied our approach to four North African countries. In the future, we plan to exploit analysis techniques about information diffusion in social networks to understand how the possible mobility of top researchers from one institution to another can impact on the quality of this latter. Furthermore, we plan to investigate the possible application of classification techniques to derive hub profile in different countries. Finally, we plan to analyze the possible application of prediction techniques to understand what kind of financial investment must be performed for maximizing the increase of both the number and the quality of hubs and publications in the countries of interest.

## References

1. Web Of Science. <http://wokinfo.com/>, 2015.
2. A. Abbasi, J. Altmann, and L. Hossain. Identifying the effects of co-authorship networks on the performance of scholars: A correlation and regression analysis of



- performance measures and social network analysis measures. *Journal of Informetrics*, 5 (4):594–607, 2011.
3. A. Abbasi, L. Hossain, and L. Leydesdorff. Betweenness centrality as a driver of preferential attachment in the evolution of research collaboration networks. *Journal of Informetrics*, 6 (3):403–412, 2012.
  4. T. Arif, R. Ali, and M. Asger. Scientific co-authorship social networks: A case study of computer science scenario in India. *Science*, 52 (12):38–45, 2012.
  5. K. Badar, J.M. Hite, and Y.F. Badir. Examining the relationship of co-authorship network centrality and gender on academic research performance: the case of chemistry researchers in Pakistan. *Scientometrics*, 94 (2):755–775, 2013. Elsevier.
  6. M. Bordons, J. Aparicio, B. González-Albo, and A.A. Díaz-Faes. The relationship between the research performance of scientists and their position in co-authorship networks in three fields. *Journal of Informetrics*, 9 (1):135–144, 2015.
  7. Z. Chinchilla-Rodriguez, A. Ferligoj, S. Miguel, L. Kronegger, and F. de Moya-Anegón. Blockmodeling of co-authorship networks in library and information science in Argentina: a case study. *Scientometrics*, 93 (3):699–717, 2012.
  8. R. Hanneman and M. Riddle. *Introduction to social network methods*. <http://faculty.ucr.edu/~hanneman/nettext/>, 2005. University of California, Riverside.
  9. J. Kim and C. Perez. Co-authorship network analysis in industrial ecology research community. *Journal of Industrial Ecology*, 19 (2):222–235, 2015.
  10. P. Liu and H. Xia. Structure and evolution of co-authorship network in an interdisciplinary research field. *Scientometrics*, 103 (1):101–134, 2015.
  11. M. Pavlov and R. Ichise. Finding Experts by Link Prediction in Co-authorship Networks. In *Proc. of the International Workshop on Finding Experts on the Web with Semantics (FEWS 2007)*, pages 42–55, Busan, Korea, 2007.
  12. S. Wasserman and K. Faust. *Social network analysis: Methods and applications*. 1994. Cambridge University Press.