

Text Mining in practice: exploring patterns in text collections of remote work offers

Karolina Kuligowska, Mirosława Lasek

Department of Information Systems and Economic Analysis, Faculty of Economic Sciences,
University of Warsaw, Warsaw, Poland
{kkuligowska,mlasek}@wne.uw.edu.pl

Abstract. The aim of this paper is to give an insight in text mining techniques in the context of unstructured text collections of location independent job offers. In order to extract useful information, uncover interesting patterns and features of remote work, we analyze five most popular and most visited websites containing job offers. We examine clusters of remote job offers, keywords describing those clusters, as well as linkages between strongly associated terms describing mobile work offers. It is interesting to observe the maturity of text mining tools which broadened their applications to new research topics and became suitable to explore new phenomena.

Keywords: text mining, text analytics, clustering, concept linking, remote work, telecommuting

1 Introduction

Since most of human knowledge is maintained in a textual form [Lin 2001], thus text analytics methods and techniques are constantly being developed, and this development accelerated recently [Agrawal 2013; Mahesh 2010; Patel 2012; Ramanathan 2013]. Previously focused on calculations, nowadays the computer power more and more serves for text issues, such as computational linguistics, natural language processing and text mining. Words patterns, context recognition, and term linkages constitute subject of insightful analysis. With emerging text mining tools scientists discover interesting results in many fields by exploring vast amount of text data.

The era of innovative technology influences as well the way of working. It is not necessarily uncommon nowadays to use the Internet connection for working instead of going to the real physical office. Remote working evolved quickly and freelancing has become an increasingly dynamically developing field. Nevertheless, there is a widespread belief that telework offers are focused mainly on those who know computer graphics and have skills in programming. Is this a true conviction?

The aim of this paper is to explore unstructured text of remote work offers by applying text mining techniques. The paper is organized as follows. Section 2 presents data source and software applied in our analysis. Section 3 introduces studies about loca-

tion independent work and reviews briefly remote work features. Section 4 combines discussed issues in one analytical frame, presents the obtained results and describes captures of tables and figures. Section 5 provides a summary of our findings. Finally, the conclusions due to this paper are considered in Section 6.

2 Data source and software applied

The remote job offers examined in our text mining analysis originate from the following websites: Careerbuilder¹, Remoteemployment², Monster³, Jobamatic⁴, Simplyhired⁵. These international websites are several of the world's most widely recognized portals with job offers for remote and regular workers. However, there are many more job portals on the Internet, beyond websites used in our analysis.

With regard to the software tool, we performed text parsing, clustering and concept linking techniques using SAS Text Miner 4.2 software within SAS Enterprise Miner 6.2 environment. By executing macro %tmfilter we extracted text of job offers from each website. Subsequently, each set of unstructured text collection was turned into an adequate structured table and analyzed within SAS business analytics flexible framework.

3 Remote work

3.1 Telework features

Remote work, as an innovative form of employment, is perceived as an alternative for the traditional in-office work environment. In the last decade, companies have slowly shifted towards a virtual workplace, and employees quickly adapted to fulfill this business demand. As technological advances have provided mobile electronic media of communication, companies have realized the benefits of the virtual workplace trends and flexible work arrangements [Busch 2011; Lister 2011]. The survey conducted in 2011 by global research company Ipsos revealed that telecommuting is primarily taking place in emerging markets of Middle East, Africa, Latin America, and Asia-Pacific [Gottfried 2012]. Moreover, according to the Forrester Research forecasts, the telework will include 43% of US workers by 2016 [Schadler 2009]. Remote work is practiced globally and there is no doubt that the telecommuting trend steadily grows.

Mobile work executed on distance is described by numerous terms and synonyms such as: remote jobs, telework, telecommuting jobs, online jobs, home based work,

¹ <http://www.careerbuilder.com/Jobs/Keyword/Remote/> (April 2013)

² <http://www.remoteemployment.com> (April 2013)

³ <http://jobsearch.monster.com/search/?q=remote> (April 2013)

⁴ <http://momstowork.jobamatic.com/a/jobs/find-jobs/q-Remote> (April 2013)

⁵ <http://www.simplyhired.com/a/jobs/list/q-remote/fjt-telecommute> (April 2013)

flexible jobs, location independent work. Taxonomy researchers distinguish three main forms of telework: fixed-site telework, mobile telework, and flexiwork [Garrett 2007]. In a broad sense, a remote location worker is an employee who performs his or her work outside the workplace, communicating results of his or her work by means of electronic communication, thus eliminating distance restrictions or any problems associated with traditional commuting practices [Watad 2010; Fuhr 2011]. This type of work is quite well suited for freelancers such as graphic designers, computer programmers and interpreters. Furthermore, especially in the U.S., remote workers offer virtual assistant services, i.e. secretarial service, customer service and sales executed by phone.

3.2 Telework research

As the phenomenon of telework is constantly spreading around the world, there are more and more research studies done on this subject. They mainly focus on measuring: the effectiveness and productivity of remote work [Bloom 2013; Dutcher 2012; Teh 2011], activation of the people at risk of unemployment or socially excluded [Baker 2006; Stroińska 2012], and psychological impact of teleworking on human stress, emotions and health [Mann 2003; Ward 2001]. Nonetheless, there is little cross-sectional research concerning features of telework accessible for freelancers willing to work remotely. Therefore we conduct our text mining study to uncover descriptive keywords, explore patterns and common features that link various remote job offers. This approach adds new research area to widely known standard text mining application fields such as bioinformatics, business intelligence and customer relationship management [Gupta 2009; Jusoh 2012].

4 Text mining patterns exploration

We started our analysis by making five independent text mining models executed separately for each dataset. First, we extracted text of job offers from each website and turned it into an adequate input dataset. Then we extracted parts of speech and gathered structured information from unstructured text, in search of relevant terms hidden in the text database [Blansch e 2010; Kaur 2013]. For this purpose we subsequently performed three phases: text parsing, text filtering and text mining, using the appropriate nodes, as shown on Figure 1.

The Text Parsing node gathered the statistical data about the terms, such as number of terms, number of documents, and term frequencies in each dataset. In text parsing process we decided to focus only on one part of speech, namely adjectives. Adjectives describe nouns in terms of many qualities and therefore they constitute relevant elements giving more information than any other part of speech, especially in the case of remote work offers. The Text Parsing node allowed us to modify our output set of parsed terms by limiting our text parsing results to terms with a role of Adjective.

Next, we used Text Filter node in order to keep only sufficiently important terms and to remove irrelevant terms. Therefore we sorted the results for each set of text data by term weight in a descending mode. Terms with weight below 0.1. were rejected from our analysis.

Finally, we applied Text Mining node to examine the information that exists in each dataset. The Text Mining node allowed us to execute categorization, clustering, and keyword extraction, which are the text mining major tasks [Gharehchopogh 2012].

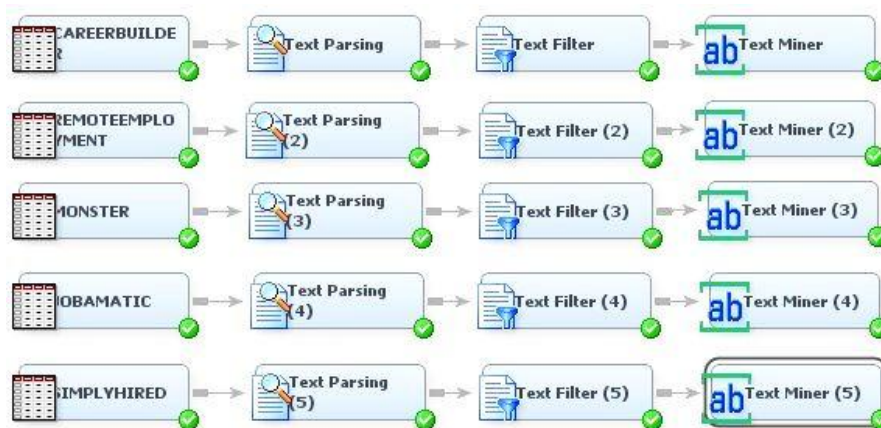


Fig. 1. Five text mining models. Source: own elaboration

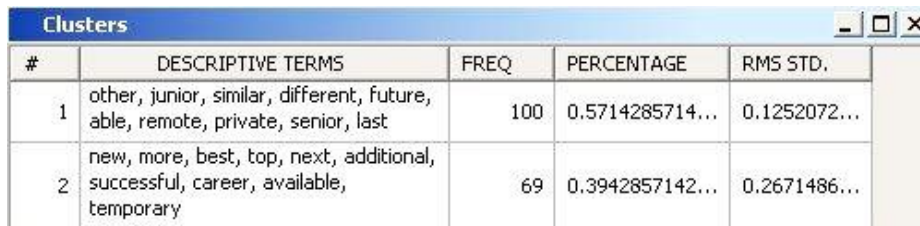
Clustering, defined as a "process of partitioning a dataset into clusters, so that elements of the same cluster are more similar to each other than to elements of different clusters" [Su 2010], constitutes a fundamental text mining method. It allows to achieve an organized overview of concepts contained in text documents and improves similar documents detection [Bolasco 2005; Vidhya 2010]. In order to find groups of alike offers in our text collection, we executed similarity based clustering. We applied the iterative Expectation Maximization algorithm, a clustering method which is most frequently used in a wide variety of applications [Aggarwal 2012]. Commonly applied Singular Value Decomposition technique [Radovanović 2008] was used for reducing the number of dimensions. Each cluster was described by 10 keywords, and we also allowed the unclustered outliers to be excluded from analysis.

Ultimately, for selected adjective with the highest weight we were eager to identify other terms that correlate the strongest with it. Therefore for each dataset we executed concept linking. This technique is used to find and present highly related terms, whereas the strength of association between these terms is measured by the chi-square statistic. Since visualization of related expressions helps to derive new insights and novel patterns in text collection [Don 2007; Gupta 2009], the linkage of correlated concepts is displayed in a form of hyperbolic tree graph with nodes that are expanded when necessary. Thus concept linking visualization enriches traditional searching methods with advanced browsing capabilities.

4.1 Careerbuilder

In Careerbuilder dataset of 175 remote work offers we identified 480 adjectives occurring in text. Three adjectives with the highest weight were: *liable* (0.927), *security* (0.921) and *advanced* (0.913).

In the result of applying Expectation-Maximization clustering algorithm, two clusters emerged. Each cluster is described by ten keywords (Descriptive terms), the number of documents in a cluster (Freq), the percentage of documents in a cluster (Percentage), and root mean squared standard deviation (RMS Std.), as presented on Figure 2. As it was stated before, the unclustered outliers were excluded from analysis. Therefore the total number of documents in clusters (169) is lower than the overall number of analyzed documents in dataset (175).



#	DESCRIPTIVE TERMS	FREQ	PERCENTAGE	RMS STD.
1	other, junior, similar, different, future, able, remote, private, senior, last	100	0.5714285714...	0.1252072...
2	new, more, best, top, next, additional, successful, career, available, temporary	69	0.3942857142...	0.2671486...

Fig. 2. Clusters sorted by Frequency. Source: own elaboration

The first cluster grouped 100 job offers described by such terms as: *other, junior, similar, different, future, able, remote, private, senior, last*. Most probably these offers refer to high school or college/university students and were aimed at 3rd, 4th or last year students able to work remotely in a private sector. The descriptive terms of second cluster, grouping 69 job offers, were: *new, more, best, top, next, additional, successful, career, available, temporary*. This may indicate a group of temporary or additional part-time jobs which may contribute to a successful career in future.

For one of the terms with the highest weight, i.e. *advanced*, we executed concept linking visualization, as presented on Figure 3.



Fig. 3. Concept linkages for the term advanced. Source: own elaboration

On the concept linking hyperbolic tree graph the term *advanced* is displayed in the center of the structure. We can observe that the term *advanced* is surrounded by the following adjectives: *new*, *best*, and *career*. The thickness of the line between concepts represents the strength of association, and a thicker line indicates a closer association. Therefore *advanced* is strongly correlated with *new* and *career*, as well as with *successful* and *additional* in expanded view mode.

4.2 Remoteemployment

In Remoteemployment database we identified 758 adjectives occurring in text of 231 remote work offers. Three adjectives with the highest weight were: *clinical* (0.955), *residential* (0.950) and *industrial* (0.936).

Once again we applied Expectation Maximization as a clustering method, and we obtained two clusters, as presented on Figure 4. The first cluster grouped 115 job offers and was described by quite irrelevant adjectives such as *both*, *unique*, *big*, *same*, *little*, *much*, *few*, *great*, *important*, and *better*. Whereas the second cluster gathered 113 documents and its descriptive terms related to contract duration, such as *part-time*, *full-time*, *interim*, *temporary*, among others.

#	DESCRIPTIVE TERMS	FREQ	PERCENTAGE	RMS STD.
1	both, unique, + big, same, + little, much, + few, + great, important, better	115	0.4978354978...	0.1667515...
2	part-time, full-time, interim, temporary, part, + fast, + large, no, successful, good	113	0.4891774891...	0.1509026...

Fig. 4. Clusters sorted by Frequency. Source: own elaboration

In case of Remoteemployment dataset, clusters look unclear and insignificant for further inference, therefore we continued to examine this dataset with concept linkage tool. For the term with the highest weight, namely *clinical*, we executed concept linking hyperbolic tree graph, as you can see on the following Figure 5.

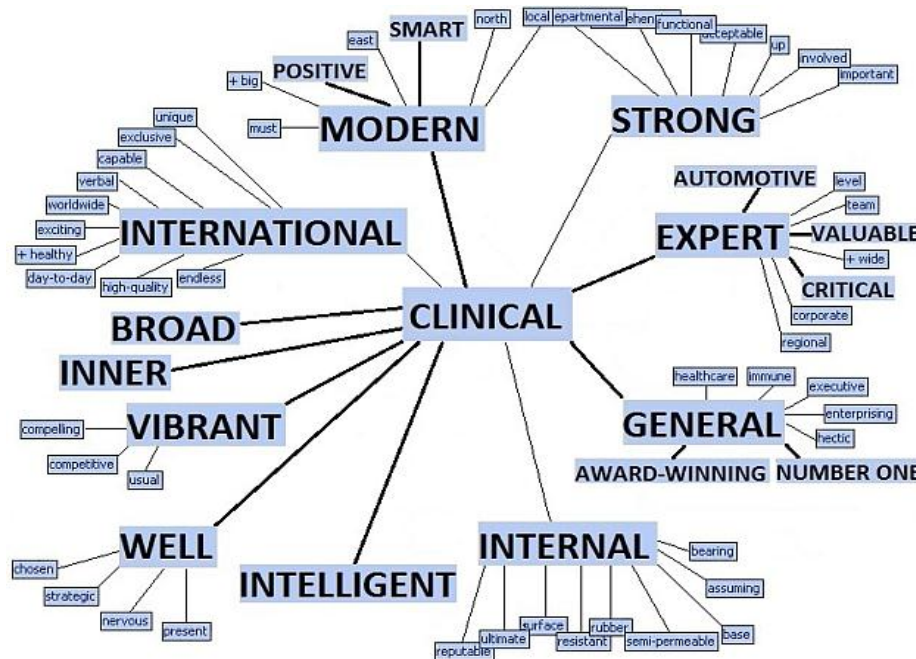


Fig. 5. Concept linkages for the term clinical. Source: own elaboration

The term *clinical* is surrounded by the following adjectives: *international*, *modern*, *strong*, *expert*, *general*, *internal*, *intelligent*, *well*, *vibrant*, *inner*, and *broad*. The term *clinical* is highly associated with *modern*, *expert*, *general*, *intelligent*, *well*, *vibrant*, *inner*, and *broad*. Further analysis of expanded nodes of terms *modern*, *expert* and *general* indicates close association with adjectives describing candidate's personal features required in work: *positive*, *smart*, *automotive*, *valuable*, *critical*, *award-winning*, *number one*.

4.3 Monster

In Monster database we identified 970 adjectives out of 179 remote work offers text. Adjectives with the highest weight were: *sole* (0.948), *automatic* (0.921), and *editorial* (0.904).

In result of clustering procedure, we obtained three clusters, as presented on Figure 6.

Clusters				
#	DESCRIPTIVE TERMS	FREQ	PERCENTAGE	RMS STD.
1	+ high, responsible, + large, excellent, + fast, common, other, daily, biweekly, monthly	74	0.4134078212...	0.1475900...
2	additional, human, real, legal, dental, free, retail, administrative, austin, top	63	0.3519553072...	0.0971502...
3	subject, invalid, marital, last, remote, veteran, similar, proud, first, united	18	0.1005586592...	0.0665655...

Fig. 6. Clusters sorted by Frequency. Source: own elaboration

The first cluster contained 74 job offers and included offers characterized by contract duration (*daily, biweekly, monthly*), as well as by candidate's features (*responsible, excellent, fast*). The second cluster gathered 63 offers and was described by terms concerning type of the job, such as *administrative, legal, dental, retail* and *additional*. Third cluster, with descriptive terms for 18 offers, turned out to be hard to interpret.

For the term *automatic*, we executed concept linking graph, as shown on Figure 7.

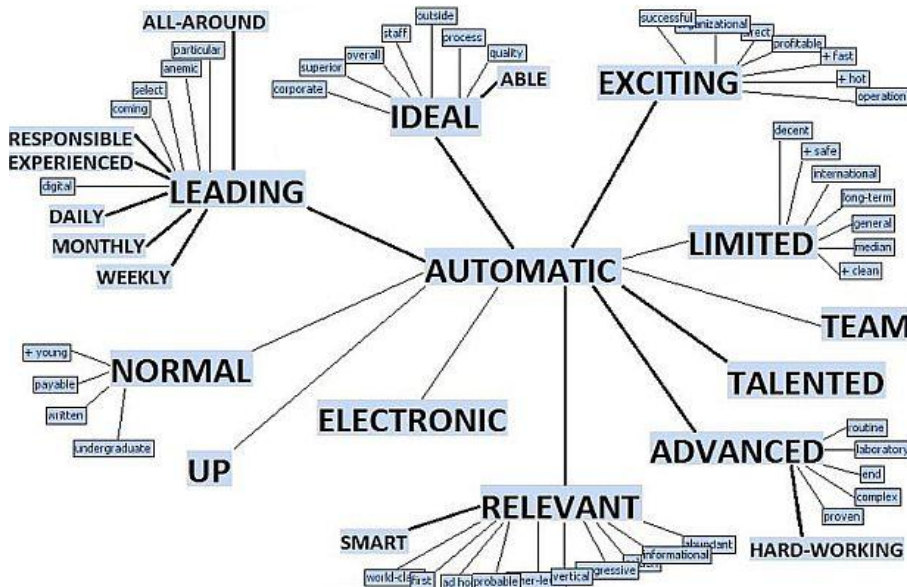


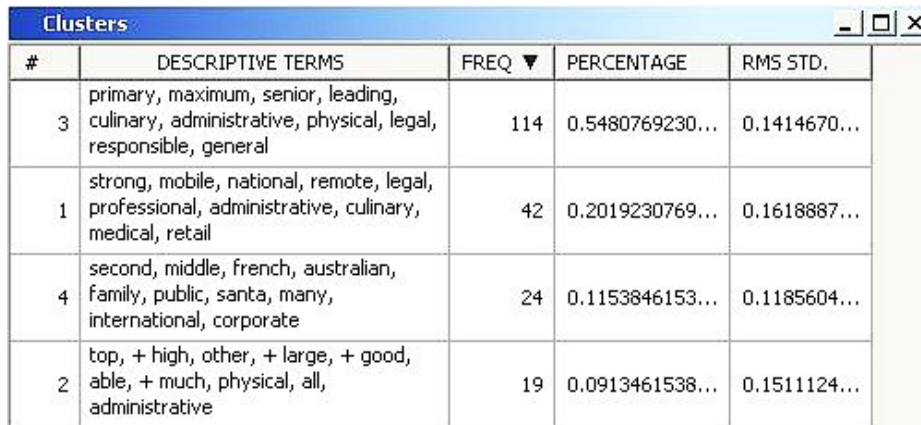
Fig. 7. Concept linkages for the term automatic. Source: own elaboration

The term *automatic* is surrounded by *leading, ideal, exciting, limited, team, talented, advanced, relevant, electronic, up, and normal*. Thick connection lines indicate close associations between the term *automatic* and adjectives: *leading, ideal, exciting, talented, advanced, and relevant*. Further expanded view once again relates to description of potential candidate's personal features (*experienced, responsible, all-around, able, hard-working, smart*), and salary period conditions (*weekly, monthly, daily*).

4.4 Jobamatic

In Jobamatic database of 208 remote work offers we identified 778 adjectives occurring in text. Among the most weighted adjectives were: *magnificent* (0.969), *statistical* (0.906), and *advisory* (0.895).

The Expectation Maximization clustering algorithm emerged four clusters, as presented on Figure 8.



#	DESCRIPTIVE TERMS	FREQ ▼	PERCENTAGE	RMS STD.
3	primary, maximum, senior, leading, culinary, administrative, physical, legal, responsible, general	114	0.5480769230...	0.1414670...
1	strong, mobile, national, remote, legal, professional, administrative, culinary, medical, retail	42	0.2019230769...	0.1618887...
4	second, middle, french, australian, family, public, santa, many, international, corporate	24	0.1153846153...	0.1185604...
2	top, + high, other, + large, + good, able, + much, physical, all, administrative	19	0.0913461538...	0.1511124...

Fig. 8. Clusters sorted by Frequency. Source: own elaboration

Third and first clusters, that grouped 156 job offers altogether, referred to miscellaneous types of job offers (*culinary, administrative, physical, legal, medical, retail*). In the fourth cluster, with descriptive terms for 24 offers, we can find geographical keywords such as French, Australian and international. Keywords of the second cluster, gathering 19 job offers, look vague and it seems hard to find any particular topic or theme for that cluster.

The adjective *statistical* occurred among terms with the highest weight, therefore we executed concept linking visualization for this term, as presented on Figure 9.



Fig. 9. Concept linkages for the term statistical. Source: own elaboration

The term *statistical* appears to be related with adjectives *culinary*, *administrative*, *responsible*, *technical*, *main*, *human*, *leading*, *social*, *legal*, *online*, and *academic*. Since it can be judged by a visual observation of the thickness of the lines, the term *statistical* reveals a higher association with terms *administrative*, *responsible*, *technical*, *leading*, *social*, *academic* than it exhibits with other terms in the text collection. Expanded subnodes are worth a closer look; in particular the term *retail*, related to the type of work, as well as other concepts associated to *statistical* (*financial*, *forensic*).

4.5 Simplyhired

In Simplyhired dataset we identified 286 adjectives occurring in text of 79 remote work offers. The most weighted adjectives were: *proud* (0.926), *associate* (0.897), and *interactive* (0.885).

In next step we executed clustering and we obtained four clusters, as presented on Figure 10. The first cluster, in spite of being the largest and gathering over 40% of documents, was unclear and therefore not particularly helpful to uncover any interesting topic in Simplyhired offers. Whereas third and second cluster, grouping altogether 32 job offers, referred to various types of jobs (*commercial*, *medical*, *local*), and the way of their performance (*simple*, *easy*, *fast*). Fourth cluster referred to the duration of the contract (*monthly*, *temporary*, *seasonal*).

#	DESCRIPTIVE TERMS	FREQ ▼	PERCENTAGE	RMS STD.
1	remote, clear, next, inaccurate, seeing, + late, similar, relevant, recent, free	32	0.4050632911...	0.1824150...
3	commercial, medical, senior, + full, new, pacific, local, all, inaccurate, seeing	20	0.2531645569...	0.2228919...
2	dream, simple, awesome, + easy, + few, single, good, + fast, important, special	12	0.1518987341...	0.1626406...
4	median, monthly, male, temporary, seasonal, nearby, marital, + little, retail, professional	11	0.1392405063...	0.1158223...

Fig. 10. Clusters sorted by Frequency. Source: own elaboration

Subsequently, we executed concept linking for one of the term with the highest weight, namely *interactive*, as presented on Figure 11.

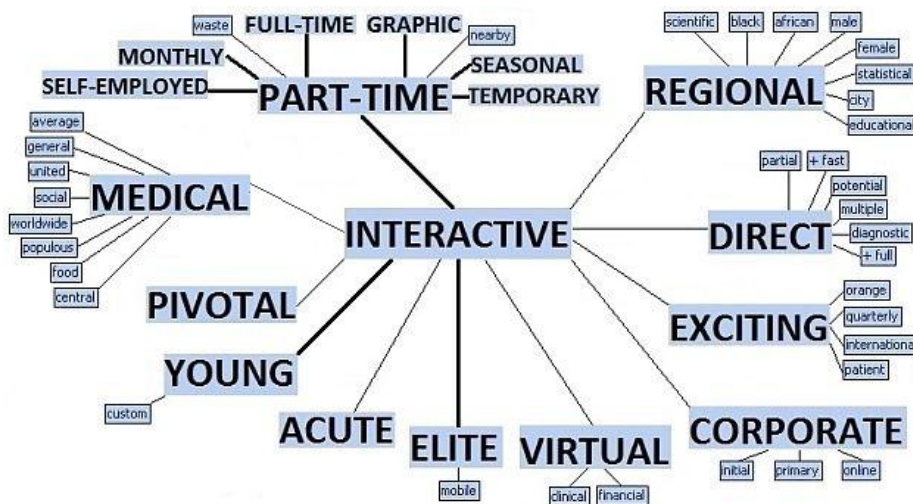


Fig. 11. Concept linkages for the term interactive. Source: own elaboration

According to the concept linking hyperbolic tree graph, the term *interactive* is associated with terms: *medical*, *part-time*, *regional*, *direct*, *exciting*, *corporate*, *virtual*, *elite*, *acute*, *young*, and *pivotal*. The width of the line outgoing from the centered adjective reveals strong correlation with the term *part-time* and its subnodes relating to contract duration (*monthly*, *full-time*, *seasonal*, *temporary*) and desirable candidate's characteristics (*self-employed*, *graphic*).

5 Discussion

We conducted our text mining study to uncover relevant features pertaining to remote job offers. Initially, we examined adjectives with the highest weights emerged in each dataset (*liable, security, advanced, clinical, residential, industrial, sole, automatic, editorial, magnificent, statistical, advisory, proud, associate, interactive*). They tend to describe personal features expected from the candidates. Without specifying a particular profession, they indicate the characteristics of particular importance in any type of work: *liable, advanced, editorial, magnificent, advisory, proud, interactive*.

Subsequently, each of five datasets was subjected to clustering. Hence, job offers were divided into sets of similar content, described by ten most relevant keywords. By analyzing them we identified four most recurrent attributes of remote work offers. The first attribute related to contract duration, and it covered terms such as *part-time, full-time, interim, temporary, daily, biweekly, monthly, seasonal*. Second group of telework offers contained names of miscellaneous types of job (*culinary, administrative, physical, legal, medical, retail, dental, commercial, local*), as well as the way of accomplishing them (*simple, easy, fast*). The third topic that appeared in advertisements concerned candidate's features: *responsible, excellent, fast, junior, able, senior, last [year], best, top, successful, available*. The fourth feature, that was distinguished among the clusters, constituted geographical and/or linguistic indicators such as *French, Australian, and international*. However, in several cases clusters turned out to be described by quite irrelevant or meaningless adjectives. It was hard to attribute any particular topic or theme for these clusters, and they were not considered as significant for further analysis.

Afterwards, for each dataset we have executed hyperbolic tree graphs in order to visually explore linkages between strongly associated terms describing mobile work offers. The adjectives with highest weights were surrounded by highly related terms, among which we could distinguish some interesting patterns. The term *advanced* was strongly correlated with *new* and *career*, as well as with *successful* and *additional* in expanded view mode. The term *clinical* was highly associated with *modern, expert, general, intelligent, well, vibrant, inner, and broad*, as well as with adjectives describing candidate's personal features required in work: *positive, smart, automotive, valuable, critical, award-winning, number one*. Close associations were detected between the term *automatic* and adjectives: *leading, ideal, exciting, talented, advanced, and relevant*. Expanded view mode for the term *automatic* related to description of desirable candidate's characteristics (*experienced, responsible, all-around, able, hard-working, smart*), and conditions of salary period (*weekly, monthly, daily*). The term *statistical* revealed high associations with terms *administrative, responsible, technical, leading, social* and *academic*. Expanded subnodes indicated *retail, financial, and forensic* as the type of work associated with statistics. *Interactive* revealed strong correlation with the term *part-time* and its subnodes relating to contract duration (*monthly, full-time, seasonal, temporary*) and to potential candidate's personal features (*self-employed, graphic*).

6 Conclusions

This paper has presented an analysis of remote work offers by using text mining techniques. The focus has been given on text parsing, text filtering, and text mining carried out for the needs of standard browsing, clustering, concept linking and further inference. Remote job offers were separated into general clusters and afterwards each cluster was automatically presented by essential descriptive keywords. By analyzing those keywords we identified the most interesting indicators for consecutive categories of job offers. By exploring interactive visualizations of highly associated terms for selected adjectives, we have examined considerable features of remote work offers text collections.

References

1. Aggarwal C. C., Zhai C., A survey of text clustering algorithms, [in:] Aggarwal C. C., Zhai C. (eds.), *Mining Text Data*, Springer, March 2012, pp. 77-128
2. Agrawal R., Batra M., A Detailed Study on Text Mining Techniques, *International Journal of Soft Computing and Engineering*, Vol. 2, Issue 6, January 2013, pp. 118-121
3. Baker P. M., Moon N. W., Ward AC., Virtual exclusion and telework: barriers and opportunities of technocentric workplace accommodation policy, *Work* 2006, Vol. 27, No. 4, pp. 421-430
4. Blansché A., Cojan J., Dufour-Lussier V., Lieber J., Molli P., Nauer E., Skaf-Molli H., Toussaint Y., Taaable 3: Adaptation of ingredient quantities and of textual preparations, , 18th International Conference on Case-Based Reasoning, Alessandria, Piemonte, Italy 2010, p. 195
5. Bloom N., Liang J., Roberts J., Ying Z. J., Does working from home work? Evidence from a Chinese experiment, *CEP Discussion Paper No. 1194*, March 2013, pp. 1-35
6. Bolasco S., Canzonetti A., Capo F. M., della Ratta-Rinaldi F., Singh B. K., Understanding Text Mining: A Pragmatic Approach, [in:] Sirmakessis S. (ed.), *Knowledge Mining*, Vol. 185 - Studies in Fuzziness and Soft Computing, Springer 2005, pp. 31-50
7. Busch E., Nash J., Bell B. S., Remote work: An examination of current trends and emerging issues, Ithaca, NY: Center for Advanced Human Resource Studies, Cornell University 2011, pp. 9-10
8. Don A., Zheleva E., Gregory M., Tarkan S., Auvil L., Clement T., Shneiderman B., Plaisant C., Discovering interesting usage patterns in text collections: integrating text mining with visualization, *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*, 2007, pp. 213-222
9. Dutcher E. G., Saral K. J., Does Team Telecommuting Affect Productivity? An Experiment, *MPRA Paper No. 41594*, Webster University Geneva, September 2012, pp. 1-29
10. Fuhr J. P., Pociask S., Broadband and telecommuting: Helping the U.S. environment and the economy, *Low Carbon Economy*, Vol. 2, No. 1, Scientific Research Publishing 2011, pp. 41-47
11. Garrett R. K., Danziger J. N., Which telework? Defining and testing a taxonomy of technology-mediated work at a distance, *Social Science Computer Review*, Vol. 25, No. 1, Spring 2007, pp. 27-47

12. Gharehchopogh F. S., Khalifehlou Z. A., Study on Information Extraction Methods from Text Mining and Natural Language Processing Perspectives, *AWERProcedia Information Technology & Computer Science*, Vol. 1, 2012, pp. 1321-1327
13. Gottfried K., The World of Work: Global Study of Online Employees Shows One in Five (17%) Work from Elsewhere, Ipsos 2012, press release 23.01.2012, <http://www.ipsos-na.com/news-polls/pressrelease.aspx?id=5486> (April 2013)
14. Gupta V., Lehal G. S., A Survey of Text Mining Techniques and Applications, *Journal of Emerging Technologies in Web Intelligence*, Vol. 1, No. 1, August 2009, pp. 60-76
15. Jusoh S., Alfawareh H. M., Techniques, Applications and Challenging Issue in Text Mining, *IJCSI International Journal of Computer Science Issues*, Vol. 9, Issue 6, No. 2, November 2012, pp. 431-436
16. Kaur R., Aggarwal S., Techniques for Mining Text Documents, *International Journal of Computer Applications*, Vol. 66, No.18, March 2013, pp. 25-29
17. Lin D., Pantel P., DIRT - Discovery of Inference Rules from Text, *Proceedings of ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2001, pp. 323-328
18. Lister K., Harnish T., The State of Telework in the U.S. - How Individuals, Business, and Government Benefit, *Telework Research Network* 2011, pp. 21-25
19. Mahesh T. R., Suresh M. B., Vinayababu M., Text mining: advancements, challenges and future directions, *International Journal of Reviews in Computing*, Vol. 3, 2010, pp. 61-65
20. Mann S., Holdsworth L., The psychological impact of teleworking: stress, emotions and health, *New Technology, Work and Employment*, Vol. 13, No. 3, Blackwell Publishing 2003, pp. 196-211
21. Patel F. N., Soni N. R., Text mining: A Brief survey, *International Journal of Advanced Computer Research*, Vol. 2, No. 4, Issue 6, December 2012, pp. 243-248
22. Radovanović M., Ivanović M., Text Mining approaches and applications, *Novi Sad Journal of Mathematics*, Vol. 38, No. 3, 2008, pp. 227-234
23. Ramanathan V., Meyyappan T., Survey of Text Mining, *International Conference on Technology and Business Management*, March 2013, pp. 508-514
24. Schadler T., Brown M., Burnes S., US Telecommuting Forecast, 2009 To 2016, *Information & Knowledge Management Professionals*, Forrester Research 2009, pp. 2-10
25. Strońska E., Elastyczne formy zatrudnienia. Telepraca. Zarządzanie pracą zdalną, *Poltext*, Warszawa 2012, pp. 115-134
26. Su Z., Kogan J., Nicholas C., Constrained clustering with k-means type algorithms, [in:] Berry M. W., Kogan J., *Text Mining: Applications and Theory*, Wiley 2010, pp. 81-103
27. Teh B. H., Ong T. S., Loh Y. L., The acceptance and effectiveness of telecommuting (work from home) in Malaysia, *Global Conference on Innovations in Management*, London 2011, pp. 34-51
28. Vidhya K. A., Aghila G., Text Mining Process, Techniques and Tools: an Overview, *International Journal of Information Technology and Knowledge Management*, Vol. 2, No. 2, 2010, pp. 613-622
29. Ward N., Shabha G., Teleworking: an assessment of socio-psychological factors, Facilities, Vol. 19, No. 1-2, Birmingham 2001, pp. 61-71
30. Watad M. M., Jenkins G. T., The Impact Of Telework On Knowledge Creation And Management, *Journal of Knowledge Management Practice*, Vol. 11, No. 4, December 2010, pp. 237-251

Tytuł artykułu:

Text Mining w praktyce: odkrywanie wzorców w tekstach ofert pracy zdalnej

Streszczenie:

Niniejszy artykuł ma na celu przedstawienie technik text mining na przykładzie nieuporządkowanych zbiorów tekstowych ofert pracy zdalnej. Przeanalizowano pięć najbardziej popularnych i najczęściej odwiedzanych portali internetowych, zawierających oferty pracy, aby wydobyć przydatne informacje, odkryć ciekawe wzorce w danych oraz wyłonić cechy telepracy. Przebadało również uzyskane klastry ofert pracy zdalnej, słowa kluczowe opisujące te klastry, a także relacje między silnie powiązаныmi terminami występującymi w różnych ofertach pracy. Dynamiczny rozwój narzędzi text mining umożliwia eksplorację nowych tematów badawczych, związanych m.in. z Internetem, freelancingiem oraz rynkiem telepracy.

Słowa kluczowe:

text mining, eksploracja danych tekstowych, klasteryzacja, drzewa powiązań, praca zdalna, telepraca