

АЛГОРИТМИЗАЦИЯ И ПРОГРАММИРОВАНИЕ

УДК 681.142.2

Yu. O. German¹, O. V. German¹, S. Nasr²

¹Belarusian State Technological University

²Belarusian State University of Informatics and Radioelectronics

INFORMATION EXTRACTION METHOD FROM A RESUME (CV)

An approach to information extraction from a short and poorly structured text document such as a resume (CV) is suggested. The computer-based resume processing is an actual interesting application problem. There are a number of web-sites for centralized CVs allocation oriented at different employers. An employer is often more interested in some peculiar features connected to professional achievements and knowledge of the applicant, not in a resume as a whole. Extraction of such peculiar information from a CV is a problem itself, especially if the CV is organized in an arbitrary form, poorly structured and contains grammatical mistakes. The suggested paper is devoted to the processing of CVs of this type. There is a short review of the existing approaches to information extraction from a CV, a keyword-based approach is selected and founded from the viewpoint of efficient information extraction the employer is interested in. The specificity of the approach is emphasized for the case when keywords define text blocks with a constant conceptual content. In this case, another problem arises, which is connected with the definition of such blocks. An approach based on a clustering technique is suggested, so each cluster is associated with a corresponding text block. At the same time, the technical realization of the approach suggested remains open for future investigations. The paper provides examples illustrating the described text extraction technique from a CV used in order to get a relevant answer to an arbitrary query.

Key words: resume, search retrieval system, text processing, key words search, clusterization.

Ю. О. Герман¹, О. В. Герман¹, С. Наср²

¹Белорусский государственный технологический университет

²Белорусский государственный университет информатики и радиоэлектроники

МЕТОД ИЗВЛЕЧЕНИЯ ИНФОРМАЦИИ ИЗ РЕЗЮМЕ

Предлагается подход для извлечения информации из коротких и плохо структурированных текстовых документов, например резюме (CV). Компьютерная обработка резюме является актуальной и интересной прикладной задачей. Имеется ряд сайтов для централизованного размещения резюме, ориентированных на различных работодателей. Работодателя часто интересуют детали, а не резюме в целом. Особенно это касается профессиональных навыков и достижений заявителя. Извлечение такого рода информации уже является проблемой, если резюме составлено в произвольной форме, плохо структурировано, содержит грамматические ошибки. Предлагаемая статья ориентирована именно на обработку такого рода заявлений. Приводится анализ существующих подходов к извлечению информации из резюме, обоснован выбор подхода, использующего ключевые слова, с помощью которых можно эффективно извлекать информацию, интересующую работодателя. Отмечена специфика подхода для случая, когда ключевые слова определяются для блоков текста с фиксированным смысловым содержанием. В этом случае возникает еще одна проблема, связанная с определением таких блоков. Предлагается подход на основе техники кластеризации, так что каждый кластер ассоциируется с соответствующим блоком текста. Вместе с тем техническая реализация этого подхода остается открытой и составляет предмет дальнейшего исследования. Приведены примеры, иллюстрирующие излагаемую технику извлечения текста из резюме как релевантного ответа на соответствующий запрос произвольной формы.

Ключевые слова: резюме, система поиска ответов, обработка текста, поиск ключевых слов, построение кластеров.

Introduction. One of the actual applied problems is an automatic resume (CV) processing. There remain a lot of job applications which are unstructured or badly organized, contain grammatical mistakes or even are incomplete. The other problem consists in potentially tremendous amount of CVs which should be processed in quite a short time interval. Finally, an employer may be inter-

ested in some specific features the applicant should possess. It follows from this that there is a need for a computer method to extract necessary data with respect to a resume with poor organization, mistakes and incompleteness.

There are three main approaches to realize such a method: keywords usage [1, 2], getting a DOM-structure of the resume considered as HTML

document [3] and a text classification approach [4]. There still remains a necessity to cope with grammatical mistakes and provide relevant answer to the query. From this viewpoint we restrict our considerations by keywords usage approach only with necessary modifications. Indeed, the approach that uses DOM-structure supposes that an original document is Html-document with a strictly defined structure. This supposition is too restrictive for our goals because our approach allows any CV form, even with semantic discrepancies. What concerns the text classification approach, it also supposes that one's CV text is divided into blocks with some strong semantics: one block is used as personal data including name, date and place of birth the other block stands for professional features and so on. We exclude this approach from our consideration.

Main part. Consider as an example the following CV:

1. *my name is Oliver Stone (A);*
2. *I am 23 years old, unmarried (B);*
3. *I am a junior researcher at university of Faradei (C);*
4. *I got a BS in computer science in 2010 from Royal College (D);*
5. *I graduated from Royal College at 2010 (D);*
6. *know such computer languages as c#, java, python (E);*
7. *my interests are connected to programming and languages (E);*
8. *my salary requirements are moderate (F);*
9. *was engaged in some projects with practical outcome (G);*
10. *I have published some articles in computer magazines in my professional areas (G);*
11. *I like to read books and listen to modern music (H);*
12. *also, I have vivid interests in business and financial programming (H);*
13. *support good relations with other people (H);*
14. *by this I apply for a vacant position of a java programmer (I);*
15. *my contact address is Belheigm city, Corwell street 10 (J);*
16. *please, use my e-mail: olstm_2123@cor.com (J).*

In round brackets we placed the semantic block identifiers (this point is discussed later).

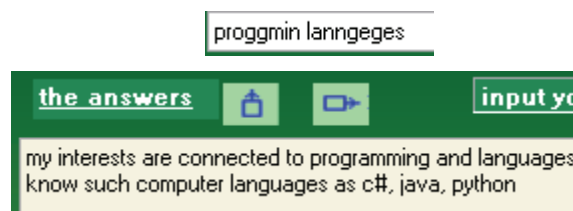
The computer program we developed provides a possibility to input any question to this text and get an answer. A question may contain distorted keywords, for example:

whats yor nname

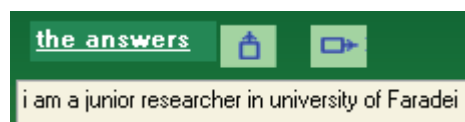
The answer is placed below



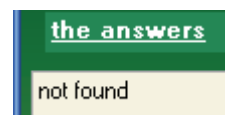
Other examples:



worked as reseraczer



A question is formulated as a set of (key)words with possible grammar mistakes. The answer is presented as a text block or even a single clause from some resume relevant to this question. If the answer is not found then we have the following resulting screenshot



We have to explain the method of CV-text processing. The entire process is divided into four stages. The first stage is to transform Word document or a pdf-file with a CV to a plain text. This can be done with a Tika system [5] which extracts a raw text and deletes unnecessary control information such as colors, fonts, and the like. The next step is to get keywords of the text (see [6]). The idea is to consider practically all text words as keywords due to not big size of a CV. Each sentence should contain at minimum one keyword. Besides, this simplifies the algorithm as it does not require to define each keyword score and test whether each sentence is covered by at least one keyword. The prepositions, conjunctions, pronouns, auxiliary and modal verbs (such as *can, have, may* etc.) are excluded. For the cursive text of the CV given as an example above, a set of keywords contains the words *name, Oliver, years, unmarried, researcher, university, Faradei, College, interests* etc. The keywords are selected in such a way that similar words are identified as the same. For instance, *programming* and *programmer* are considered as one keyword. Thus each keyword labels one or more sentences. To realize the keyword selection we use the Dice metrics (measure) given by the formula

$$P = \frac{2 \cdot |X \cap Y|}{|X| + |Y|},$$

where $|X|$ ($|Y|$) stands for the set X (Y) size. For example, for $X = programming$ and $Y = programmer$

one has $P = 2 \cdot 8 / (11 + 10) \approx 0.8$. We adopt the rule according to which two words are considered alike if the Dice measure is 0.5 or greater. From this, in the example above the words *programming* and *programmer* are considered one keyword.

It should be noticed here that fuzzy sequential text searching methods represent evident practical interest especially if they are applied on servers with very large amount of CVs to be processed [7, 8].

The next stage consists of constructing two directories. Each directory represents a collection of $\langle \text{key}, \text{value} \rangle$ pairs. The first directory contains the pairs of items and each pair represents a record with *key* standing for a keyword and *value* representing a set of numbers of the sentences labelled with the keyword *key*. The second dictionary consists of the pairs representing the numbers of the sentences and their texts. Now let us explain how to extract an answer to the question represented by a set of words (some of words are keywords and some are not (these latter words do not belong to the CV text)). For each keyword k_i in the question the set of numbers of the sentences $N_i = \{n_{i1}, n_{i2}, \dots, n_{iz}\}$ is defined. Then for all N_i the most frequently encountered number(s) n_w is (are) defined. This number n_w defines a sentence to be displayed as an answer. If there are more than one candidate to be an answer then all candidates are displayed.

It may be the case when a CV contains quite big blocks of semantically connected information. This situation is somewhat wider than that one considered above. Each semantically united block may consist of one or more sentences. The block specifies some concept that may be considered as a semantical whole. Practically, we associate such blocks with paragraphs (indentations) or the text segments separated from the others with empty lines or by spaces. As above, we should define the keywords for the sentences: $k_1, k_2, k_3, \dots, k_z$.

We introduce the third dictionary with *keys* representing the numbers of the sentences and *values* standing for text blocks. Thus, we get the following chain: keyword \rightarrow sentence(s) \rightarrow text block(s). The searching procedure remains as before. Keywords from the query are used to find the numbers of the sentences and the corresponding text blocks. The text block which is referred to by the majority of keywords is then selected. To be more understandable, let us take as an example the query:

which programming languages are you interested in?

The keywords *programmer*, *language*, *interests* all are presented in the query with the corresponding sentences numbers:

programmer: 7(E), 12(H), 14(I)

language: 6(E), 7(E)

interests: 7(E), 12(H)

In round brackets the block identifiers are placed. As one can see, the block *E* is referred to most often and is therefore selected. It includes the sentences 6 and 7. The final step consists of definition of the text diapason to be displayed. Evidently, the rational position is to display all the sentences with the numbers starting from the minimum and ending with the maximum number among the numbers of the sentences selected. That is, in our case the sentences 6 and 7 will be displayed both:

know such computer languages as c#, java, python;

my interests are connected to programming and languages.

Let us consider one more example of that kind. Suppose that a query is the following:

what are your interests?

Here the only keyword is presented:

interests: 7(E), 12(H)

This time the answer is composed of the 7th and the 12th sentences.

The problem of blocks definition has a special scientific meaning. The simplest case consists in dividing the text by paragraphs separated with indentions or empty lines. Evidently, each block should contain connected notions (keywords). One says that related keywords form a cluster. So, it is necessary to build a number of clusters, combining different but semantically related keywords. The first problem with clusters is their number. Obviously, one needs to apply the existing clustering means to define an optimal cluster structure. For these aims a Python statistical modules can be used, such as, for example, SciLearn [9] or R language [10] or specialized programming packages such as [11]. As an alternative approach one can use a correlation-based technique [12]. The second problem is to map the cluster structure to blocks. The idea is to apply covering procedure to associate the clusters with the sentences labelled (covered) with the keywords from that cluster. One, however, should keep in mind that the covering procedure must select the sentences with sequentially ordered numbers. As this is quite a special issue going beyond the frames of the paper, we omit it.

In statistics, the measure of items connectivity may be characterized with the correlation coefficient [12]. The correlation coefficient between two keywords x and y is defined as below:

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \cdot \sqrt{\sum_i (y_i - \bar{y})^2}}$$

With the help of the above formula one is in position to find pairwise correlations of the keywords. Here, $x_i = 1$ if keyword x occurs in the i_{th}

text block (the same is right with other keyword y). The cluster structure we are interested in should provide maximum value of the correlation coefficient weighted through the entire CV.

Conclusion. The above given material can be directly used in automatic CV processing. There are some main application problems, which can be solved with the outline technique. The most important is that one connected to extracting necessary information from the CV text which represents interest to employer. The approach we used here do not restrict the type of information. It is an evident advantage of it. The user can use mistakenly written queries and apply them to mistakenly written CVs. The second type of the problems to deal with is a CV automatic structuring to write them into database. Our approach makes one enable to get a structured CV by cluster definition and labelling clusters with keywords.

There remain some open questions though. The most important one is a mapping the cluster structure to blocks structure. This issue requires further considerations. In general, one may admit that the same sentence belongs to different blocks. We use the approach described below. To define text blocks we select verbs supposing that each verb stands for some relation(s) in the domain of the text and plays a role of some semantic function defined on the domain. The commonly used verbs (such as have, can, be, get etc.) should be excluded. Also it is important that total number of English verbs is not so big and gives a possibility to use them from a previously built and indexed database. Let us address the abstract of this article as an example. The set of verbs consists of the next items: *suggest* {1, 6, 10, 11}, *process* {2, 6}, *extract* {1, 5, 7, 12}, *orient at* {3}, *select* {7}, *is interested* {2, 4, 7}, *employ* {3, 7}, *define* {8, 9}, *connect* {4, 9}, *arise* {9}, *base* {2, 10}, *cluster* {10}, *organize* {5}, *associate* {10}, *correspond* {10}, *realize* {11}, *apply* {2, 4}, *illustrate* {12}, *address* {12}, *resume* {2, 4}, *investigate* {11}, *contain* {5}, *devote* {6}. In figure brackets there are the numbers of the sentences. Our supposition is that verbs define the sense of the text, so they can be used to select text blocks (as segments with some completed conception). The next step is to define the set of verbs covering with minimum number of verbs selected. A covering π consists of the verbs and each sentence from abstract contains at least one verb from π . For instance, one of the possible minimum-size covers is represented by the set $\pi = \{extract, is interested, define, employ, suggest\}$. We should notice that finding a minimum-size covering is a hard computational complexity problem, so one can use any good heuristic method to solve it. The size of π represents the number of certain text blocks.

In our example the abstract of the paper should be divided into 5 blocks. The simplest way to define them is to use the numbers of the sentences covered by the verbs. For instance, verb *extract* covers the sentences 1, 5, 7, 12. So, the sentences with these numbers define a text block in the approach we describe. We have 5 blocks here: $A = \{1, 5, 7, 12\}$, $B = \{2, 4, 7\}$, $C = \{8, 9\}$, $D = \{3, 7\}$, $E = \{1, 6, 10, 11\}$. From this, one can see that some sentence may belong to different blocks at the same time. Consider how a query “*extract info employer*” is treated. The procedure results in

extract: 1, 5, 7, 12(A)
info: 1(A, E), 5(A), 7(A, B, D)
employer: 3, 7(D)

Non-keyword item *info* refers to different blocks. The most referenced are blocks A and D . So the answer is either A or D or both of them. Each of the blocks A and D should be scanned separately in order to define the most relevant answer. It should be noted here that additional block processing is needed as the most relevant answer may correspond even to one part of a block. Thus, in the example we have

(1) *An approach to information extraction from a short and poorly structured text document such as a resume (CV) is suggested.*

(5) *Extraction of such peculiar information from a CV is a problem itself especially if the CV is organized in an arbitrary form, poorly structured and contains grammatic mistakes.*

(7) *There is a short review of the existing approaches to information extraction from a CV, a keyword-based approach is selected and founded from the viewpoint of efficient information extraction the employer is interested in.*

(12) *The paper provides examples illustrating the described text extraction technique from a CV used in order to get a relevant answer to an arbitrary query.*

From these sentences only (7) and (12) should be taken as an answer from the block A . Consider then the block D :

(3) *There are a number of web-sites for centralized CVs allocation oriented at different employers.*

(7) *There is a short review of the existing approaches to information extraction from a CV, a keyword-based approach is selected and founded from the viewpoint of efficient information extraction the employer is interested in.*

Sentence 7 is the most relevant. Combining both blocks A and D the final answer is represented by the sentence 7.

Evidently, one can apply the total approach not to CVs only. Any type of not long texts can be processed in the same way. These include paper abstracts, e-mail contents, patent formula descriptions and the other short document types. As a further research direction we point to the problem of short

document clustering. The idea to use keywords is not the only one. Text ontology tools, semantic networks

and the other means also may be used. It is also necessary to consider the usage of synonyms.

References

1. Maheshwari S., Sainani P., Reddy K. An approach to extract special skills to improve the performance of resume selection. In: Databases in Networked Information Systems. Lecture Notes in Computer Science. Vol. 5999, pp. 256–273. Berlin, Germany, Springer, 2010.
2. Kopparapu S. K. Automatic extraction of usable information from unstructured resumes to aid search. Proceedings of the 1st IEEE International Conference “Progress in Informatics and Computing (PIC ’10)”. 2010, vol. 1, pp. 99–103.
3. X Ji, Zeng J., Zhang S., Wu C. Tagtree template for Web information and schema extraction. Expert Systems with Applications. 2010, vol. 37, no.12, pp. 8492–8498.
4. Yu K., Guan G., Zhou M. Resume information extraction with cascaded hybrid model. In Proceedings of the 43rd Annual Meeting “Association for Computational Linguistics (ACL’05)”. 2005 June, pp. 499–506.
5. The Apache Software foundation. Apache Tika – a content analysis toolkit. Available at <http://tika.apache.org> (accessed 12.02.2019).
6. Buttcher S., Clarke C. L. A., Cormack G. V. Information retrieval. Implementing and evaluating search engines. The Mit Press. London. England. 2010. 606 p.
7. Tarhio J., Ukkonen E. Approximate Boyer-Moore String matching. SIAM Computing. Vol. 22, 1993, pp. 243–260.
8. Anu S., Joby G. Fuzzy pattern matching algorithm for location based approximate strings. International journal of scientific and engineering research. Vol. 7, issue 7, 2016, pp. 583–587.
9. Sheppard K. Introduction to Python for econometrics, statistics and data analysis. University of Oxford, 2014. 394 p.
10. Trevor M. Undergraduate Guide to R. Available at: <http://www.biostat.jhsph.edu/~ajaffe/docs/under-gradguidetoR.pdf> (accessed 12.02.2019).
11. Kirilov A. Projects and applications using AForge.NET Framework. Available at: <https://www.codeproject.com/Articles/16859/AForge-NET-open-source-framework> (accessed 12.02.2019).
12. Cox D. R., Snell E. J. Applied statistics. Principles and examples. Chapman & Hall. 1981. 190 p.

Information about the authors

German Yulia Olegovna – PhD (Engineering), Senior Lecturer, the Department of Information Systems and Technology. Belarusian State Technological University (13a, Sverdlova str., 220006, Minsk, Republic of Belarus). E-mail: juliagerman@tut.by

German Oleg Vitoldovicz – PhD (Engineering), Assistant Professor, the Department of Information Systems and Technology. Belarusian State Technological University (13a, Sverdlova str., 220006, Minsk, Republic of Belarus). E-mail: ovgerman@tut.by

Nasr Sara – PhD student, the Department of Information Technologies in Automated System. Belarusian State University of Informatics and Radioelectronics (6, P. Brovki str., 220600, Minsk, Republic of Belarus). E-mail: sara.nasrh@gmail.com

Поступила 15.01.2019