



Natural Language Processing in assistance to Inventive Design activities

Daria Berdyugina, Denis Cavallucci

► To cite this version:

Daria Berdyugina, Denis Cavallucci. Natural Language Processing in assistance to Inventive Design activities. 32nd CIRP Design Conference (CIRP Design 2022) - Design in a changing world, March 28th to March 30th, 2022, Paris, Mar 2022, Paris, France. pp.7-12, <10.1016/j.procir.2022.05.206>. <hal-03797640>

HAL Id: hal-03797640

<https://hal.science/hal-03797640v1>

Submitted on 4 Oct 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

6th CIRP Conference on Surface Integrity

Natural Language Processing in assistance to Inventive Design activities

Daria Berdyugina¹, Denis Cavallucci¹^aINSA Strasbourg/ICube Lab., 24, Boulevard de la Victoire, Strasbourg 67000, France* Corresponding author. Tel.: +33-7-5561-0209 ; E-mail address: daria.berdyugina@insa-strasbourg.fr**Abstract**

Design activity requires engineers to use their knowledge to understand, formulate and solve industrial problems whose complexity is constantly increasing. Our work, at the frontier of linguistics, computer science and engineering sciences, consists in making optimal use of the information contained in texts useful to designers, particularly patents. Existing method to extract information from patents already exist, but without taking into consideration the use of their output as a mean to populate Inventive Design ontology. This paper describes how we use Natural Language Processing to extract information useful to designers in the context of an inventive activity. Our ontology inspired by the Theory of Inventive Problem Solving (TRIZ) helps us to categorize the information collected in corpus consisting of thousands of unstructured texts, to redistribute it to the designer and thus increase his ability to formulate and solve the problems that his subject of study confronts him with. This study concerns a new way of categorizing patents into newly defined technical domains so as to better perceived a prior art during the early stages of an inventive process. Our results attest significant progresses both in terms of the speed to understand data complexity and in terms of the accuracy of the perception of an explored field prior to enter into a resolution logic.

© 2022 The Authors. Published by Elsevier B.V.

This is an open access article under the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nc-nd/4.0>)

Peer-review under responsibility of the scientific committee of the 32nd CIRP Design Conference

Keywords: Inventive Design Method; TRIZ; Natural Language Processing; Augmented creativity**1. Introduction**

Nowadays, we are witnessing an acceleration of scientific and technological progresses. Consequently, it is important to be capable of dealing with a huge amount of information appearing every day. Among this information, patents represent one of the most important sources of inventive information, featuring a constantly revised representation of human inventive activity. In Research and Design domain, the information taken from patent texts has two major application areas. Firstly, these texts are used to find a technical issue to solve. Secondly, patents contain an updated information about new methods and techniques that could be applied to resolve inventive problems.

Patent documents are intended to indicate the limits of a given invention (including its technical description) in order to protect the intellectual property of its owners. Therefore, the patent texts are edited using peculiar style of writing that is distinctive for legal documents. This style is manifested in long and complex phrases, repetitions and anaphora that causes difficulties

for understanding the text (especially for non-specialists in legal science).

On the other hand, patent documents contain unstructured technical information such as the state of the art of an invention's field or detailed technical description of claimed invention. In spite of the experience of patent's reader to work with technical texts from their domains, the double nature of patent documents creates obstacles for text understanding. Henceforth, the work with patents demands a lot of human resources and time-consuming.

However, with the development of computer-science branch, text analysis becomes more accessible thanks to, notably, Natural Languages Processing (NLP) [15, 17, 12, 9]. The application of NLP methods and techniques allows to automatically analyse unstructured text data (including patents) to facilitate its understanding and save the human resources.

Patent text contents represent an important amount of inventive information and could be used to feed existing creative techniques. There are a lot of creative techniques aiming to help companies in their inventive activities. For example, the brainstorming or synectics [19]. All these methods focus on structured creative exercises in order to trigger inventive ideas. However, regardless of their fame, these methods have a significant

drawback: they do not permit to break the boundaries of the knowledge and the expertise of participants of an inventive exercise [1].

Contrary to the cited creative techniques, TRIZ [2] (The Theory of Inventive Problem Solving through its Russian acronym) permits to overcome the above-mentioned limitation. The creator of TRIZ, the Soviet Engineer G. S. Altshuller, after analyzing manually 40,000 of patent texts, deduced that all inventions are made according to recurrent laws of evolution. Therefore, the inventive process could be formalized in order to facilitate the inventive activity. TRIZ methodology offers the possibility to find the most suitable inventive solution, without appealing to a compromise solution.

Combining NLP techniques for information manipulation on one of the main sources of inventive information (patent texts) allows to facilitate and accelerate the inventive process within TRIZ methodology. Notably, the automatic clustering of patent documents of a given field permits to group the information by the most representative terms of a corpus. The existing clustering algorithms permits to group the patent texts among the corpus data in order to improve the understanding of main subjects of every document and to mine some additional information about a technology expressed in patents.

In the Section 2, we describe TRIZ and Inventive Design Methods (IDM). Moreover, we explore the existing methods of clustering. In the Section 3, we present the methodology that we apply to achieve our goal. In the Section 4, we describe an example of the application of our methodology on a real-case study. Section 5 contains the discussion about the advantages and disadvantages of applied methodology and further work. Finally, in Section 6, we present a conclusion of the present article.

2. State of the art

2.1. TRIZ and IDM

TRIZ starts to gain its popularity in the 1990s [29]. This popularity is based on the capacity of TRIZ to ameliorate an inventive process making it easier and faster. This theory is applicable in all areas of human activity.

However, despite its effectiveness, classic TRIZ methods are difficult to understand and to apply, especially for non-experienced users. This fact is explainable by the absence of a formalized ontology, therefore, it is problematic to perform any computation on its abstract concepts [8].

IDM [7], based on TRIZ, aims to overcome the above-mentioned drawback. In the IDM ontology, there are 3 basic elements that participate to problem-solving process: problem, partial solution and parameters.

A problem “describes a situation where an obstacle prevents progress, an advance or the achievement of what has to be done” [22]. A partial solution “expresses a result that is known in the domain and verified by experience” [22].

The parameters, as one of the basic concepts of IDM ontology, form a contradiction. According to its definition, a con-

tradiction “[...] characterized by a set of three parameters and where one of the parameters can take two possible opposite values Va and \bar{Va} ” [22]. Thus, we distinguish two groups of parameters: action parameter (AP) and evaluation parameter (EP). An AP, “[...] is characterized by the fact that it has a positive effect on another parameter when its value tends to Va and that it has a negative effect on another parameter when its value tends to \bar{Va} (that is, in the opposite direction)” [22]. An EP could be characterized as a parameter that “[...] can evolve under the influence of one or more action parameters” and which make possible to “evaluate the positive aspect of a choice made by the designer” [22].

The goal of our research consists of automatic extraction of parameters from the domain-specific patent corpus in order to facilitate the use of IDM and save the time of users. By applying the NLP techniques, such as clustering, we could mine the information about the main concepts described in documents. Moreover, the clustering algorithms help to group the documents in order to ease the further application of information retrieval methods.

The clustering techniques, which is used for categorization, in the context of IDM-related studies, aims to help engineers to quickly take inventory of their subject upstream of a study. The clustered information facilitates the formulation of contradictions, which represent an important element of problem-solving process according to IDM techniques.

2.2. Clustering algorithms

The analysis of patent information represents one of the most important branches of technologies mapping. Patent map is a visualization method which allows to represent the patent text graphically and analyze it. These techniques are widely used in a lot of scientific approaches, for example [13, 28, 14]. The most popular approaches in patent data analysis and patent mapping are based on structured information such as dates, assignees, or citations. These data can be analyzed by traditional bibliometric techniques [5].

Nowadays, the techniques for information extraction out of unstructured patent text data become more applied thanks to the emergence of numerous data-mining methods [26]. On the other hand, other researches are based on automatic summarization methods for patent texts [25, 27].

However, the visualization of patent analysis raises an interest among the data mining community. One of the most used techniques for visualisation of patents is clustering.

In the context of present research, we consider K-means clustering algorithm which uses a distance between data points for clusterization computation. The K-Means algorithm groups data by trying to separate the samples into n groups of equal variance, thereby minimizing a criterion known as inertia or putting it within the sum of squares of the cluster. This algorithm requires to specify the number of clusters. It supports a large number of samples and is used in different applications.

The K-means algorithm divides a set of N samples X into K relatively disjoint clusters C . Each cluster is described by the mean μ_j . The mean value is commonly referred to centroid of

the cluster. The K-means algorithm (1) aims to choose centroids that minimize the inertia.

$$\sum_{i=0}^n \min_{\mu_j \in C} (\|x_i - \mu_j\|^2) [6] \quad (1)$$

Inertia can be thought of as a measure of how many internally coherent clusters are. This has some drawbacks [6].

- Inertia assumes that the cluster is convex and isotropic, but this is not always the case. It does not react to elongated clusters or irregularly shaped manifolds.
- Inertia is not a normalized metric: smaller values are better, and zero has been found to be optimal. However, in very high dimensional spaces, the Euclidean distance tends to expand (this is an example of the so-called "curse of dimensionality").

Basically, K-Means algorithm comprises three steps. By the first step, the algorithm chooses the initial centroids. The easiest way is to select k samples from the data set X . Once initialized, K-means consists of a loop between the other two steps. In the first step, each sample is assigned to the closest centroid. The second step creates a new centroid by averaging all the samples associated with each previous centroid. The difference between the old and new centroids is calculated and the algorithm repeats these last two steps until this value is less than the threshold. In other words, it repeats until the centroids does not move much [6]. K-means supports an expected value maximization algorithm that uses a small, all-same, diagonal co variance matrix.

3. Methodology

In this section, we describe the methodology used for automatic clustering of the patent documents belonging to a given field. It is important to precise that we apply an existing method of clustering in order to compare if the automatic document categorization is suitable for the task performing by humans in the context of inventive process. The actual contribution of our work is to prove that an existing clustering algorithm may establish a survey of subjects expressed in patent texts.

3.1. Corpus preprocessing

For performing our task, first, we need to transform the textual data into the suitable input for clustering algorithms. Most of the known algorithms of clustering use vectorized textual data (Word Embeddings [18]) as an input. Before performing this transformation, it is required nonetheless to apply a preprocessing pipeline in order to get the best result.

Foremost, we extend the list of stop words for patent texts in order to eliminate the noise from the final result that is caused by specific-patent words that are frequently repeated in almost

all texts. For formulation of this list, we use a tool for statistical corpus analysis Antconc [3, 4]. Thanks to this tool, we extract the list of the most common used words on the patent texts and add them to the classical English-language stop words list¹.

Secondly, we tokenize and stem the input text. The stemming [16] represents the process of finding a word base for a given source word. A word base is not necessarily the same as the morphological root of a word. Stemming is part of the text normalization process. The stemming allows to better group a semantically close words in order to ameliorate the final result of clustering of the documents.

The third step consists on grouping the poly lexical terms in order to capture multiwords terms. Thus, we search the words with the highest probability to appear together in text. We are interested in bi- and trigrams identification. The identification of poly lexical terms permits to enrich the final result for presenting not only simple terms but also wider range of lexic related to every cluster node.

3.2. Clustering of the patent document

After applying a preprocessing pipeline, we have a set of cleaned and stemmed vocabulary which could be transformed to vectors. The clustering algorithms aim to segregate groups with similar traits (for textual data, it is words) in order to assign these groups to the clusters.

For identification of the similarity, we use tf-idf² vectorizer. Tf-idf [23] is a measure of originality of a word. This originality is measured by comparison of between the frequency of an apparition of a word in a corpus with a number of documents containing this word. Term frequency (tf) is the ratio of the number of occurrences of a word to the total number of words in the document. Thus, the importance of a word t_i within a single document d is estimated (2).

$$tf(t, d) = \frac{n_t}{\sum_k n_k} [11] \quad (2)$$

where n_t is the number of occurrences of words t in the document and the denominator represent the total number of words in the document.

Inverse document frequency (idf) is the inversion of the frequency with which a word appears in collection documents (3). Idf consideration reduces the weight of commonly used words. For each unique word within a particular document collection, there is only one Idf value.

$$idf(t, D) = \log \frac{|D|}{|\{d_i \in D \mid t \in d_i\}|} [11] \quad (3)$$

¹ Available at <https://www.nltk.org/book/ch02.html>

² Term Frequency Inverse Document Frequency

where $|D|$ is the number of documents in a corpus and $|\{d_i \in D | t \in d_i\}|$ is the number of documents in collection D in which t occurs (when $n_i \neq 0$).

Thus, the tf-idf measure is the product of two factors (4):

$$tf - idf(t, d, D) = tf(t, d) \times idf(t, D) [11] \quad (4)$$

A high weight in tf-idf is given to words with a high frequency within a particular document and with a low frequency of use in other documents.

After the vectorization step, we compute a cosine similarity [24] between every possible pair of words in order to obtain a semantic distance between them. Performing this step allows to classify the document words and calculate the cluster members. The last step of the methodology is to apply k-means clustering algorithm [6] for documents classification³. As a final result, our method shows the grouped patent documents by the terms as well as the most representative terms for every cluster.

4. Case study

In order to validate our approach, we use a corpus composed of 15 patents that treat a subject of automobile's door latch mechanism. This corpus is grouped manually by the technology that is described in the text. The categorization of the texts by an expressed technology permits to investigate the main subjects upstream of each study.

The human annotation defines 11 groups. The process of human categorization took 3 days of meetings, reading and exchanges between 3 experts to reach this level of categorization of the domain. However, thanks to the clustering algorithm, described in Section 3, the same exercise took approximately 2 minutes.

We try to artificially reproduce this grouping using our methodology described in Section 3. K-means clustering algorithm requires the number of clusters to define. Thus, we use the number of groups chosen by human annotation. The Fig. 1 shows the result the of application of our methodology. Every sphere on the Fig. 1 corresponds to the patent document. The members of a same cluster are coloured using the same colour. At the marginal notes, we provide the key terms for every cluster, i.e. a subject that is mainly discussed by the documents of every cluster.

Compared to the human clustering (Table 1), the clustering algorithm succeeded to identify correctly 6 groups out of 11. However, it is important to note that human annotation is based on a common technology exposed in the patent texts. Our approach (see 3) consists of finding the similar texts with common vocabulary. In theory, the use of common vocabulary means a description of a similar technology. However, the patent writer could not resort to the use of common terms, i.e., explain the

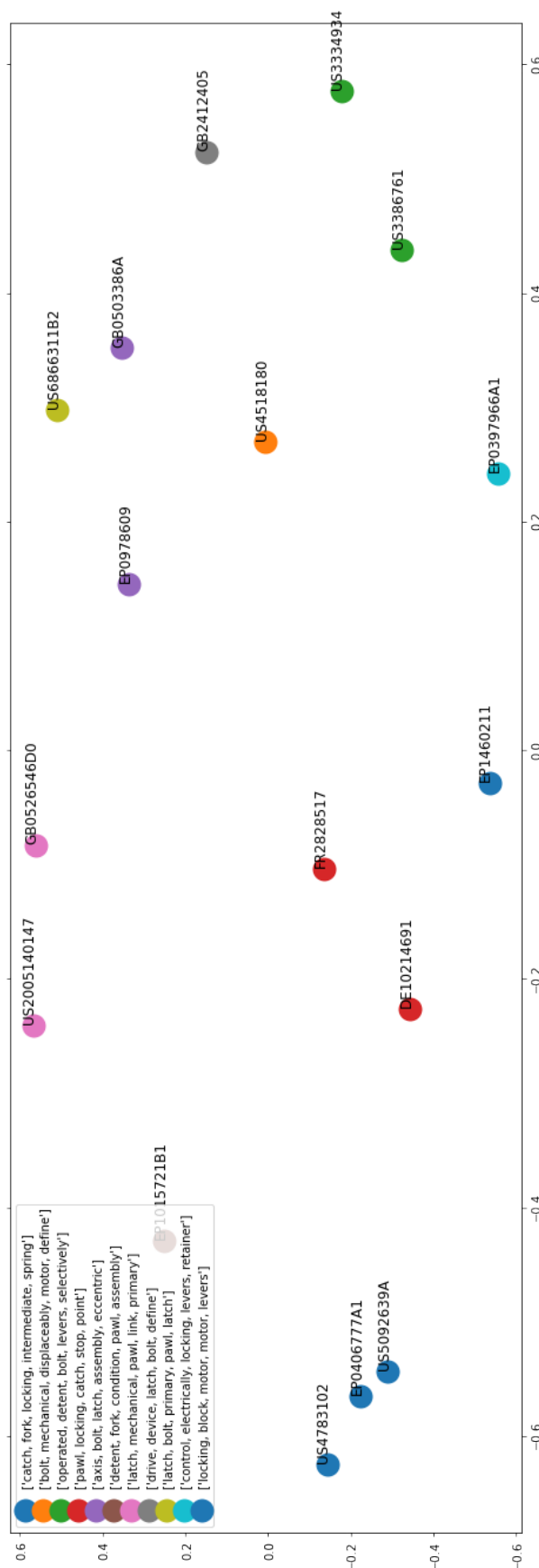


Fig. 1. Box plot of number of positions sent per iteration using this scheme

³ We use scikit-learn library available at <https://scikit-learn.org/stable/modules/clustering.html>

Table 1. Comparison of human and automatic clustering

Human selected group	Algorithm identified group	Automatic associated terms
Electrically operated lock, particularly for automotive application (EP0397966)*, Latch assembly (EP0978609) Locking device (DE10214691)*	Electrically operated lock, particularly for automotive application (EP0397966)* Locking device (DE10214691)*. Two catch automobile lock comprises rotating bolt and two movable catches able to lock bolt in closed position (FR2828517) Lock requiring reduced opening force (EP0406777)* Latch, in particular for a motor vehicle door (US4783102), Lock requiring reduced opening force (US5092639)* Latch bolt for vehicles (GB2412405)*	control, electrically, locking, levers, retainer pawl, locking, catch, stop, point
Lock requiring reduced opening force (US5092639)*, Lock requiring reduced opening force (EP0406777)* Closure latch (US3334934). Latch bolt for vehicles (GB2412405)* Vehicle body door latch and locking system (US3386761)*, Latch, in particular for a motor vehicle door (US4783102), Two catch automobile lock comprises rotating bolt and two movable catches able to lock bolt in closed position (FR2828517) Vehicle door latch with reduced release effort (EP1015721)*	 Closure latch (US3334934). Vehicle body door latch and locking system (US3386761)* Vehicle door latch with reduced release effort (EP1015721)*	catch, fork, locking, intermediate, spring drive, device, latch, bolt, define operated, detent, bolt, levers, selectively
Motor vehicle lock (EP1460211)* Low release energy latch mechanism (US2005140147)*, Latch mechanism (GB0526546D0)* Latch assembly (GB0503386) Latch mechanism for a vehicle (US6866311B2)* Automobile power door latch (US4518180)*	Motor vehicle lock (EP1460211)* Latch mechanism (GB0526546D0)*, Low release energy latch mechanism (US2005140147)* Latch assembly (EP0978609), Latch assembly (GB0503386) Latch mechanism for a vehicle (US6866311B2)* Automobile power door latch (US4518180)*	detent, fork, condition, pawl, assembly locking, block, motor, motor, levers latch, mechanical, pawl, link, primary axis, bolt, latch, assembly, eccentric latch, bolt, primary, pawl, latch bolt, mechanical, displaceably, motor, define

same subject differently. The Table 1 shows the result of comparing human grouping and automatic clustering. We provide the title of every patent text and its reference number, as well as the most representative terms for every automatic clustered group. We tag with '*' the patents in every group which match the human judgement.

5. Discussion and result evaluation

As we mentioned in Section 4, thanks to the use of our methodology, we succeed to automate the identification of 6 groups of patent documents out of 11 compared to the human annotation. However, these results could be improved in adding an automatic algorithm of technology detection.

According to distributional theory [10], the words with similar meaning are tendencies to appear in similar context. For the cluster's distribution, there is a same logic: more terms of the documents are similar, more the clusters are close.

As it discussed, the clustering algorithm permits to obtain the key terms for every cluster. Therefore, for a specialist of the given domain, it is simple to deduce the technology discussed in every cluster. Hence, these terms may be helpful in contradiction identification since some of them could be viewed as an AP or even EP.

One of the drawbacks of our method is the requirement to enter the exact number of clusters to extract. This mechanic could be inconvenient for a user because it forces him to anticipate how many subjects are exposed in his documents' collection. However, this method is useful to pretreat the corpus of a given domain in order to reveal an approximate thematic distribution, which aids to choose the exact documents in which the user is interested.

With regard to advantages, our methodology suits to the IDM-related process and allows not only to filter the questionable patent documents, but also to extract a specific terminology which could be a part of domain-specific contradiction. It is important to say that domain-specific contradictions ease the main domain-problem comprehension. Moreover, its representation in matrix form facilitate the whole problem-solving process because it allows to detect the main zones which could be ameliorated and where there is no patent application yet.

In order to objectively evaluate the results of clustering, we compute some of the most evaluation measures: Rand index [20], homogeneity, completeness and V-measure [21]. As it shown in Table 2, our algorithm demonstrates a good result.

Table 2. Evaluation score computation

Evaluation measure	Score
Rand index	0.926
Homogeneity	0.845
Completeness	0.85
V-measure	0.845

6. Conclusion

In the present article, we exposed a methodology of automatic clustering algorithm for a domain-specific patent corpus. This methodology allows categorizing the documents of a domain-specific corpus. Such categorization facilitates a problem-solving process in the context of IDM. The automatic analysis of words' distribution shows the close to human decisions' result. we consider the further improvement of our methodology in order to achieve better results and to investigate larger corpora.

The presented method feed an IDM problem-solving process which aids to find the most suitable inventive solution. The development of NLP techniques allows automating the processes of human-reading and understanding of large corpus of unstructured texts and significantly save time and human resources.

The experiment conducted in the context of the present article shows drawbacks of the exposed method. Therefore, as a further work, we aim to improve our algorithm in order to present a better result in the point of view of accuracy and also of a user experience. We project to verify the performance of other clustering technique in order to overcome the inconveniences presented in Section 5. Moreover, we prepare an algorithm for an automatic labelling of clusters, which eases the memorization activity and allows saving more time and human resources.

References

- [1] Aiken, M., Krosp, J., Shirani, A., Martin, J., 1994. Electronic brainstorming in small and large groups. *Information & Management* 27, 141–149. URL: <https://www.sciencedirect.com/science/article/pii/0378720694900426>, doi:[https://doi.org/10.1016/0378-7206\(94\)90042-6](https://doi.org/10.1016/0378-7206(94)90042-6).
- [2] Altshuller, G., 1984. *Creativity As an Exact Science*. Taylor & Francis. Google-Books-ID: Jd0OAAAAQAAJ.
- [3] Anthony, L., 2005. Antconc: design and development of a freeware corpus analysis toolkit for the technical writing classroom, in: IPCC 2005. Proceedings. International Professional Communication Conference, 2005., pp. 729–737. doi:[10.1109/IPCC.2005.1494244](https://doi.org/10.1109/IPCC.2005.1494244).
- [4] Anthony, L., 2020. Antconc (version 3.5.9) [computer software]. Available from <https://www.laurenceanthony.net/software>.
- [5] Archibugi, D., Planta, M.M.G., 1996. Measuring technological change through patents and innovation surveys. *Technovation* 16, 451–519.
- [6] Arthur, D., Vassilvitskii, S., 2007. K-means++: The advantages of careful seeding, in: Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, Society for Industrial and Applied Mathematics, USA. p. 1027–1035.
- [7] Cavallucci, D., 2012. From triz to inventive design method (idm): towards a formalization of inventive practices in r&d departments.
- [8] Cavallucci, D., Khomenko, N., 2007. From TRIZ to OTSM-TRIZ: Addressing complexity challenges in inventive design. *International Journal of Product Development* 4, 1477– 9056. URL: <https://hal.archives-ouvertes.fr/hal-00686768>, doi:[10.1504/IJPD.2007.011530](https://doi.org/10.1504/IJPD.2007.011530).
- [9] Guida, G., Mauri, G., 1986. Evaluation of natural language processing systems: Issues and approaches. *Proceedings of the IEEE* 74, 1026–1035. doi:[10.1109/PROC.1986.13580](https://doi.org/10.1109/PROC.1986.13580).
- [10] Harris, Z., 1954. Distributional structure. *Word* 10, 146–162. URL: https://link.springer.com/chapter/10.1007/978-94-009-8467-7_1, doi:[10.1007/978-94-009-8467-7_1](https://doi.org/10.1007/978-94-009-8467-7_1).
- [11] Jones, K.S., 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation* 28, 11–21.
- [12] Jurafsky, D., Martin, J., 2008. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. volume 2.
- [13] Kim, Y.G., Suh, J.H., Park, S.C., 2008. Visualization of patent analysis for emerging technology. *Expert Systems with Applications* 34, 1804–1812. URL: <https://www.sciencedirect.com/science/article/pii/S0957417407000577>, doi:<https://doi.org/10.1016/j.eswa.2007.01.033>.
- [14] Lee, S., Yoon, B., Park, Y., 2009. An approach to discovering new technology opportunities: Keyword-based patent map approach. *Technovation* 29, 481–497.
- [15] Lewis, D.D., Jones, K.S., 1996. Natural language processing for information retrieval. *Communications of the ACM* 39, 92–101.
- [16] Lovins, J.B., 1968. Development of a stemming algorithm. *Mech. Transl. Comput. Linguistics* 11, 22–31.
- [17] Manning, C., Schütze, H., 1999. *Foundations of statistical natural language processing*. MIT press.
- [18] Mikolov, T., Sutskever, I., Chen, K., Corrado, G., Dean, J., 2013. Distributed representations of words and phrases and their compositionality, in: *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*, Curran Associates Inc., Red Hook, NY, USA. p. 3111–3119.
- [19] Prince, G., 1970. The practice of creativity : a manual for dynamic group problem solving. undefined URL: [/paper/The-practice-of-creativity-%3A-a-manual-for-dynamic-Prince/32e17319c36c26b56a5e3aeb24044052eb44763d](https://paperkit.net/paper/The-practice-of-creativity-%3A-a-manual-for-dynamic-Prince/32e17319c36c26b56a5e3aeb24044052eb44763d).
- [20] Rand, W.M., 1971. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association* 66, 846–850. URL: <http://www.jstor.org/stable/2284239>.
- [21] Rosenberg, A., Hirschberg, J., 2007. V-measure: A conditional entropy-based external cluster evaluation measure, in: *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pp. 410–420.
- [22] Rousselot, F., Zanni-Merk, C., Cavallucci, D., 2012. Towards a formal definition of contradiction in inventive design. *Comput. Ind.* 63, 231–242. URL: <https://doi.org/10.1016/j.compind.2012.01.001>, doi:[10.1016/j.compind.2012.01.001](https://doi.org/10.1016/j.compind.2012.01.001).
- [23] Salton, G., Buckley, C., 1988. Term-weighting approaches in automatic text retrieval. *Information Processing and Management* 24, 513–523. URL: <https://www.sciencedirect.com/science/article/pii/0306457388900210>, doi:[https://doi.org/10.1016/0306-4573\(88\)90021-0](https://doi.org/10.1016/0306-4573(88)90021-0).
- [24] Singhal, A., 2001. Modern information retrieval: A brief overview. *IEEE Data Eng. Bull.* 24, 35–43. URL: <http://dblp.uni-trier.de/db/journals/debu/debu24.html#Singhal01>.
- [25] Tseng, Y.H., Lin, C.J., Lin, Y.I., 2007a. Text mining techniques for patent analysis. *Inf. Process. Manag.* 43, 1216–1247.
- [26] Tseng, Y.H., Wang, Y.M., Lin, Y.I., Lin, C.J., Juang, D.W., 2007b. Patent surrogate extraction and evaluation in the context of patent mapping. *Journal of Information Science* 33, 718 – 736.
- [27] Yoon, B., Phaal, R., 2013. Structuring technological information for technology roadmapping: data mining approach. *Technology Analysis & Strategic Management* 25, 1119 – 1137.
- [28] Yoon, B., Yoon, C., Park, Y., 2002. On the development and application of a self-organizing feature map-based patent map. *R & D Management* 32, 291–300.
- [29] Zlotin, B., Zusman, A., Kaplan, L.C., Visnepolschi, S., Proseanic, V., Malkin, S., 2005. *Triz beyond technology: The theory and practice of applying triz to non-technical areas*.