



Discovering Associations Among Technologies Using Modified Term Frequency Inverse Document Frequency

Mohammadhadi Alaeiyan, Mehdi Alaeiyan* and Abolfazl Salemi

ABSTRACT: In the current era of rapid technological advancement, the generation of innovative and efficient ideas requires a thorough understanding of the interrelationships among various technologies. Owing to the vast volume of data embedded within technological domains, automated methods are essential for uncovering meaningful associations. This paper introduces a methodology based on a generalized Term Frequency-Inverse Document Frequency (TFIDF) model to extract associations among technologies from a diverse corpus of textual sources, including scholarly publications. The resulting associations are represented as a weighted graph, where nodes denote individual technologies and edge weights reflect the strength of their co-occurrence. This structure, termed the “association tech-graph,” serves as a valuable tool for analyzing trends and guiding innovation within industrial sectors. By adjusting model parameters, multiple graph variants can be generated, allowing deeper analytical insights. The findings suggest that combinations such as Aerial Robotics with Advanced Driver Assistance enhance autonomy, while Actuators with Adjustable Hoisting Machines improve operational efficiency in heavy-duty systems.

Key Words: Text mining, weighted graphs, technology association, data-driven innovation, TFIDF.

Contents

1 Introduction	1
1.1 The rest of this paper	2
2 Related Works	2
3 Proposed method	3
3.1 Web Crawling	4
3.2 Cleaning and Tokenization	5
3.3 Aggregation	5
3.4 Term Frequency	5
3.5 Joining	6
3.6 TFIDF	7
3.7 Distance/Cosine Similarity	7
4 Graph Analysis	8
5 Concluding remarks	9

1. Introduction

Recent decades have accelerated the advancement of knowledge and technology. Each domain of science and technology now encompasses numerous subsections, replete with data, techniques, formulas, and a plethora of content. The sheer volume of this information presents significant challenges for manual management, even within a specialized field. Concurrently, a prevalent approach to fostering new technological developments or scientific breakthroughs is the amalgamation of existing technologies. Emerging technologies frequently represent a synthesis of innovative methods and techniques drawn from diverse disciplines, as highlighted in the literature [1]. Consequently, for the conception of an idea, the devising of an innovative method, or the development of a new technology, a robust understanding of related technologies is essential. It is within this context that the association between technologies assumes critical importance [2,3].

* Corresponding author.

2010 *Mathematics Subject Classification*: 05C90, 68T50.

Submitted July 08, 2025. Published September 22, 2025

The aim of this paper is to assign an accurate vector to each technology, thereby deriving the associations among these technologies through the relationships between the vectors. An extensive corpus of textual data is required to construct these vectors, leading to the amassing of a considerable number of articles and patents. Many words within these documents are extraneous to the pertinent technology, necessitating their removal. The documents are primed for vector extraction following a thorough cleansing and tokenization process. For this purpose, we introduce and apply a generalized TFIDF model. Subsequently, the relationships among the vectors were discerned by calculating their distances or the cosine of the angles between them. These relationships are ultimately represented as a weighted graph.

Contributions. Among the novel contributions of our work are the following:

1. Assigning a real vector to each technology to discern their interrelations based on vector analysis (See Section 3).
2. Employing a generalized TFIDF model with dual parameters to streamline computations and enhance precision (Described in Subsection 3.6).
3. Representing these relationships through a weighted graph (presented in Subsection 3.7).

1.1. The rest of this paper

The rest of this article is organized as follows: Section 2 provides a review of related work. Section 3 introduces the proposed method, gives an overview of the study, and presents the basic definition and meaning of TFIDF. Subsection 3.1 explains how data related to different technologies and their volumes were collected. Subsection 3.2 describes how irrelevant words and phrases were removed and how the remaining data were preprocessed. The steps for calculating TFIDF on large datasets are covered in Subsections 3.3, 3.4, and 3.5. Section 4 shows a graph where each node represents a specific technology and the weighted edges show how strongly they are related. Section 3.6 further explains the detailed process of computing TFIDF. Finally, Section 5 concludes with a concise summary.

2. Related Works

Azimi et al. [4] proposed a neural network-based model for discerning technological associations; however, their study was confined to a mere 19 technologies. The concept of tech-mining for associations was introduced in [5], yet it lacked a sufficiently empirical methodology and comprehensive quantitative results. A method for identifying relationships between technology and products via keyword analysis of documents was outlined in [6], even though keywords alone do not provide a comprehensive analysis. Furthermore, Mahgoub [7] introduced a text mining approach, termed Extracting Association Rules from Text (EART), which amalgamates XML technology with the TFIDF Information Retrieval (IR) scheme to pinpoint distinctive keywords or features, subsequently employing data mining techniques to unearth association rules. While combining IR schemes like TFIDF with data mining techniques such as association rule mining is a burgeoning technique for extracting knowledge from text, the focus remains solely on keyword sets during the association rule mining phase. This method was applied to web page news documents concerning the bird flu outbreak, and the efficacy of the EART system was benchmarked against related works utilizing the Apriori algorithm, showcasing superior execution time. Alaeiyan et al. [8] proposed machine learning methods to compute a graph's girth, maximum clique number, and maximum independent set number. Random graph generation techniques were used to create numerous graph instances. Fourteen features were then presented and employed for training and testing in the classification and prediction of these graph properties [9].

In contrast, the present research encompasses a substantially larger set of technologies, totaling 198. The analysis extends beyond mere keywords and abstracts to include all words within the articles, yielding numerical results that facilitate the straightforward determination of technological associations. Moreover, the innovative utilization of free parameters within the generalized TFIDF method enables enhanced precision of the resultant findings, tailored to specific requirements. Regarding future work, considering that the title and keywords of articles hold greater significance than other sections in uncovering technological relationships, an area for further investigation could involve assigning a coefficient to terms found in titles and keywords during TFIDF calculation.

3. Proposed method

This section provides an overview and details the methodology for preparing the technology association graph (tech-graph).

This study examines the interrelationships among various technologies within an industry and uses texts and documents related to each technology as the primary data source. Figure 1 depicts the block diagram of the proposed methodology. Data retrieval was performed via the web, employing the Scopus API¹, as explicated in Subsection 3.1.

Consider the set of words that are frequently repeated within industry I . This set is called the basic words of industry I , denoted by the symbol $BW(I)$.

$$BW(I) := \{w_1, \dots, w_n\},$$

where the w_i 's are the most repeated words of industry I . Let T be a technology in industry I . We define the characteristic vector(CV) of T as follows:

$$CV(T) := \{F_T(w_1), \dots, F_T(w_n)\},$$

where F_T can be any real function.

$$\begin{aligned} F_T : BW(I) &\longrightarrow \mathbb{R} \\ w_i &\longmapsto F_T(w_i). \end{aligned}$$

Utilizing the abovementioned definitions, each technology is allocated a vector comprising real elements, termed the 'Characteristic Vector.' Subsequently, a metric for their association can be ascertained by analyzing these vectors. A key challenge, however, arises in selecting F . In this study, a generalized form of TFIDF is employed. As depicted in Figure 1, the manuscript contents are sourced from Elsevier, and texts not pertinent to the context are excised. Following this, terms and words are individually isolated, purified, and tokenized, as delineated in Section 3.2.

To begin, let us delve into the concept of Term Frequency-Inverse Document Frequency (TFIDF) and its role in evaluating word importance within documents. TFIDF is a widely recognized method for assessing the significance of words in a document collection [10,11,12]. Although various definitions exist, the most prevalent one has been adopted for this discussion. Let t represent a term or word, d denote a document (such as a book, article, or magazine), and N signify the total number of documents. Term Frequency (TF) and Document Frequency (DF) are defined as follows:

$$\begin{aligned} TF(t, d) &:= \frac{\text{count of } t \text{ in } d}{\text{number of terms in } d}, \\ DF(t, N) &:= \text{number of occurrence of } t \text{ in } N \text{ documents.} \end{aligned}$$

It is widely acknowledged that a word or phrase holds greater significance for a document if it is frequently repeated within that document and less so in others. Consequently, the term t is deemed more important for the document d when the value of $TF(t, d)$ is high and the value of $DF(t, N)$ is low. For practical purposes, IDF is employed in lieu of DF , which is defined as follows:

$$IDF(t) := \text{Log}\left(\frac{N}{DF + 1}\right).$$

In this definition, the logarithm is employed to facilitate measurement on a more manageable scale. Ultimately, TFIDF is formally defined as:

$$TFIDF(t, d) := TF(t, d) * IDF(t).$$

By definition, the importance of a term t within a document d is directly proportional to its $TFIDF(t, d)$ value [13]; a higher value signifies greater significance.

¹ https://dev.elsevier.com/api_docs.html

Subsequently, a modified TFIDF is computed, which integrates two specific parameters and yields a series of TFIDF tables (further detailed in Section 3.6). Additionally, two distinct metrics are employed to assess the association among technologies: Distance and Cosine similarity, both elaborated upon in Section 3.7. The application of these two metrics results in the generation of two separate sets of tables, facilitating comprehensive comparative analysis.

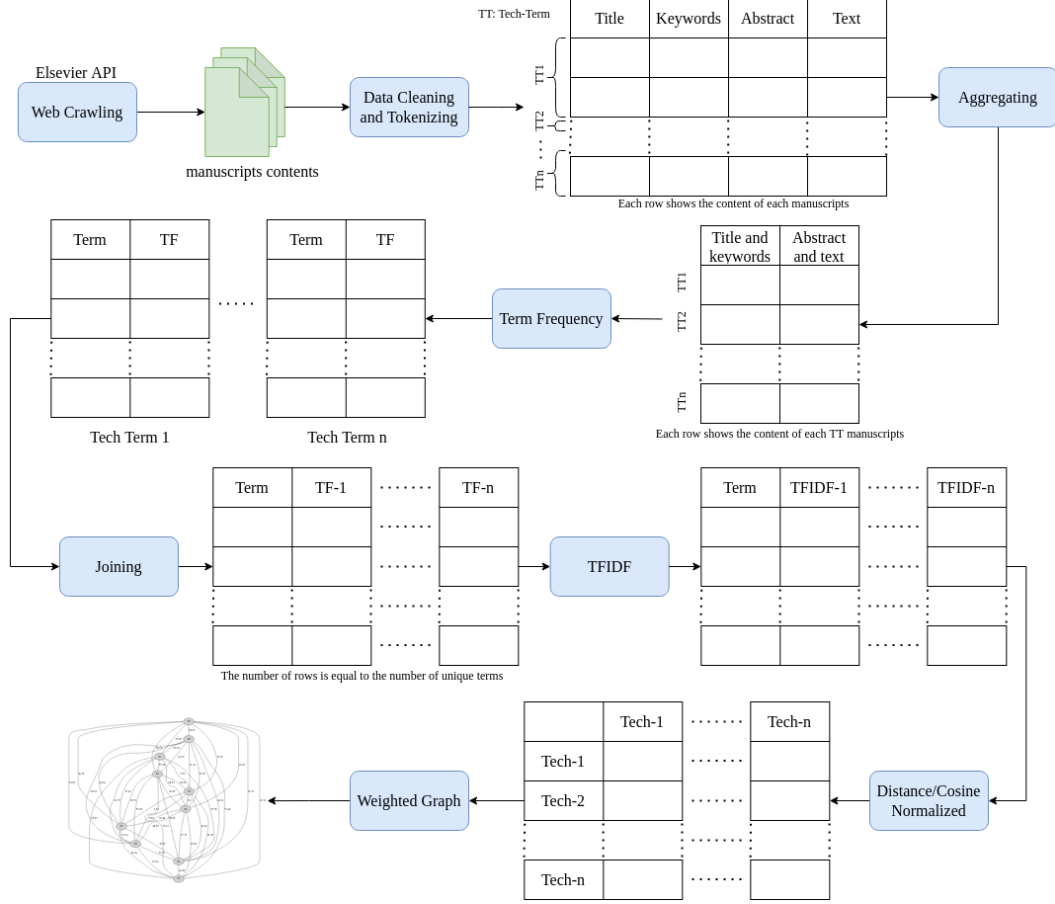


Figure 1: General block diagram of proposed approach.

3.1. Web Crawling

To effectively elucidate the relationships among technologies within an industry, it is essential to analyze a substantial volume of textual data pertinent to each technology. The internet serves as the most accessible source for acquiring such information. Initially, attempts were made to download relevant data directly from various scientific journal websites, including *ScienceDirect*² and *Nature*³, alongside leveraging Google search results. Unfortunately, this approach yielded limited success, primarily due to platform restrictions that hindered the acquisition of a sufficient quantity of documents and texts.

To overcome these limitations and secure a substantial volume of textual data from reputable scientific journals, the Scopus API was utilized. Through the execution of Python scripts, between two and five thousand articles were amassed for each technology, culminating in approximately 21GB of text, which predominantly encompassed abstracts and introductions. Empirical evidence suggests this amount of data proved sufficient for analytical purposes. The collected data was systematically stored in a database

² <https://www.sciencedirect.com>

³ <https://www.nature.com>

No	symbol	text symbol
1	C++	c-plus-plus
2	C#	c-sharp
3	.Net	dot-net
4	F#	F-sharp

Table 1: Symbols

using the .db format. As depicted in Figure 1, this database archives technology terms and manuscript content within two distinct tables. The first table, dedicated to technology terms, features two columns (**id** and **term**) and contains two hundred rows. The second table, housing the downloaded articles, is structured with six columns: **id**, **tech-term ID** (corresponding to the **id** in the technology terms table), **title**, **keywords**, **abstract**, and **main text**. Given the voluminous nature of the data, it was partitioned into four separate tables, each encapsulating articles pertinent to approximately fifty technologies. These tables collectively contain 846,497 rows, representing the total number of articles collected .

3.2. Cleaning and Tokenization

The dataset procured comprised numerous extraneous elements, such as editorial information about publications and authors, which were extrinsic to the core document text and thus necessitated removal. Additionally, words and phrases serving purely grammatical functions, like prepositions and auxiliary verbs, were deemed irrelevant to our objective of discerning semantic connections between technologies. Efforts also included excluding generic expressions unrelated to the specific technology being analyzed, including numerical data and locational or temporal indicators. Beyond the standard stop words identified by Python libraries such as *nltk*⁴ and *spacy*⁵, specific non-essential terms pertinent to our research were excised, a selection of which is presented in Table 2. The procedural steps of data cleansing are outlined in Table 3. After cleansing, the text underwent tokenization, transforming it into an array of terms. These tokens, each corresponding to a particular technology, were cataloged in a table featuring four columns encompassing three to five thousand rows. The aggregate count of tokens post-cleansing stood at 1,243,889,278. Table 1 enumerates the symbols employed during this phase.

3.3. Aggregation

Consolidating all pertinent data into a singular repository is essential to facilitate the TFIDF computation process. This involves amalgamating data associated with each specific technology and unifying all technologies with their respective datasets. For each technology, tokens, originally spanning three to five thousand rows, are consolidated into a single row and then bifurcated into two distinct columns: 'title-keywords' and 'abstract-text.' Subsequently, four individual tables, each representing data for approximately fifty technologies, are concatenated to form a comprehensive single table. This results in a comprehensive table comprising 198 rows and two columns, with each row encapsulating the cleansed and tokenized data pertinent to a single technology (refer to Figure 1). Specifically, the first column consolidates all titles and keywords, while the second column encompasses all abstracts and the full texts of the articles.

3.4. Term Frequency

For the computation of TFIDF values, the frequency of each term within the dataset relevant to each technology (as represented by each row in the output table from the preceding step) is tallied and recorded in separate tables. This process consequently generates 198 tables—one for each technology under consideration—each containing two columns: 'term' and 'TF'. The number of rows in each table corresponds to the count of unique tokens associated with its respective technology.

⁴ <https://www.nltk.org>

⁵ <https://spacy.io>

No	Category	removed word
1	Prepositions	aboard, about, above, across, after, against, along, amid, among, anti, around, as, at, before, behind, below, beneath, beside, besides, between, beyond, but, by, concerning, considering, despite, down, during, except, excepting, excluding, following, for, from, in, inside, into, like, minus, near, of, off, on, onto, opposite, outside, over, past, per, plus, regarding, round, save, since, than, through, to, toward, towards, under, underneath, unlike, until, up, upon, versus, via, with, within, without
2	Verbs	be, am, is, are, was, were, may, shall, should, will, would, can, could, have, has, do, does, doesn't
3	Whs	what, where, who, whom, whose, when, which
4	Suffix	.jpg, .jpeg, .png, .svg, .pdf
5	Letter numbers	one, two, three, four, five, six, seven, eight, nine, ten, eleven, twelve, thirteen, fourteen, fifteen, sixteen, seventeen, eighteen, nineteen, twenty, thirty, forty, fifty, sixty, seventy, eighty, ninety, hundred
6	Date Time	time, now, yesterday, tomorrow, day, days, week, weeks, month, months, year, years

Table 2: Stop Words

No	step title	description or pseudo code
1	Change to lower case	for each term in text term = lowercase(term)
2	Replace symbols with text	for each symbol, text_symbol in ValidSymbolsTable replace(symbol, text_symbol) E.g.: c++ —¿ c-plus-plus
3	Join words of tech term	for each tech_term in DB replace(tech_term, joined_tech_term) E.g.:art robotics —¿art-robotics
4	Removing URLs	removing any url in text
5	Removing Emails	removing any text like something@something
6	Removing Apostrophe	E.g: removing n't
7	Removing any non-alphabetic	for each term in text if term != [A-Z][a-z][0-9] term = ""
8	Removing extra space	
9	Removing extra dash	
10	Removing stop words	for each term in text if StopWordSet contain term delete term

Table 3: Cleaning Steps

3.5. Joining

Given the prevalence of shared terms across the 198 individual tables and the necessity of amalgamating all terms into a singular structure for IDF computation, these tables have been consolidated. This integration yields a comprehensive table featuring 198 columns, with each column representing the term frequency corresponding to a specific technology. The number of rows in this consolidated table is commensurate with the total count of unique terms identified across all articles.

3.6. TFIDF

Our preliminary objective was to allocate a vector of real numbers to each technology, termed the 'characteristic vector' of technology, and to deduce the association between technologies predicated upon these vectors. This aim was realized through the computation of TFIDF. In the process of calculating TFIDF, we encountered two primary challenges.

The initial challenge encountered was the sheer volume of individual words across all articles, amounting to several million, which is reflected in the output table from the previous step. This vast number results in an exceedingly lengthy characteristic vector for each technology, complicating inter-technological relationships' computation. To address this issue, we introduced a parameter, `min_sum_freq`, assigned an arbitrary value, and proceeded to eliminate terms whose cumulative frequency of occurrence *TF* across all technologies fell below the `min_sum_freq` threshold. This method effectively expunges numerous non-essential terms, diminishing the overall term count and, consequently, the length of each technology's characteristic vector.

The second challenge pertains to the calculation of *IDF*. Traditionally, the frequency of a term's occurrence in other technologies is deemed irrelevant; only its presence or absence is considered. However, given the extensive corpus of articles per technology, numerous terms are present at least once in the datasets of all technologies, rendering their *IDF* and, by extension, their *TFIDF* values null. This results in a preponderance of zero-valued elements within the characteristic vectors of technologies. To rectify this, we revised the *IDF* computation to account for the frequency of terms across different technologies. We introduced another parameter, `min_doc_freq`, with an arbitrary range from 0 to 100. For the calculation of a term *t*'s *IDF* from technology *F* with the `min_doc_freq` parameter, we adopt the following approach. If the frequency of term *t* in technology exceeds `min_doc_freq` percent of term *t*'s frequency in technology *F*, we increment the *DF* by one.

As both parameters are variable, an array of distinct *TFIDF* tables can be generated. Each table encapsulates the term and *TFIDF* values corresponding to a specific technology.

A lower `min_sum_freq` value results in fewer terms being discarded, which may enhance the precision of the association but at the cost of increased computational demand. A zero value implies no terms are omitted, aligning with the original *TFIDF* methodology. Conversely, a higher `min_doc_freq` value improves the accuracy of the outcomes, albeit with escalated computational complexity. Opting for a zero value equates to the conventional *TFIDF* calculation.

3.7. Distance/Cosine Similarity

In the preceding phase, the TFIDF for each technology was computed. This has yielded 198 characteristic vectors, each symbolizing a distinct technology, thereby enabling the ascertainment of associations between these technologies through vector comparison. Two metrics were adopted to evaluate these inter-vector relationships: the Euclidean distance and the cosine similarity. The application of these metrics to each output table from the prior step resulted in two distinct matrices, each comprising 198 rows and 198 columns. The cell values within these matrices reflect the degree of association between the technologies corresponding to their respective row and column. A higher value indicates a more substantial relation when employing the distance metric, whereas a lower value signifies a closer association under the cosine metric. To facilitate comparison and express these associations in percentage terms, the values were normalized using the following formula:

$$\bar{x} = \frac{x - MinValue}{MaxValue - MinValue} \times 100.$$

The culmination of the analysis is represented as a weighted graph. These graphs were constructed using GraphvizOnline⁶. Specifically, Figure 2 illustrates a graph generated from a table defined by the parameters `min_sum_freq = 2000` and `min_doc_freq = 20`, applying the distance criterion for a set of ten technologies.

⁶ <https://dreampuf.github.io/GraphvizOnline>

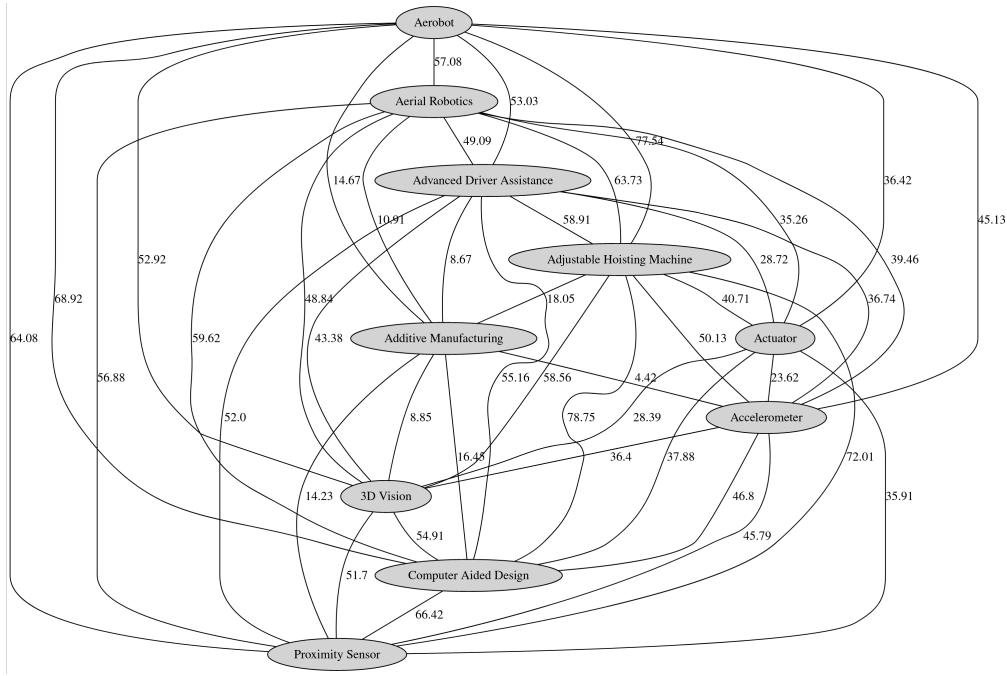


Figure 2: The weighted graph of 10 technologies with $\text{min_sum_freq} = 2000$, $\text{min_doc_freq} = 20$ and distance criteria.

4. Graph Analysis

As illustrated in Figure 2, the graph's vertices represent distinct technologies, while its weighted edges quantify the degree of association between these technologies. The key observations derived from this 10-vertex graph are as follows:

1. **Interconnectivity and Relationship Strength:** The graph vividly illustrates the interconnectivity among various technologies, where nodes denote individual technologies and weighted edges quantify their relational closeness. Higher edge values signify stronger or more significant relationships.
2. **Significant Correlations and Their Implications:** Certain edges exhibit high values, indicative of strong correlations. For instance, 'Aerobot' and 'Aerial Robotics' share a close relationship (57.08), suggesting that progress in one will significantly influence the other, enhancing applications in logistics and surveillance. The convergence of 'Aerial Robotics' with 'Advanced Driver Assistance' points to future innovations in autonomous vehicles, enabling seamless navigation across ground and aerial domains and fostering multi-modal transport solutions. The strong link between 'Additive Manufacturing' and 'Adjustable Hoisting Machines' (55.16) implies a transformative shift in manufacturing, potentially leading to fully automated, on-site construction systems. Furthermore, the integration of 'Accelerometers' with 'Actuators' foresees a future of real-time adaptive responsive systems, expanding applications in robotics and wearable technology. The connection between 'Proximity Sensors' and 'Computer-Aided Design (CAD)' (51.7) signals a trend towards embedding sensor technology into design, thereby improving simulation accuracy with real-world data.
3. **Emerging Technologies and Strategic Integrations:** The graph highlights the potential for new technologies through strategic combinations. Emerging sophisticated integrations, such as 'Aerial Robotics' with 'Advanced Driver Assistance,' could lead to advanced autonomous systems capable of navigating complex environments. Similarly, the merging of 'Actuators' with 'Adjustable

Hoisting Machines’ promises improved efficiency and safety in heavy-duty applications, demonstrating their mutual enhancement.

4. **Potential Future Technologies:** Future technological advancements include ‘Integrated Autonomous Systems,’ potentially arising from the combination of aerial technologies with diverse sensors for operation in varied environments. Additionally, ‘Smart Manufacturing Machines’ may emerge from blending additive manufacturing with adjustable hoisting technologies, facilitating autonomous design and assembly of complex structures.
5. **Technological Convergence:** The overarching theme identified is technological convergence, characterized by increasing interaction and shared development goals across various disciplines. This convergence is poised to foster entirely new fields, leveraging the combined strengths of multiple technologies to create innovative solutions.

5. Concluding remarks

This paper discerned the associations between technologies by attributing a vector to each and computing both the distance and the cosine similarity between these vectors. This methodology was applied to 198 technologies within the robotics sector; however, it was found to be versatile enough to be adapted for any quantity of technologies across various industries. The vector calculations were performed using the generalized TFIDF method, which was characterized by two adjustable parameters. The precision of the resultant output was found to be directly proportional to the meticulous selection of these parameter values, albeit at the expense of increased computational complexity and duration. Given the arbitrary nature of parameter selection and the choice between the two evaluative criteria, a multitude of graphs could be generated, each with its unique level of accuracy. A prospective avenue for research was identified as involving the determination of optimal parameter values while balancing computational costs. Since the title and keywords of articles were considered more important than other sections in discovering the relationship between technologies, assigning a coefficient to the terms in the title and keywords during TFIDF calculation was suggested as a future work. Overall, the method was found to be valuable for mapping technological convergence and identifying promising areas for future development, though room for refinement in parameter selection and the integration of broader datasets was acknowledged.

Acknowledgements

The authors sincerely appreciate the Editor and anonymous reviewers for their valuable insights and constructive suggestions, which have significantly improved the clarity and quality of this manuscript.

Competing interests

The authors confirm that they have no competing interests to disclose.

Funding

There are no funding sources available for this study.

Availability of data and materials

No data were generated or analyzed during the course of this research.

References

1. P. Brey, *Ethics of emerging technology*, in *The Ethics of Technologies: Methods and Approaches*, Rowman & Littlefield International, Lanham, MD, USA, 2017.
2. S. Graham, *Software-Sorted Geographies*, *Progress in Human Geography*, vol. 29, pp. 562–580, (2005). doi:10.1191/0309132505ph568oa
3. A. Sharma and G. R. Iyer, *Technology-Driven Business Model Innovation: A Review*, *International Journal of Management Reviews*, vol. 14, pp. 467–484, (2012). doi:10.1111/j.1468-2370.2011.00320.x
4. S. Azimi, H. Veisi, M. Fateh-rad, and R. Rahmani, *Discovering Associations Among Technologies Using Neural Networks for Tech-Mining*, *IEEE Transactions on Engineering Management*, vol. 69, pp. 1394–1404, (2022). doi:10.1109/TEM.2020.2981284

5. B. Qasemizadeh, P. Buitelaar, and F. Monaghan, *Developing a Dataset for Technology Structure Mining*, in *Proc. 2010 IEEE Fourth Int. Conf. Semantic Computing*, Pittsburgh, PA, USA, pp. 32–39, (2010). [doi:10.1109/ICSC.2010.73](https://doi.org/10.1109/ICSC.2010.73)
6. B. Yoon, I. Park, and B. Coh, *Exploring technological opportunities by linking technology and products: Application of morphology analysis and text mining*, *Technological Forecasting and Social Change*, vol. 86, pp. 287–303, (2014). [doi:10.1016/j.techfore.2013.10.013](https://doi.org/10.1016/j.techfore.2013.10.013)
7. H. Mahgoub, D. Rosner, N. Ismail, and F. Torkey, *A Text Mining Technique Using Association Rules Extraction*, in *Proc. World Academy of Science, Engineering and Technology*, vol. 2, June 2008.
8. M. Alaeiyan, M. Alaeiyan, and K. K. Obayes, *Machine Learning Predicts Graph Properties: Clique, Girth, and Independent Numbers*, *Discrete Mathematics, Algorithms and Applications*, vol. 2024, p. 2450113, (2024). [doi:10.1142/S1793830924501136](https://doi.org/10.1142/S1793830924501136)
9. M. Alaeiyan, *Characteristics and eigenvalues of the newly defined Ala graph*, *Physica Scripta*, vol. 100, no. 5, p. 055201, (2025). [doi:10.1088/1402-4896/ad42f6](https://doi.org/10.1088/1402-4896/ad42f6)
10. C. Zhai and S. Massung, *Text Data Management and Analysis: A Practical Introduction to Text Mining*, Morgan and Claypool Publishers, (2016). [doi:10.2200/S00721ED1V01Y201505HLT027](https://doi.org/10.2200/S00721ED1V01Y201505HLT027)
11. Y. Zhang and L. Zhao, *A Modified TF-IDF Method for Text Classification*, *International Journal of Computer Applications*, vol. 181, pp. 18–24, (2018). [doi:10.5120/ijca2018917395](https://doi.org/10.5120/ijca2018917395)
12. S. Hussain and N. Chaudhry, *Enhanced TF-IDF Vectorization for Text Mining*, *Journal of Computer and Communications*, vol. 8, pp. 52–60, (2020). [doi:10.4236/jcc.2020.87005](https://doi.org/10.4236/jcc.2020.87005)
13. C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*, MIT Press, (2008). [doi:10.1017/CBO9780511809071](https://doi.org/10.1017/CBO9780511809071)

Mohammadhadi Alaeiyan,
 Faculty of Computer Engineering,
 K. N. Toosi University of Technology,
 Seyed Khandan, Shariati Ave, Tehran 16317-14191, Tehran, Iran.
 E-mail address: m.alaeiyan@kntu.ac.ir

and

Mehdi Alaeiyan and Abolfazl Salemi,
 School of Mathematics and Computer Science,
 Iran University of Science and Technology,
 Narmak, Tehran 16846, Tehran, Iran.
 E-mail address: alaeiyan@iust.ac.ir
 E-mail address: abolfazl.salemi@gmail.com