

# Citation Analysis on Scientific Articles Using Cosine Similarity

Contents:

- Artikel
- Cover Prosiding
- Editorial
- Daftar isi
- Peer Review
- Korespondensi
- Similarity

# Citation Analysis on Scientific Articles Using Cosine Similarity

Ulfa Mardatillah

Department of Informatics  
UN Sunan Gunung Djati  
Bandung, Indonesia  
ulfamardatillah@gmail.com

Wildan Budiawan Zulfikar

Department of Informatics  
UN Sunan Gunung Djati  
Bandung, Indonesia  
wildan.b@uinsgd.ac.id

Aldy Rialdy Atmadja

Department of Informatics  
UN Sunan Gunung Djati  
Bandung, Indonesia  
aldy@if.uinsgd.ac.id

Ichsan Taufik

Department of Informatics  
UN Sunan Gunung Djati  
Bandung, Indonesia  
ichsan@uinsgd.ac.id

Wisnu Uriawan

Department Informatics and  
Mathematics  
INSA, Lyon, France  
wisnu\_u@uinsgd.ac.id

**Abstract**— *Citations have a vital role in scientific articles. Every sentence that is listed must be obtained from the reference. There is a measure that is similarity. The problem is, the likeness of a sentence in a scientific article to the source it refers to can have different meanings. The high level of similarity can undoubtedly increase the similarity value if it is checked using the plagiarism (similarity checker) tool, where the higher the value is more similar to other articles on the contrary with low similarity remaining. However, a low level of similarity can also mean that a statement in a scientific paper has referred to a completely irrelevant source. This research aims to analyze the level of similarity between scientific articles and other articles or sources it refers to. Checking the input text starts with text preprocessing, weighting the TF-IDF, and determining the similarity level using Cosine Similarity. Based on the test results, it is found that the results in a scientific article document have a high level of similarity with the article it refers to, and the other 30% have a low level of similarity, or in other words, there is no relationship.*

**Keywords**—*scientific articles, journals, TF-IDF, cosine, text mining, citation, incorrect citation*

## I. INTRODUCTION

The development of information technology is increasing and has a positive impact. One of the positive impacts of information technology development is the ease of exchanging information. This convenience is often misused by someone or several people in completing work. The misused usually occurs in Scientific Documents. Scientific documents are documents of reasoning and research results using a variety of objects and methods. A good document has explicit references and sources as a reference, for example, a scientific paper or journal that has explicit references and citations and does not misuse the original citation of the document [1]–[4].

Scientific journal manuscripts are the author's scientific research manuscripts and consist of direct and indirect quotations. Direct quotations are quotations that are written exactly like the source, both in language and spelling [5], [6]. An indirect quotation is a quotation that is not the same as the original, and the citer only takes the main idea of the quoted source to be restated with the sentence compiled by the quote. Quotations are considered important as valid

references and are included in the writing of scientific papers [2], [3], [7], [8].

Mistakes that harm scientific journal scriptwriters are called plagiarism. Based on previous research that uses a cosine similarity method to check for plagiarism only, this research complements and refines it by weighing each word in scientific documents and testing the text-similarity value. So, documents that have a high similarity value are abused quotes and a high similarity value. The low similarity is a document that has been cited with little abuse or towards correctness.

The website-based bibliographic reference citation checking system for scientific journals will be used in checking scientific documents. Then, this study uses the TF-IDF weighting method, which functions to determine the weighting value of the validity of scientific journal manuscripts, and the Cosine Similarity Algorithm, which looks for word equations according to scientific journal manuscripts.

The basic principle of the method used in this research is Text Preprocessing using Case Folding, Tokenizing, Filtering, Stemming. Furthermore, the value of word weights in scientific documents, namely TF-IDF, is carried out after the stemming. Also, stopword removal stages are calculated on the value or weight of a word (term) in the document. After obtaining the word weight value, then testing the level of similarity using Cosine similarity, which compares the similarity between documents. In this case, what is being compared is a query with a training document. In calculating cosine similarity, the first thing to do is do a scalar multiplication between the query and the document, add it up, multiply the length of the document and the square of the query length, and then calculate the square root. Furthermore, the scalar multiplication result is divided by the multiplication result of the length of the document and the query. In related works, this method is widely used and shows a good performance in several cases [9]–[15]. Also, it possible to add any improvement to increase its accuracy and performance [16]–[19].

## II. METHODOLOGY

### A. Text Processing

Text mining or by other names such as intelligent text analysis, text data mining, or knowledge discovery in a text can simply be interpreted as the process of finding patterns that were not previously seen in certain text documents or sources.

The initial step of text mining is preprocessing. Some of them are case folding or changing text to lowercase and without punctuation or special characters, filtering, namely by removing conjunctions, conjunctions, tokenizing the separation of a document or sentence into each word or term, and Stemming or separating affixes in words into basic word [15], [20].

TABLE I. CASE FOLDING

Document	Result
D1	ada beberapa penelitian terdahulu yang memanfaatkan mikrokontroler arduino uno dan sensor seperti untuk pengontrolan suhu dan ruang
D2	perancangan dan implementasi pengontrol suhu ruangan berbasis mikrokotoller arduino uno dengan sensor suhu lm3s layak digunakan dan diaplikasikan sebagai sistem pengontrol suhu ruangan

The next step is tokenizing. This process converts sentences into words with spaces as the separating feature as described in Table 2.

TABLE II. TOKENIZING

Document 1 (D1)	Document 2 (D2)
ada	perancangan
beberapa	dan
penelitian	implementasi
terdahulu	pengontrol
yang	suhu
memanfaatkan	ruangan
mikrokontroler	berbasis
arduino	mikrokontroler
uno	arduino
dan	uno
sensor	dengan
seperti	sensor
untuk	suhu
pengontrolan	lm3s
suhu	layak
ruang	digunakan
	dan
	diaplikasikan
	sebagai
	sistem

Filtering is removing unnecessary words in the input text. So, the filtered text will be easily processed at a later stage. This study uses 758 words stored in an external dataset. After that is stemming. This step changes the filtered word into a root word and removes the affixes. The results of the filtering and stemming processes can be seen in table 3.

TABLE III. D1 AND D2 AFTER FILTERING AND STEMMING

Document 1 (D1)	Document 2 (D2)
teliti	rancang
manfaat	implementasi
mikrokontroler	control
arduino	suhu
uno	ruang
sensor	basis

control	mikrokontroler
Suhu	arduino
Ruang	uno
	sensor
	suhu
	lm3s
	guna
	Aplikasi
	Sistem

### B. TF-IDF and Cosine Similarity

TF is obtained from the number of words that appear in a text, so you can find out how many times the word appears. And after this TF calculation will be continued with the calculation of DF. Table 4 is the process of calculating the weight value to determine the number of words that appear in the citation text and the source of the quote so that it can be calculated by calculating the cosine similarity [21], [22].

TABLE IV. TERM FREQUENCY

No	Word	Term Frequency (TF)	
		Document 1 (D1)	Document 2 (D2)
1	Teliti	1	0
2	Manfaat	1	0
3	mikrokontroler	1	1
4	Arduino	1	1
5	uno	1	1
6	sensor	1	1
7	control	1	2
8	suhu	1	3
9	ruangan	1	2
10	rancang	0	1
11	implementasi	0	1
12	basis	0	1
13	lm3s	0	1
14	guna	0	1
15	aplikasi	0	1
16	sistem	0	1

Then calculate the value of DF (document frequency), which is the number of words that appear in the two documents tested as described in table 5.

TABLE V. DOCUMENT FREQUENCY

No	Word	Document Frequency
1	teliti	1
2	manfaat	1
3	mikrokontroler	2
4	arduino	2
5	uno	2
6	sensor	2
7	control	3
8	suhu	4
9	ruangan	3
10	rancang	1
11	implementasi	1
12	basis	1
13	lm3s	1
14	guna	1
15	aplikasi	1
16	sistem	1

In table 6 are the results of the IDF (*inverse document frequency*) which is the value of the word analysis using (1) which is the value of the times with the TF (*term frequency*) in order to get the weight value.

$$1 + \log(N/df + 1) \quad (1)$$

TABLE VI. IDF RESULT

Word	DF	IDF 1	IDF 2
Teliti	1	1	0
Manfaat	1	1	0
Mikrokontroller	2	0.5	0.5
Arduino	2	0.5	0.5
Uno	2	0.5	0.5
Sensor	2	0.5	0.5
Control	3	0.5	0.823908741
Suhu	4	0.5	1.301
Ruangan	3	0.5	0.823908741
Rancang	1	0	1
Implementasi	1	0	1
Basis	1	0	1
lm3s	1	0	1
Guna	1	0	1
Aplikasi	1	0	1
Sistem	1	0	1

After calculating the IDF (inverse document frequency), the next step is to analyze the IDF, the TF value, and the IDF value that have been obtained from the above calculations, then the TFIDF analysis is a calculation to get the weighted value of the quoted text and the source of the citation tested. The following table 7 is a TFIDF analysis table.

TABLE VII. TF-IDF RESULT

No	Word	TF		IDF		TF IDF	
		D1	D2	D1	D2	D1	D2
1	Teliti	1	1	1	1	1	0
2	Manfaat	1	1	1	1	1	0
3	mikrokontroller	1	0.5	0.5	0.5	0.5	0.5
4	Arduino	1	0.5	0.5	0.5	0.5	0.5
5	Uno	1	0.5	0.5	0.5	0.5	0.5
6	Sensor	1	0.5	0.5	0.5	0.5	0.5
7	Control	1	0.5	0.5	0.5	0.5	0.82
8	Suhu	1	0.5	0.5	0.5	0.5	1.30
9	Ruangan	1	0.5	0.5	0.5	0.5	0.82
10	Rancang	0	0	0	0	0	1
11	Implementasi	0	0	0	0	0	1
12	Basis	0	0	0	0	0	1
13	lm3s	0	0	0	0	0	1
14	Guna	0	0	0	0	0	1
15	Aplikasi	0	0	0	0	0	1
16	Sistem	0	0	0	0	0	1

The cosine similarity calculation is described in table 8 where document 1 is the citation text and document 2 is the source text of the citation. The citation text is multiplied by the citation source text, the citation text and the citation source are squared so that the total is obtained to find the text-similarity value.

TABLE VIII. CALCULATION OF COSINE SIMILARITY

Word	TF IDF				
	D2	D2	D1.D2	D1^2	D2^2
Teliti	1	0	0	1	0
Manfaat	1	0	0	1	0
mikrokontroller	0.5	0.5	0.25	0.25	0.25
Arduino	0.5	0.5	0.25	0.25	0.25
Uno	0.5	0.5	0.25	0.25	0.25
Sensor	0.5	0.5	0.25	0.25	0.25
Control	0.5	0.82	0.41	0.25	0.67
Suhu	0.5	1.30	0.65	0.25	1.69
Ruang	0.5	0.82	0.41	0.25	0.67
Rancang	0	1	0	0	1
implementasi	0	1	0	0	1
Basis	0	1	0	0	1
lm3s	0	1	0	0	1
Guna	0	1	0	0	1
Aplikasi	0	1	0	0	1

Sistem	0	1	0	0	1
			2.47	3.75	11.05

$$Crossproduct = 6.13 \quad (2)$$

$$A.B = 2.4744, |A|, |B| = 6.13 \quad (3)$$

$$Cos(A,B) = 0.4036 \quad (4)$$

Based on the results of the manual algorithm calculation above (4), the similarity value in the text is 0.4036. In this study, the level of similarity will be divided into 3 categories, namely high, medium, and low. The low level is in the range of 0-0.3, the medium level is in the range of 0.3-0.6, and the high level is in the range of 0.6-1.0. The result of (4) is 0.4036 so that it is known that the text has a medium level.

### III. RESULT AND DISCUSSION

The system test scenario defines the amount of data used and what processes are carried out in the system to determine the error in a citation of this scientific journal. The process carried out in testing this system is data collection, text preprocessing, weighting, and text similarity Cosine Similarity. The data used in testing this system is an online journal available on the internet by looking for references to the bibliography listed in the document to be tested and inputting the citation text. 10 of the 50 test samples can be seen in the table below.

TABLE IX. TEST RESULT

No.	Citation	Result	Similarity Level
1	T001	0.43352308	Medium
2	T002	0.12427994	Low
3	T003	0.07988254	Low
4	T004	0.16726948	Low
5	T005	0.76533020	High
6	T006	0.08811075	Low
7	T007	0.02087507	Low
8	T008	0.05383300	Low
9	T009	0.16374133	Low
10	T010	0.24542823	Low

In the table above, there is a threshold, that is, if the similarity of the text is 0-0.3, it means that the similarity level is low. If the similarity of the text is 0.3-0.6 then the similarity is medium. If the similarity text results in a value of 0.6-1.0, the similarity level is high as described in the section on methodology.

Based on the 50 samples the data tested has a different similarity value for each text input. The thing that affects the similarity value is the number of word weights tested. In the first data sample, the number of text 1 has 9 words and the input text 2 has 15 words, with the number of word weights that are not too far in the range. In the second input text test, between text 1 and text 2, the number of word weights is not proportional, so it is dominated by text 2. And the similarity value is tested based on the weight of the words that appear in the two input texts. Documents with medium and high levels concerning the suitability of the quotation and the source of the citation are 70%, and the low level is 30%, unequal documents are documents that do not have a match between the quotation and the source of the citation. The

results of data validation checking scientific journal documents are that 70% of the citations in the article are true and 30% of the citations are incorrect.

#### IV. CONCLUSION

Checking scientific journal citation in testing sample data, it is found that the text is not similar. The citation text that does not have a similarity value, the text close to similar is the text that has a value close to similar, and similar is the text that has a similarity value. TF-IDF weighting in testing one citation text with one citation source text that comes from the bibliography, then the weighting greatly affects the value of the similarity of the text. The similarity of Cosine Similarity in testing one citation text with one citation source text that comes from the bibliography. Then the text-similarity value is obtained by the threshold value and knowing the text is included in several categories, namely not similar, close to similar, and similar with similar accuracy and Approaching similar, which is 70% and not 30% similar because of similarity and approaching similar. There is an agreement between the quote and its source and the word weight value that affects it.

#### REFERENCES

- [1] M. Golosovsky, *Citation Analysis and Dynamics of Citation Networks*. Cham: Springer International Publishing, 2019.
- [2] H. F. Moed, "Citation analysis in research evaluation," in *Proceedings of ISSI 2005: 10th International Conference of the International Society for Scientometrics and Informetrics*, 2005, vol. 2, pp. 437–441.
- [3] A. Shabani and R. Saadat, "Examining the citations received by DOAJ's journals from ISI Web of Science's articles (2003-2008)," in *Proceedings - 3rd International Conference on Information Sciences and Interaction Sciences, ICIS 2010*, 2010, pp. 109–113.
- [4] X. Su, A. Prasad, M. Y. Kan, and K. Sugiyama, "Neural multi-task learning for citation function and provenance," in *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries*, 2019, vol. 2019-June, pp. 394–395.
- [5] M. Roman, A. Shahid, S. Khan, A. Koubaa, and L. Yu, "Citation Intent Classification Using Word Embedding," *IEEE Access*, vol. 9, pp. 9982–9995, 2021.
- [6] Y. Yang and Y. Zheng, "The effect of open access journals on citation impact: A citation analysis of open access journals using Google Scholar," in *4th International Conference on Cooperation and Promotion of Information Resources in Science and Technology, COINFO 2009*, 2009, pp. 278–280.
- [7] H. H. Bi, J. Wang, and D. K. J. Lin, "Comprehensive citation index for research networks," *IEEE Trans. Knowl. Data Eng.*, vol. 23, no. 8, pp. 1274–1278, 2011.
- [8] P. A. H. Kazi and M. S. Patwardhan, "Context based citation summary of research articles: A step towards qualitative citation index," in *IEEE International Conference on Computer Communication and Control, IC4 2015*, 2016.
- [9] P. P. Gokul, B. K. Akhil, and K. K. M. Shiva, "Sentence similarity detection in Malayalam language using cosine similarity," in *RTEICT 2017 - 2nd IEEE International Conference on Recent Trends in Electronics, Information and Communication Technology, Proceedings*, 2017, vol. 2018-January, pp. 221–225.
- [10] W. Darmalaksana, C. Slamet, W. Budiawan Zulfikar, F. Fadillah, D. Sa'adillah Maylawati, and H. Ali, "Latent semantic analysis and cosine similarity for hadith search engine," vol. 18, no. 1, pp. 217–227, 2020.
- [11] S. Roopak and T. Thomas, "A novel phishing page detection mechanism using HTML source code comparison and cosine similarity," in *Proceedings - 2014 4th International Conference on Advances in Computing and Communications, ICACC 2014*, 2014, pp. 167–170.
- [12] S. Pattnaik and A. K. Nayak, "Summarization of odia text document using cosine similarity and clustering," in *Proceedings - 2019 International Conference on Applied Machine Learning, ICAML 2019*, 2019, pp. 143–146.
- [13] S. Zhu, J. Wu, and G. Xia, "TOP-K cosine similarity interesting pairs search," in *Proceedings - 2010 7th International Conference on Fuzzy Systems and Knowledge Discovery, FSKD 2010*, 2010, vol. 3, pp. 1479–1483.
- [14] D. Soyusiawaty and Y. Zakaria, "Book data content similarity detector with cosine similarity (case study on digilib.uad.ac.id)," in *Proceeding of 2018 12th International Conference on Telecommunication Systems, Services, and Applications, TSSA 2018*, 2018.
- [15] B. Li and L. Han, "Distance weighted cosine similarity measure for text classification," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2013, vol. 8206 LNCS, pp. 611–618.
- [16] K. C. Kirana, S. Wibawanto, N. Hidayah, and G. P. Cahyono, "The improved artificial neural network based on cosine similarity in facial emotion recognition," in *International Conference on Electrical Engineering, Computer Science and Informatics (EECSI)*, 2019, pp. 45–48.
- [17] P. Sitikhu, K. Pahi, P. Thapa, and S. Shakya, "A Comparison of Semantic Similarity Methods for Maximum Human Interpretability," in *International Conference on Artificial Intelligence for Transforming Business and Society, AITB 2019*, 2019.
- [18] C. Beşiktaş and H. A. Mantar, "Real-time traffic classification based on cosine similarity using sub-application vectors," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2012, vol. 7189 LNCS, pp. 89–92.
- [19] M. Kryszkiewicz, "Determining Cosine Similarity Neighborhoods by Means of the Euclidean Distance," *Intell. Syst. Ref. Libr.*, vol. 43, pp. 323–345, 2013.
- [20] M. Biba and M. Mane, "Sentiment analysis through machine learning: An experimental evaluation for Aslbanian," *Adv. Intell. Syst. Comput.*, vol. 235, pp. 195–203, 2014.
- [21] M. Irfan, Jumadi, W. B. Zulfikar, and Erik, "Implementation of Fuzzy C-Means algorithm and TF-IDF on English journal summary," in *2017 Second International Conference on Informatics and Computing (ICIC)*, 2017, pp. 1–5.
- [22] M. M. Saad, N. Jamil, and R. Hamzah, "Evaluation of Support Vector Machine and Decision Tree for Emotion Recognition of Malay Folklores," *Bull. Electr. Eng. Informatics*, vol. 7, no. 3, pp. 479–486, Sep. 2018.

Proceeding of  
**2021 The 7th International Conference on Wireless and Telematics  
(ICWT2021)**

19-20 August 2021  
Bandung, Indonesia

ISBN:  
978-1-6654-4402-6 (XPLORE COMPLIANT)  
978-1-6654-4401-9 (USB)

## Table of Content

Implementation of C-Band Antenna for Man-Portable Ground Surveillance Radar  
*Wervyan Shalannanda and Tommi Hariyadi*

A smart waste management under a wireless communication system  
*Marwan Marwan, Nurhayati Nurhayati, Nur Aminah, Rusdi Wartapane, Muhammad Syarif and Ibrahim Abduh*

Augmented Reality using Natural Feature Tracking Method to Introduce Science Verses in Qur'an  
*Dian Sa'Adillah Maylawati, Khusnul Khotimah, Diena Rauda Ramdania, Muhammad Ali Ramdhani, Yana Aditia Gerhana, Mohamad Irfan and Rosihon Anwar*

Remote Medical Check-up Kit System for Communication Platform in Rural Area  
*Moh Rizki Hidayatulloh and Iskandar*

A Deep Learning Approach To Analyze The Sentiment Of Online Game Users  
*Agung Wahana, Wildan Budiawan Zulfikar, Wildan Najah Wildiansyah, Aldy Rialdy Atmadja, Diena Rauda Ramdania and Beki Subaeki*

The Comparation of Steepest Ascent Hill Climbing and A-star for Classic Game  
*Wildan Budiawan Zulfikar, Mohamad Irfan, Rijal Muharom Yahya, Diena Rauda Ramdania and Jumadi Jumadi*

Transducer Array System Design for Underwater Communication Systems Using Audio Frequency  
*Viandra Nurmalita, Ni Nyoman Dheanty Maharani, Ian Joseph Matheus Edward and Iskandar*

Expert System to Detect Social Anxiety in Student using the Rule-based Reasoning Method  
*Cepy Slamet, Nida Maria Ulfah, Dian Sa'Adillah Maylawati, Rachmat Jaenal Abidin, Khaerul Manaf and Rosleny Marliani*

A DApp Architecture for Personal Lending on Blockchain  
*Wisnu Uriawan, Cepy Slamet, Agung Wahana and Vemy Suci Asih*

Text Summarization using TextRank for Knowledge Externalization from Indonesian Online Discussion Forums  
*Yana Aditia Gerhana, Ichsan Taufik, Raka Iqbal Syamsuddin, Undang Syaripudin, Dian Sa'adillah Maylawati, Titik Khawa Abdul Rahman and Muhammad Ali Ramdhani*

Implementation of Channel Coding System As Error Checking in The Underwater Communication System Using Audio Frequency  
*Ni Nyoman Dheanty Maharani, Viandra Nurmalita, Ian Joseph Matheus Edward and Iskandar*

Topic Clusterization of Indonesian Journal Article based Topic Modeling  
*Ali Rahman, Ana Hadiana, Agung Wahana, Wahyudin Darmalaksana, Muhammad Ali Ramdhani and Rizky Sam Pratama*

Similarity Level Analysis of the Voices of Twins Using the Analysis of Variance and Likelihood Ratio Methods  
*Mohamad Irfan, Diena Rauda Ramdania, Nurul Hasni, Ichsan Budiman, Dian Sa'Adillah Maylawati and Khaerul Manaf*

Citation Analysis on Scientific Articles Using Cosine Similarity

*Ulfa Mardatillah, Wildan Budiawan Zulfikar, Aldy Rialdy Atmadja, Ichsan Taufik and Wisnu Uriawan*

Augmented Reality using Features Accelerated Segment Test to Introduce Batik Garutan

*Nur Lukman, Rizaldi Andriansyah, Ichsan Taufik, Agung Wahana, Dian Sa'Adillah Maylawati, Beki Subaeki and Putri Diesy Fitriani*

Making Augmented Reality Learning Media In Stereochemistry of Carbohydrate Concepts Oriented Multiple Representations

*Ilham Fitrotul Hayat, Ferli Septi Irwansyah and Neneng Windayani*

Development e-module on the concept of reduction-oxidation (redox) oriented towards chemical literacy

*Adi Setya Permana, Ferli Septi Irwansyah and Cucu Zenab Subarkah*

Web Crawler Development to Optimize the Sentiment Analysis Process of Indonesian News Article Opinions

*Ichsan Budiman, Muhammad Deden Firdaus, Muhammad Insan Al-Amin, Eva Nurlatifah and Jumadi M.Cs*

Factors Affecting The Success of ICT Management Between Pesantren-Based Private Islamic Religious College

*Cecep Nurul Alam, Yaya Suryana, Khaerul Manaf, Beki Subaeki, Faiz M Kaffah, Syopiansyah Jaya Putra and Muhammad Indra Nurardy S*

IoT Monitoring System for Solar Power Plan Based on MQTT Publisher/Subscriber Protocol

*Januar Muhamad Ramadhan, Rina Mardiaty and Irsyad Nashirul Haq*

Fuzzy Logic Control for Avoiding Static Obstacle in Autonomous Vehicle Robot

*Fikri Ahmad Fauzi, Edi Mulyana, Rina Mardiaty and Aan Eko Setiawan*

Design of Semi-Autonomous Navigation Robot Using Integrated Remote Control And Fuzzy Logic

*Farhan Wildani, Rina Mardiaty, Edi Mulyana and Aan Eko Setiawan*

Control System for Electrical Conductivity in Aquaponic Cultivation Using Fuzzy Logic

*Rina Yuhasari, Rina Mardiaty, Nanang Ismail and Setia Gumilar*

Design of Automatic Watergate System Using ESP32 Microcontroller Based on Fuzzy Logic Method

*Fadhilah Ramdhani, Rina Mardiaty, Edi Mulyana and Ah Fathonih*

Design and Simulation Transfer Learning on Image Processing for Determining Condition of Robot Based on Neural Network

*Asep Khoerudin, Rina Mardiaty, Edi Mulyana and Aan Eko Setiawan*

Implementation of Fuzzy Logic Control to Maintain pH Content in Aquaponic

*Muhammad Daffa Fadilah, Rina Mardiaty, Nanang Ismail and Ading Kusdiana*

Design of Voltage and Flow Monitoring System for PJU-TS Using the Internet of Things (IoT)

*Agi Sahrul Pradana, Adam Faroqi, Edi Mulyana and Fauzan Ali Rasyid*



Telescop and E-Tara: Online Application For Detecting The Authentic And Fraud Divorce Certificates In Religious Court

*Ramdani Wahyu Sururie, Yoghi Arief Susanto and Lia Kamelia*

The Prototype of Smart Garden Fertigation System with Solar Photovoltaic System Based on IoT  
*Reno Muhammad Fadilla, Nabila Safitri Dwi Oktafiandini, Muhitha Adha, Lia Kamelia and Akmaliyah*

Survey on Hybrid Techniques in The Classification of Nutrient Deficiency Levels in Citrus Leaves  
*Lia Kamelia, Titik Khawa Abdul Rahman, Hoga Saragih and Saepul Uyun*

IoT-Based Battery Monitoring System in Solar Power Plants with Secure Copy Protocol (SCP)  
*Mufid Ridlo Effendi, Ridwan Sidik Al-Falah, Sarbini Sarbini and Nanang Ismail*

Growth Monitoring System and Automatic Watering for Chili Plants Based on Internet of Things  
*Dzulfikri Hanafi, Rina Mardiaty, Toni Prabowo and Cecep Hidayat*

AMC-Based Multiband Microstrip Antenna for Wireless Communication  
*Muhammad Farhan Maulana, Eki Ahmad Zaki Hamidi, Achmad Munir and Nanang Ismail*

Characteristic Performance of L-Band Waveguide BPF Made of Substrate Integrated Structure  
*Aulia Fathoni, Nanang Ismail, Hardi Nusantara and Achmad Munir*

Characterization of SIW-Based 5.8 Ghz Monopulse Antenna Using Type-V Linear Slot  
*Ebin Novendra, Achmad Munir, A. Bachrun Rifa'i and Nanang Ismail*

Design of Humadity Control with Automated Drip Irrigation System Using Node-RED and MQTT on Cactus Plant

*Rachma Dianty, Rina Mardiaty, Edi Mulyana and Aep Kusnawan*

A Systematic Literature Review on IoT-based Smart Grid

*Nike Sartika, Yuda Sukmana, Mufid Ridlo Effendi, Iu Rusliana, Khomisah and Yuyun Yuningsih*

The Development of The Home Electrical Power Consumption System Prototype in Real-Time  
*Yusuf Maulana, Riski Darmawan, Lia Kamelia, Edi Mulyana, Mufid Ridlo Effendi and Muhammad Ibnu Pamungkas*

The Prototype of Smart Power Meter at Home Based on Internet of Things

*Andi Irwansyah Satriananda, Lia Kamelia, Mufid Ridlo Effendi and Aep Kusnawan*

AB MIX Hydroponics Nutrient Solution Concentration Control Using Microcontroller Based On-Off Control Method

*Disan Alvin Nur Fadillah, Adam Faroqi, Lia Kamelia and Ah Fathonih*

Optimization of Ring Topology Routes Using Node and Link Swapping Methods

*Dimas Aji Pangestu, Radifa Akbar Abhesa, Zahrina Maryam and Nana Rachmana Syambas*

Travel Salesman Problem Algorithm Using Proposed Algorithm Development with K-Smallest City Method

*Arba Robbani, Ayu Latifah, Devani Claudia and Nana Rachmana Syambas*

Hardware Design and Implementation of Long Range Communication System for Rural Area  
*I Kompiang Gede Wirahita Putramas, Ian Joseph Matheus Edward, Iskandar, and Dharma Favitri Hariyanto*

Design and Implementation of Software and Web Dashboard on Long Range Communication Systems for Rural Area  
*Putu Priyana Pradipta, Ian Joseph Matheus Edward, Iskandar, and Dharma Favitri Hariyanto*

Prototype Smart Fish Farm in Fish Farming Koi  
*Pupug Ginanjar, Sarah Opipah, Dadan Rusmana, Muhlas Muhlas, Mufid Ridlo Effendi and Eki Ahmad Zaki Hamidi*

Design and Implementation of The Blind Navigation Aids Using Ultrasonic Sensor  
*Muhammad Yasir, Indri Nurfazri Lestari, Cucu Setiawan, Ulfiah Ulfiah, Mufid Ridlo Effendi and Eki Ahmad Zaki Hamidi*

Design and Implementation of Clothesline and Air Dryer Prototype Base on Internet of Things  
*Mohammad Haekal Gifari, Irfan Fahmi, Ajid Thohir, Abdullah Syafei, Mufid Ridlo Effendi and Eki Ahmad Zaki Hamidi*

Ball Identification and Localization using Regression for Wheeled Robot Soccer  
*Annisa Firasanti, Ahmad Fahrurrozi, Aeri Sujatmiko, Putra Wisnu Agung Sucipto and Eki Ahmad Zaki Hamidi*

ANN Design Model to Recognize the Direction of Multi-Robot AGV  
*Aan Eko Setiawan, Angga Rusdinar, Rina Mardiaty and Eki Ahmad Zaki Hamidi*

Fast Heuristic Algorithm Optimization for Travelling Salesman Problem  
*Ahmad Fauzi Iskandar, Ali Farhani Sani, Risyad Riyadi, Shafira Febriani and Nana Rachmana Syambas*

Power Allocation Effect on Capacity of Single Carrier Power Domain Non-Orthogonal Multiple Access (NOMA)  
*Muhammad Avi Majid Kaaffah and Iskandar*

Parameter-Based Clustering Algorithm for Clustered Wireless Sensor Networks with High Altitude Platforms as the Base Station  
*Dimas Aji Pangestu, Hendrawan and Iskandar*

Effect of Imperfect SIC in Non-Orthogonal Multiple Access (NOMA) over Rayleigh Fading Channels  
*Shita Herfiah and Iskandar*

Performance Comparison Between MIMO-NOMA 4x4 and MIMO-OMA 4x4  
*Ana Farahdiba and Iskandar*

Multiple User Fairness Power Allocation for Downlink NOMA  
*Novelita Rahayu and Iskandar Iskandar*

Arabic Part of Speech (POS) Tagging Analysis using HMM Trigram method on Al-Qur'an Ayah Sentences  
*Arief Fatchul, Izki Zakiyah Al-Hamro, Asep Solih Awalluddin and Muhammad Ibnu Pamungkas*

Arabic Part of Speech (POS) Tagging Analysis using Bee Colony Optimization (BCO) Algorithm on Quran Corpus

*Arief Fatchul, Fauziah Fauziah, Elis Ratnawulan and Dian Rahmat Gumelar*

Analysis Part of Speech Tagging Using Hidden Markov Model on Qur'an Data

*Arief Fatchul Huda, Muhammad Hafidz Naufal Hilal, Aep Saepuloh and Dedi Supriadi*

Development of Intelligent Telegram Chatbot using Natural Language Processing

*Asaad Balla Fadlelmula Babiker, Teddy Surya Gunawan, Nanang Ismail and Mufid Ridlo Effendi*

Analyzing Interference Problems in 5G Wireless Network

*Karisa Ardelia Hanifah, Faizal Husni, Ali Farhani Sani, Weryyan Shalannanda, Irma Zakia and Tutun Juhana*

Comparative Analysis of Cloud RAN Model on Low-latency Application for IoT-based 5G Networks

*Putu Priyana Pradipta, Gelar Pambudi Adhiluhung, Hafizh Mulya Harjono, Weryyan Shalannanda, Irma Zakia and Tutun Juhana*

Path Loss Estimation of 5G Millimeter Wave Propagation Channel–Literature Survey

*Muhammad Brata, Irma Zakia*



UINSGDNet  
Services

Wildan Budiawan Zulfikar <wildan.b@uinsgd.ac.id>

---

## ICWT 2021 submission 18

1 message

---

**ICWT 2021** <icwt2021@easychair.org>  
To: Wildan Budiawan Zulfikar <wildan.b@uinsgd.ac.id>

Thu, Apr 15, 2021 at 9:04 PM

Dear authors,

We received your submission to ICWT 2021 (The 7th International Conference on Wireless and Telematics 2021):

Authors : Ulfa Mardatillah, Wildan Budiawan Zulfikar, Aldy Rialdy Atmadja, Ichsan Taufik and Wisnu Uriawan  
Title : Citation Analysis on Scientific Articles Using Cosine Similarity  
Number : 18

The submission was uploaded by Wildan Budiawan Zulfikar <[wildan.budiawan.z@gmail.com](mailto:wildan.budiawan.z@gmail.com)>. You can access it via the ICWT 2021 EasyChair Web page

<https://easychair.org/conferences/?conf=icwt2021>

Thank you for submitting to ICWT 2021.

Best regards,  
EasyChair for ICWT 2021.



Wildan Budiawan Zulfikar &lt;wildan.b@uinsgd.ac.id&gt;

## ICWT 2021 notification for paper 18

2 messages

**ICWT 2021** <icwt2021@easychair.org>

Sat, Jul 31, 2021 at 8:15 PM

To: Wildan Budiawan Zulfikar <wildan.b@uinsgd.ac.id>

The 7th International Conference on Wireless and Telematics (ICWT) 2021 Editorial Committee has completed the reviewing process, and we are pleased to inform you that your manuscript,

Title : Citation Analysis on Scientific Articles Using Cosine Similarity  
Paper Number : 18

has been ACCEPTED for the VIRTUAL presentation in the 7th International Conference on Wireless and Telematics (ICWT) 2021 on 19-20 August 2021. We require the author(s) to revise the full paper according to the reviewers' comments (if any) and prepare 13-15 minutes presentation video of your paper (voiceover and talking head on presentation slides--the template can be accessed via <https://bit.ly/icwt2021pptx>).

The author shall submit the camera-ready full paper and the video presentation BEFORE Tuesday, 10 August 2021 23.59 GMT+7 by filling a form in <https://bit.ly/icwt2021>.

The full paper MUST STRICTLY pdf-express compatible (instructions in the postscript) and comply with the guidelines for Camera-Ready Submission at <https://icwt-seei.org/2021/submission-guidelines/>.

The proceedings of ICWT 2015, ICWT 2016, ICWT 2017, 2018, ICWT 2019, and ICWT2020 have been uploaded to the IEEE Xplore. For this year, accepted and presented papers will be included in the proceeding of ICWT 2021 and submitted to the IEEE Xplore. Failure in complying with the ICWT 2021 conference template may cause the full paper to be EXCLUDED for submission to IEEEExplore. For your consideration in presenting a paper at the 7th ICWT in Bandung (in person or virtually), at least one author of each accepted paper MUST register at the FULL registration fee, and the required FULL registration fee MUST be paid before the final revision deadline. Visit our website or contact us for additional registration information.

We apologize for the late notification.

Sincerely yours,  
ICWT 2021 Committee

P.S.:

Please use the IEEE pdf-express site to develop your compatible pdf file as follow:

- Add the ICWT2021 copyright notice on the bottom left part of the first page: 978-1-6654-4402-6/21/\$31.00 ©2021 IEEE
- Log in to the following IEEE pdf Express site: <https://ieee-pdf-express.org/>
- Select "Create account" link to create a new account.
- Enter "52862X" as the Conference ID of the ICWT2021, enter your email address, enter your password, and continue to enter information as prompted.
- Login to your account, select the "Create New Title", and upload your paper.
- After the process is completed, your IEEE Xplore compatible PDF will be available shortly. A copy of this PDF file will be sent to your email.

SUBMISSION: 18

TITLE: Citation Analysis on Scientific Articles Using Cosine Similarity

----- REVIEW 1 -----

SUBMISSION: 18

TITLE: Citation Analysis on Scientific Articles Using Cosine Similarity

AUTHORS: Ulfa Mardatillah, Wildan Budiawan Zulfikar, Aldy Rialdy Atmadja, Ichsan Taufik and Wisnu Uriawan

----- Overall evaluation -----

SCORE: 0 (borderline paper)

----- TEXT:

1. Explain, what is the new contribution proposed in this paper?
2. Is the cosine similarity method in determining the level of similarity your proposal, if not then you must include a reference to the method.
3. Briefly explain the basic principle of TF-IDF and cosine similarity in the Section before Section II (Methodology)
4. Provide an explanation of the validation of the level of success that you did in processing data.
5. Revise your manuscript from grammatical errors and typos and also table and figure position in the page should be improved.

----- REVIEW 2 -----

SUBMISSION: 18

TITLE: Citation Analysis on Scientific Articles Using Cosine Similarity

AUTHORS: Ufa Mardatillah, Wildan Budiawan Zulfikar, Aldy Rialdy Atmadja, Ichsan Taufik and Wisnu Uriawan

----- Overall evaluation -----

SCORE: 0 (borderline paper)

----- TEXT:

Similarity check result is over 20%. Please rewrite/paraphrase texts from references, make sure the revised similarity check result is 20% maximum. Proofread to make sure there is no typographical errors, grammatical errors, or untranslated words from foreign languages.

---

**Wildan Budiawan Zulfikar** <wildan.b@uinsgd.ac.id>  
To: Ulfamardatillah@gmail.com

Sun, Aug 1, 2021 at 9:28 AM

[Quoted text hidden]

--

Regards

**Wildan Budiawan Z**

Department of Informatics

UIN Sunan Gunung Djati Bandung



---

## Copyright Pending Notice for Article: Citation Analysis on Scientific Articles Using Cosine Similarity

1 message

---

**ECopyright@ieee.org** <ECopyright@ieee.org>  
To: wildan.b@uinsgd.ac.id

Mon, Nov 1, 2021 at 11:00 AM

Username and password to access the Electronic Publication Agreement for: Citation Analysis on Scientific Articles Using Cosine Similarity

Publication Title: 2021 7th International Conference on Wireless and Telematics (ICWT)  
Article Title: Citation Analysis on Scientific Articles Using Cosine Similarity  
Author E-mail: [wildan.b@uinsgd.ac.id](mailto:wildan.b@uinsgd.ac.id)  
eCF Paper Id: 341859

Dear Sir or Madam:

The above-titled article has been submitted to IEEE for publication in an upcoming journal or conference proceedings. The publication volunteer for that journal or conference has indicated to us that you are authorized to complete the article's publishing agreement with IEEE. Further, he/she has provided your name and e-mail address.

URL: <https://ecopyright.ieee.org/ECTT/login.do>  
Username: [wildan.b@uinsgd.ac.id](mailto:wildan.b@uinsgd.ac.id)  
Password: 1635739247692

Please be aware that the manuscript submission process will not be able to continue until this important step is completed. Finally, if you are not authorized to complete the publishing agreement, I hope you will forward the e-mail on to the person you believe has the necessary authority to sign. Thank you, in advance, for your consideration and effort.

IEEE IPR Office

PLEASE DO NOT RESPOND TO THIS EMAIL.

For technical assistance or to search our knowledge base, please visit our support site at :  
<http://ieee.custhelp.com/app/answers/list/p/82>

# Citation Analysis on Scientific Articles Using Cosine Similarity

*by N/a N/a*

---

**Submission date:** 19-Apr-2023 10:58AM (UTC+0700)

**Submission ID:** 2069019562

**File name:** Citation\_Analysis\_on\_04.pdf (186.7K)

**Word count:** 3354

**Character count:** 16102



# Citation Analysis on Scientific Articles Using Cosine Similarity

**2** Ulfa Mardatillah  
Department of Informatics  
UIN Sunan Gunung Djati  
Bandung, Indonesia  
ulfamardatillah@gmail.com

Wildan Budiawan Zulfikar  
Department of Informatics  
UIN Sunan Gunung Djati  
Bandung, Indonesia  
wildan.b@uinsgd.ac.id

Aldy Rialdy Atmadja  
Department of Informatics  
UIN Sunan Gunung Djati  
Bandung, Indonesia  
aldy@if.uinsgd.ac.id

Ichsan Taufik  
Department of Informatics  
UIN Sunan Gunung Djati  
Bandung, Indonesia  
ichsan@uinsgd.ac.id

Wisnu Uriawan  
Department Informatics and  
Mathematics  
INSA, Lyon, France  
wisnu\_u@uinsgd.ac.id

**Abstract**— Citations have a vital role in scientific articles. Every sentence that is listed must be obtained from the reference. There is a measure that is similarity. The problem is, the likeness of a sentence in a scientific article to the source it refers to can have different meanings. The high level of similarity can undoubtedly increase the similarity value if it is checked using the plagiarism (similarity checker) tool, where the higher the value is more similar to other articles on the contrary with low similarity remaining. However, a low level of similarity can also mean that a statement in a scientific paper has referred to a completely irrelevant source. This research aims to analyze the level of similarity between scientific articles and other articles or sources it refers to. Checking the input text starts with text preprocessing, weighting the TF-IDF and determining the similarity level using Cosine Similarity. Based on the test results, it is found that the results in a scientific article document have a high level of similarity with the article it refers to, and the other 30% have a low level of similarity, or in other words, there is no relationship.

**Keywords**—scientific articles, journals, TF-IDF, cosine, text mining, citation, incorrect citation

7

## I. INTRODUCTION

The development of information technology is increasing and has a positive impact. One of the positive impacts of information technology development is the ease of exchanging information. This convenience is often misused by someone or several people in completing work. The misused usually occurs in Scientific Documents. Scientific documents are documents of reasoning and research results using a variety of objects and methods. A good document has explicit references and sources as a reference, for example, a scientific paper or journal that has explicit references and citations and does not misuse the original citation of the document [1]–[4].

Scientific journal manuscripts are the author's scientific research manuscripts and consist of direct and indirect quotations. Direct quotations are quotations that are written exactly like the source, both in language and spelling [5], [6]. An indirect quotation is a quotation that is not the same as the original, and the citer only takes the main idea of the quoted source to be restated with the sentence compiled by the quote. Quotations are considered important as valid

references and are included in the writing of scientific papers [2], [3], [7], [8].

Mistakes that harm scientific journal scriptwriters are called plagiarism. Based on previous research that uses a cosine similarity method to check for plagiarism only, this research complements and refines it by weighing each word in scientific documents and testing the text-similarity value. So, documents that have a high similarity value are abused quotes and a high similarity value. The low similarity is a document that has been cited with little abuse or towards correctness.

The website-based bibliographic reference citation checking system for scientific journals will be used in checking scientific documents. Then, this study uses the TF-IDF weighting method, which functions to determine the weighting value of the validity of scientific journal manuscripts, and the Cosine Similarity Algorithm, which looks for word equations according to scientific journal manuscripts.

The basic principle of the method used in this research is Text Preprocessing using Case Folding, Tokenizing, Filtering, Stemming. Furthermore, the value of word weights in scientific documents, namely TF-IDF, is carried out after the stemming. Also, stopword removal stages are calculated on the value or weight of a word (term) in the document. After obtaining the word weight value, then testing the level of similarity using Cosine similarity, which compares the similarity between documents. In this case, what is being compared is a query with a training document. In calculating cosine similarity, the first thing to do is do a scalar multiplication between the query and the document, add it up, multiply the length of the document and the square of the query length, and then calculate the square root. Furthermore, the scalar multiplication result is divided by the multiplication result of the length of the document and the query. In related works, this method is widely used and shows a good performance in several cases [9]–[15]. Also, it possible to add any improvement to increase its accuracy and performance [16]–[19].

## II. METHODOLOGY

### A. Text Processing

Text mining or by other names such as intelligent text analysis, text data mining, or knowledge discovery in a text can simply be interpreted as the process of finding patterns that were not previously seen in certain text documents or sources.

The initial step of text mining is preprocessing. Some of them are case folding or changing text to lowercase and without punctuation or special characters, filtering, namely by removing conjunctions, conjunctions, tokenizing the separation of a document or sentence into each word or term, and Stemming or separating affixes in words into basic word [15], [20].

TABLE I. CASE FOLDING

Document	Result
D1	ada beberapa penelitian terdahulu yang memanfaatkan mikrokontroller arduino uno dan sensor seperti untuk pengontrolan suhu dan ruang
D2	perancangan dan implementasi pengontrol suhu ruangan berbasis mikrokontroller arduino uno dengan sensor suhu lm3s layak digunakan dan diaplikasikan sebagai sistem pengontrol suhu ruangan

The next step is tokenizing. This process converts sentences into words with spaces as the separating feature as described in Table 2.

TABLE II. TOKENIZING

Document 1 (D1)	Document 2 (D2)
ada	perancangan
beberapa	dan
penelitian	implementasi
terdahulu	pengontrol
yang	suhu
memanfaatkan	ruangan
mikrokontroller	berbasis
arduino	mikrokontroller
uno	arduino
dan	uno
sensor	dengan
seperti	sensor
untuk	suhu
pengontrolan	lm3s
suhu	layak
ruang	digunakan
	dan
	diaplikasikan
	sebagai
	sistem

Filtering is removing unnecessary words in the input text. So, the filtered text will be easily processed at a later stage. This study uses 758 words stored in an external dataset. After that is stemming. This step changes the filtered word into a root word and removes the affixes. The results of the filtering and stemming processes can be seen in table 3.

TABLE III. D1 AND D2 AFTER FILTERING AND STEMMING

Document 1 (D1)	Document 2 (D2)
teliti	rancang
manfaat	implementasi
mikrokontroller	control
arduino	suhu
uno	ruang
sensor	basis

control	mikrokontroller
Suhu	arduino
Ruang	uno
	sensor
	suhu
	lm3s
	guna
	Aplikasi
	Sistem

### B. TF-IDF and Cosine Similarity

TF is obtained from the number of words that appear in a text, so you can find out how many times the word appears. And after this TF calculation will be continued with the calculation of DF. Table 4 the process of calculating the weight value to determine the number of words that appear in the citation text and the source of the quote so that it can be calculated by calculating the cosine similarity [21], [22].

TABLE IV. TERM FREQUENCY

No	Word	Term Frequency (TF)	
		Document 1 (D1)	Document 2 (D2)
1	Teliti	1	0
2	Manfaat	1	0
3	mikrokontroller	1	1
4	Arduino	1	1
5	uno	1	1
6	sensor	1	1
7	control	1	2
8	suhu	1	3
9	ruangan	1	2
10	rancang	0	1
11	implementasi	0	1
12	basis	0	1
13	lm3s	0	1
14	guna	0	1
15	aplikasi	0	1
16	sistem	0	1

Then calculate the value of DF (document frequency), which is the number of words that appear in the two documents tested as described in table 5.

TABLE V. DOCUMENT FREQUENCY

No	Word	Document Frequency
1	teliti	1
2	manfaat	1
3	mikrokontroller	2
4	arduino	2
5	uno	2
6	sensor	2
7	control	3
8	suhu	4
9	ruangan	3
10	rancang	1
11	implementasi	1
12	basis	1
13	lm3s	1
14	guna	1
15	aplikasi	1
16	sistem	1

In table 6 are the results of the IDF (inverse document frequency) which is the value of the word analysis using (1) which is the value of the times with the TF (term frequency) in order to get the weight value.

$$1 + \log(N/df + 1) \quad (1)$$

TABLE VI. IDF RESULT

Word	DF	IDF 1	IDF 2
Teliti	1	1	0
Manfaat	1	1	0
Mikrokontroler	2	0.5	0.5
Arduino	2	0.5	0.5
Uno	2	0.5	0.5
Sensor	2	0.5	0.5
Control	3	0.5	0.823908741
Suhu	4	0.5	1.301
Ruangan	3	0.5	0.823908741
Rancang	1	0	1
Implementasi	1	0	1
Basis	1	0	1
lm3s	1	0	1
Guna	1	0	1
Aplikasi	1	0	1
Sistem	1	0	1

After calculating the IDF (inverse document frequency), the next step is to analyze the IDF, the TF value, and the IDF value that have been obtained from the above calculations, then the TFIDF analysis is a calculation to get the weighted value of the quoted text and the source of the citation tested. The following table 7 is a TFIDF analysis table.

TABLE VII. TF-IDF RESULT

No	Word	TF		IDF		TF IDF	
		D1	D2	D1	D2	D1	D2
1	Teliti	1	1	1	1	1	0
2	Manfaat	1	1	1	1	1	0
3	mikrokontroler	1	0.5	0.5	0.5	0.5	0.5
4	Arduino	1	0.5	0.5	0.5	0.5	0.5
5	Uno	1	0.5	0.5	0.5	0.5	0.5
6	Sensor	1	0.5	0.5	0.5	0.5	0.5
7	Control	1	0.5	0.5	0.5	0.5	0.82
8	Suhu	1	0.5	0.5	0.5	0.5	1.30
9	Ruangan	1	0.5	0.5	0.5	0.5	0.82
10	Rancang	0	0	0	0	0	1
11	Implementasi	0	0	0	0	0	1
12	Basis	0	0	0	0	0	1
13	lm3s	0	0	0	0	0	1
14	Guna	0	0	0	0	0	1
15	Aplikasi	0	0	0	0	0	1
16	Sistem	0	0	0	0	0	1

The cosine similarity calculation is described in table 8 where document 1 is the citation text and document 2 is the source text of the citation. The citation text is multiplied by the citation source text, the citation text and the citation source are squared so that the total is obtained to find the text-similarity value.

TABLE VIII. CALCULATION OF COSINE SIMILARITY

Word	TF IDF				
	D2	D2	D1.D2	D1^2	D2^2
Teliti	1	0	0	1	0
Manfaat	1	0	0	1	0
mikrokontroler	0.5	0.5	0.25	0.25	0.25
Arduino	0.5	0.5	0.25	0.25	0.25
Uno	0.5	0.5	0.25	0.25	0.25
Sensor	0.5	0.5	0.25	0.25	0.25
Control	0.5	0.82	0.41	0.25	0.67
Suhu	0.5	1.30	0.65	0.25	1.69
Ruang	0.5	0.82	0.41	0.25	0.67
Rancang	0	1	0	0	1
implementasi	0	1	0	0	1
Basis	0	1	0	0	1
lm3s	0	1	0	0	1
Guna	0	1	0	0	1
Aplikasi	0	1	0	0	1

Sistem	0	1	0	0	1
			2.47	3.75	11.05

$$\text{Crossproduct} = 6.13 \quad (2)$$

$$A.B = 2.4744, |A|, |B| = 6.13 \quad (3)$$

$$\text{Cos}(A,B) = 0.4036 \quad (4)$$

Based on the results of the manual algorithm calculation above (4), the similarity value in the text is 0.4036. In this study, the level of similarity will be divided into 3 categories, namely high, medium, and low. The low level is in the range of 0-0.3, the medium level is in the range of 0.3-0.6, and the high level is in the range of 0.6-1.0. The result of (4) is 0.4036 so that it is known that the text has a medium level.

### III. RESULT AND DISCUSSION

The system test scenario defines the amount of data used and what processes are carried out in the system to determine the error in a citation of this scientific journal. The process carried out in testing this system is data collection, text preprocessing, weighting, and text similarity Cosine Similarity. The data used in testing this system is an online journal available on the internet by looking for references to the bibliography listed in the document to be tested and inputting the citation text. 10 of the 50 test samples can be seen in the table below.

TABLE IX. TEST RESULT

No.	Citation	Result	Similarity Level
1	T001	0.43352308	Medium
2	T002	0.12427994	Low
3	T003	0.07988254	Low
4	T004	0.16726948	Low
5	T005	0.76533020	High
6	T006	0.08811075	Low
7	T007	0.02087507	Low
8	T008	0.05383300	Low
9	T009	0.16374133	Low
10	T010	0.24542823	Low

In the table above, there is a threshold, that is, if the similarity of the text is 0-0.3, it means that the similarity level is low. If the similarity of the text is 0.3-0.6 then the similarity is medium. If the similarity text results in a value of 0.6-1.0, the similarity level is high as described in the section on methodology.

Based on the 50 samples the data tested has a different similarity value for each text input. The thing that affects the similarity value is the number of word weights tested. In the first data sample, the number of text 1 has 9 words and the input text 2 has 15 words, with the number of word weights that are not too far in the range. In the second input text test, between text 1 and text 2, the number of words is not proportional, so it is dominated by text 2. And the similarity value is tested based on the weight of the words that appear in the two input texts. Documents with medium and high levels concerning the suitability of the quotation and the source of the citation are 70%, and the low level is 30%, unequal documents are documents that do not have a match between the quotation and the source of the citation. The



results of data validation checking scientific journal documents are that 70% of the citations in the article are true and 30% of the citations are incorrect.

#### IV. CONCLUSION

Checking scientific journal citation in testing sample data, it is found that the text is not similar. The citation text that does not have a similarity value, the text close to similar is the text that has a value close to similar, and similar is the text that has a similarity value. TF-IDF weighting in testing one citation text with one citation source text that comes from the bibliography, then the weighting greatly affects the value of the similarity of the text. The similarity of Cosine Similarity in testing one citation text with one citation source text that comes from the bibliography. Then the text-similarity value is obtained by the threshold value and knowing the text is included in several categories, namely not similar, close to similar, and similar with similar accuracy and Approaching similar, which is 70% and not 30% similar because of similarity and approaching similar. There is an agreement between the quote and its source and the word weight value that affects it.

#### REFERENCES

- [1] M. Golosovsky, *Citation Analysis and Dynamics of Citation Networks*. Cham: Springer International Publishing, 2019.
- [2] H. F. Moed, "Citation analysis in research evaluation," in *Proceedings of ISSI 2005: 10th International Conference of the International Society for Scientometrics and Informetrics*, 2005, vol. 2, pp. 437–441.
- [3] A. Shabani and R. Saadat, "Examining the citations received by DOAJ's journals from ISI Web of Science's articles (2003-2008)," in *Proceedings - 3rd International Conference on Information Sciences and Interaction Sciences, ICIS 2010*, 2010, pp. 109–113.
- [4] X. Su, A. Prasad, M. Y. Kan, and K. Sugiyama, "Neural multi-task learning for citation function and provenance," in *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries*, 2019, vol. 2019-June, pp. 394–395.
- [5] M. Roman, A. Shahid, S. Khan, A. Koubaa, and L. Yu, "Citation Intent Classification Using Word Embedding," *IEEE Access*, vol. 9, pp. 9982–9995, 2021.
- [6] Y. Yang and Y. Zheng, "The effect of open access journals on citation impact: A citation analysis of open access journals using Google Scholar," in *4th International Conference on Cooperation and Promotion of Information Resources in Science and Technology, COINFO 2009*, 2009, pp. 278–280.
- [7] H. H. Bi, J. Wang, and D. K. J. Lin, "Comprehensive citation index for research networks," *IEEE Trans. Knowl. Data Eng.*, vol. 23, no. 8, pp. 1274–1278, 2011.
- [8] P. A. H. Kazi and M. S. Patwardhan, "Context based citation summary of research articles: A step towards qualitative citation index," in *IEEE International Conference on Computer Communication and Control, IC4 2015*, 2016.
- [9] P. P. Gokul, B. K. Akhil, and K. K. M. Shiva, "Sentence similarity detection in Malayalam language using cosine similarity," in *RTEICT 2017 - 2nd IEEE International Conference on Recent Trends in Electronics, Information and Communication Technology, Proceedings*, 2017, vol. 2018-January, pp. 221–225.
- [10] W. Darnalaksana, C. Slamet, W. Budiawan Zulfikar, F. Fadillah, D. Sa'adillah Maylawati, and H. Ali, "Latent semantic analysis and cosine similarity for hadith search engine," vol. 18, no. 1, pp. 217–227, 2020.
- [11] S. Roopak and T. Thomas, "A novel phishing page detection mechanism using HTML source code comparison and cosine similarity," in *Proceedings - 2014 4th International Conference on Advances in Computing and Communications, ICACC 2014*, 2014, pp. 167–170.
- [12] S. Pattnaik and A. K. Nayak, "Summarization of odia text document using cosine similarity and clustering," in *Proceedings - 2019 International Conference on Applied Machine Learning, ICAML 2019*, 2019, pp. 143–146.
- [13] S. Zhu, J. Wu, and G. Xia, "TOP-K cosine similarity interesting pairs search," in *Proceedings - 2010 7th International Conference on Fuzzy Systems and Knowledge Discovery, FSKD 2010*, 2010, vol. 3, pp. 1479–1483.
- [14] D. Soyusyawaty and Y. Zakaria, "Book data content similarity detector with cosine similarity (case study on digilib.uad.ac.id)," in *Proceeding of 2018 12th International Conference on Telecommunication Systems, Services, and Applications, TSSA 2018*, 2018.
- [15] B. Li and L. Han, "Distance weighted cosine similarity measure for text classification," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2013, vol. 8206 LNCS, pp. 611–618.
- [16] K. C. Kirana, S. Wibawanto, N. Hidayah, and G. P. Cahyono, "The improved artificial neural network based on cosine similarity in facial emotion recognition," in *International Conference on Electrical Engineering, Computer Science and Informatics (EECSI)*, 2019, pp. 45–48.
- [17] P. Sitikhu, K. Pahi, P. Thapa, and S. Shakya, "A Comparison of Semantic Similarity Methods for Maximum Human Interpretability," in *International Conference on Artificial Intelligence for Transforming Business and Society, AITB 2019*, 2019.
- [18] C. Beşiktaş and H. A. Mantar, "Real-time traffic classification based on cosine similarity using sub-application vectors," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2012, vol. 7189 LNCS, pp. 89–92.
- [19] M. Kryszkiewicz, "Determining Cosine Similarity Neighborhoods by Means of the Euclidean Distance," *Intell. Syst. Ref. Libr.*, vol. 43, pp. 323–345, 2013.
- [20] M. Biba and M. Mane, "Sentiment analysis through machine learning: An experimental evaluation for Aslibanian," *Adv. Intell. Syst. Comput.*, vol. 235, pp. 195–203, 2014.
- [21] M. Irfan, Jumadi, W. B. Zulfikar, and Erik, "Implementation of Fuzzy C-Means algorithm and TF-IDF on English journal summary," in *2017 Second International Conference on Informatics and Computing (ICIC)*, 2017, pp. 1–5.
- [22] M. M. Saad, N. Jamil, and R. Hamzah, "Evaluation of Support Vector Machine and Decision Tree for Emotion Recognition of Malay Folklores," *Bull. Electr. Eng. Informatics*, vol. 7, no. 3, pp. 479–486, Sep. 2018.

# Citation Analysis on Scientific Articles Using Cosine Similarity

## ORIGINALITY REPORT

10%

SIMILARITY INDEX

8%

INTERNET SOURCES

7%

PUBLICATIONS

3%

STUDENT PAPERS

## PRIMARY SOURCES

1	<a href="http://ur.aeu.edu.my">ur.aeu.edu.my</a> Internet Source	2%
2	Submitted to School of Business and Management ITB Student Paper	1%
3	Agung Wahana, Wildan Budiawan Zulfikar, Wildan Najah Wildiansyah, Aldy Rialdy Atmadja, Diena Rauda Ramdania, Beki Subaeki. "A Deep Learning Approach To Analyze The Sentiment Of Online Game Users", 2021 7th International Conference on Wireless and Telematics (ICWT), 2021 Publication	1%
4	Submitted to Universitas Putera Batam Student Paper	1%
5	<a href="http://docplayer.net">docplayer.net</a> Internet Source	1%
6	<a href="http://thesai.org">thesai.org</a> Internet Source	1%
7	<a href="http://ijasce.org">ijasce.org</a>	

Internet Source

1 %

8

[join.if.uinsgd.ac.id](http://join.if.uinsgd.ac.id)

Internet Source

1 %

9

Submitted to Panjab University

Student Paper

1 %

10

Narongkorn Uthathip, Pornrapeepat Bhasaputra, Woraratana Pattaraprakorn. "Application of ANFIS Model for Thailand's Electric Vehicle Consumption", Computer Systems Science and Engineering, 2022

Publication

<1 %

11

"Table of Content", 2021 7th International Conference on Wireless and Telematics (ICWT), 2021

Publication

<1 %

12

[www.hindawi.com](http://www.hindawi.com)

Internet Source

<1 %

13

[www.slideshare.net](http://www.slideshare.net)

Internet Source

<1 %

Exclude quotes Off

Exclude matches Off

Exclude bibliography On