# WEIGHTED TREE SIMILARITY SEMANTIC SEARCH FOR E-COMMERCE CONTENT

[1]**I MADE SUKARSA,** [2] **I MADE BAYU DWI PUTRA** [3]**I GUSTI MADE ARYA SASMITA**

[123]Department of Information Technology, Udayana University, Indonesia

E-mail: [1]e_arsa@yahoo.com , [2] madebayudwiputra@gmail.com , [3] Aryasasmita@yahoo.com

**ABSTRACT**

This paper is to explain about e-commerce semantic search content. This Algorithm is using weight tree representation, TF-IDF calculation and cosine similarity. This method made the output of search result be more detail. Searched article is inputed to the system manually. Then, registered user or not, is able to do the article search on the system. Input from keywords inserted by the user processed by weighted tree calculation TF-IDF, and Cosine Similarity. The highest content weight will be displayed in on the highest list.

**Keywords**: *Semantic Search, Search Engine, Semantic Similarity ,TF-IDF, Weighted Tree Similarity*

## 1. INTRODUCTION

The development of technology today is growing very fast. It caused by the nature of human beings who require a device that able to simplify their lives. And with the technology development, people are getting easier to do their activities without doing it manually.

One of the developments in information technology is semantic search, many definitions that explain the meanings of semantic search, but in general semantic search is a science of representing and understanding of a word and search engine data. On semantic search there were many things that can be searched in the search process, such as content, images, and the results are varied one example of a link [1].

This research especially used for e-commerce content. The steps are preprocessing, Tf-Idf calculation and normalization with the cosine similarity. Preprocessing step consists of case folding, tokenizing, stop word, and stemming. After the basic word was formed, then the weighting step of the keyword will be done by checking the frequency of word appearance in the content. The results are sorted from the highest frequency of the appearance of word to the lowest one.

Semantic search can be applied to the e-commerce content search. Semantic technology is a smart application that gives simplicity to people to do the process of a wanted content searching.

## 2. THEORITICAL BACKGROUND

### 2.1 Semantic Search

Semantic search engine is a search technology that is being developed at this time were very influential and useful for the user to obtain a more detailed search results [2].In a semantic search many processes that must be passed, but the core of the search is the process of matching between data and keywords entered by the user [3]. Applications semantic search must have a user friendly display, where the latter process user searches by entering keywords and semantic search applications will eventually provide a search result page ranking list of recommended web [7]. In general a semantic search keeps the information content of a semantic web and the search is not only content that can be searched through the search process, in addition to the search results can be content in the form of images and the others [4].

### 2.2 E-Commerce

E-commerce have a wide definition, but in general e-commerce activity is sales, purchases, which uses the internet. E-commerce consists of some of the most common models include commercial sales made directly to consumers and sales to consumers among users [5].There are few examples of e-commerce models

a). Business-to-consumer

An often carried out trading, where the selling process made directly.

b). Business-to-business

A trading process with another company, for example by doing trading from the production company to the distribution company.

c) Business-to-government

A trading process between the private company with the government.

d) Consumer- to-government

Government's trading process to do the payment of many things, such as tax payment.

e) Consumer-to-consumer

A trading process where user is able to sell the product to another user.

### 2.3 Weighted Tree Similarity Algorithm

Weighted tree similarity algorithm is one of the unique algorithm because this algorithm consists of a tree thinking about the node labeled, labeled branch, and the branch weight [6]. At the algorithm level of similarity is between 0 and 1. Where the value is closer to 1, the value has the same level of similarity [1]. The steps that used in the calculation is to calculate the term or the occurrence of the word the most or the most common, having known a couple who have a branch it will proceed with the calculation of similarity with cosine similarity.

### 2.4 Related Research

The application that designed had some comparison and unique to general semantic search. The result of content is especially about electronic such as mobile phones, laptop, TV and the other. The method that used in search process is weighted tree similarity, where content uploaded manually and stored in a e-commerce database.

There are earlier research about semantic search. Sofi Silvia [11] has conducted similar research, in indonesian language entitled "Rancang Bangun Search Engine Tafsir AL-Quran Yang Mampu Memproses teks Bahasa Indonesia Menggunakan Metode Jaccard similarity" where the research aims to facilitated to find verse of Al-quran by measuring similarity of documents. In Sofi Silvia research was used Jaccard similarity method.

*Table 1. Comparative Study Of Semantic Search*

| Feature | Related Research | Our Research |
|---|---|---|
| Search column | ✓ | ✓ |
| Upload document | ✓ | ✓ |
| Grab document | - | ✓ |
| Preprocessing text | ✓ | ✓ |
| Weight calculation | ✓ | ✓ |
| List document/article history | - | ✓ |

Table 1 show the comparative between this research and related semantic search. Data of table showed that there were many different about both of research. This research had more feature then related research, there were: Grabbed document and list of document/article history feature.

### 3. SYSTEM OVERVIEW

This section explains about the overview or scheme of the system, where at every search engine consists of many features, but basically the main features of the system are search feature and search source selection feature. System schemes designed as good as possible so the system able to running well and the resulting output was as expected.
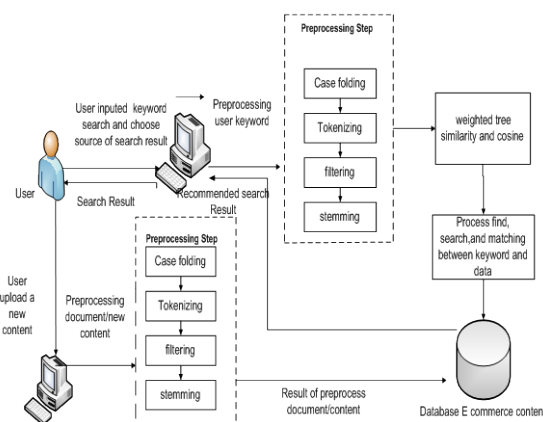


*Figure 1 System overview diagram*

Just like another system or application, there is general description of the system that is useful to give a description of the flow system to be made. These are explanations of system overview on figure 1 :

1). User was able to input keywords in the system, than that user can upload new content to system.

2). Search keywords or new content uploaded by the user was get to the preprocessing step and calculation results will stored in database.

3). Calculation step done by counted the number of TF, IDF, and cosine.

4).System will provide search results based on the highest weight and the system will recommend the search results to the user.

### 3.1 Semantic Search interface

Semantic search application must have a good user interface in order to help users to interact with the semantic search applications easier. The user interface is the most important part of the section, where the user inputs a query that will be processed and the results will be displayed in the semantic interface. The search results will be displayed at the interface application, in the list of web content rankings form, according to input keywords on the search column. At this semantic search application ,       the user interface consists           of:

1). Search column interface, where in this interface is the most important part of the semantic search applications. In this part, the user can input keywords that are searched in semantic search application.

2). Search source interface, where in this interface the user can choose the source of the search results and will be displayed in the search results on the semantic search applications.

3). Content upload interface, which in this interface users can upload content that will be stored in the e-commerce content database.

4). Article list interface, where in this article list interface we can see the list of uploaded articles by the user.

### 3.2 Example of manual calculation and tree scheme

This section is describes the tree scheme of e-commerce content and an example manual calculations between tree contend and user's input tree. A tree article has 3 parents: keyword, categories, and sources. Each parent has a weight, which if added together would be worth 1. The weight of each parent is obtained from the equation:

$$W\ i = 1/n \qquad (1)$$

**Description:**

$W_i$        : the weight of the i-th parent
n        : total existing parent

Each parent has a child called identifier. The identifier of the keyword parents are the keyword1 and 2. Identifier of the category parents are a phone and tablet, while the identifier for the source parents are google, yahoo, and bing. Each identifier also has a weight derived from the equation:

$$W\ ind\ =\ frek\ /\ n \qquad (2)$$

$W_{ind}$       : the i-th identifier weights
frek       : the number of appearance identifier
n        : total existing identifier

After finding the weight of the parent and the identifier, the next step can be done by knowing the tree input from the user.
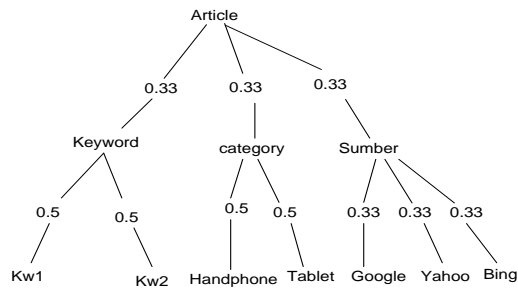


*Figure 2 Article tree scheme design*
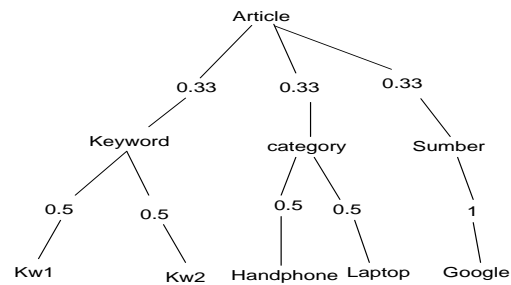


*Figure 3 User's input tree scheme design*

The calculation of the similarity can be done by comparing two trees, that has the same parent name. The calculation done by comparing each similar parent with its identifier. The calculation can be seen in the equation:

$$BK = \sum \Big( \big( W_{(i)} * W_{(j)} \big) * \big( \sum W\ ind_{(i)} * W\ ind_{(j)} \big) \Big) \qquad (3)$$

**Description:**

BK     : similarity weights
$W_{(i)}$     : the weight of the i-th primary parent
$W_{(j)}$     : the weight of the j-th input parent
$Wind_{(i)}$ : the weight of the i-th primary identifier
$Wind_{(j)}$ : the weight of the j-th input identifier

Based on the two trees, the weights of similarity between input user and article tree can be known on the following calculation:

$$[(0.33 * 0.33) * [(0.33 * 1) + (0.33 * 1)]]$$
$$= [ 0.1089 * [ 0.33 + 0.33]]$$
$$= [ 0.1089 * 0.66 ]$$
$$= 0.071874$$

And can be found the weight of similarities between input user and article at 0.071874

### 3.3 Weighting Combination With TF, IDF, And Cosine Similarity.

This section describes the semantic search's step in combination with cosine similarity, which the steps are, counting the number of frequency term, calculate the Inverse Document Frequency TF-IDF) and cosine similarity calculation to calculation of semantic search. Below is an example of the calculations where there are 5 documents and keyword test is "blackberry onyx tipe baru warna putih garansi tam".

Document 1 : Blackberry Gemini Yang Berwarna Merah.
Document 2 : Nokia Lumia Berfitur Kamera.
Document 3 : Blackberry Bold Yang Bergaransi TAM.
Document 4: Samsung Galaxy Yang Bertipe Layar Sentuh.
Document 5 : Nokia Yang Berwarna Putih.

From the used six documents, the document that is used as a test document or keyword is document 2. After the test documents that will be used is decided, the chosen document will go to the preprocessing step. In the process retrieval data, preprocessing step applied to the data in database[9].

a). C*ase folding* step

The step to change the uppercase to lowercase.

Document 1 :blackberry gemini yang berwarna merah.
Document 2 : nokia lumia berfitur kamera.
Document 3 : blackberry bold yang bergaransi tam.

Document 4 :samsung galaxy yang bertipe layar sentuh.
Document 5 :nokia yang berwarna putih.

b). S*topword* step

This step is to eliminate words that are not important, which does not affect the search results.

Document 1 :blackberry gemini berwarna merah.
Document 2 :nokia lumia berfitur kamera.
Document 3 :blackberry bold bergaransi tam.
Document 4 :samsung galaxy bertipe layar sentuh.
Document 5 :nokia berwarna putih.

c). *Stemming* step

This step is to eliminates prefix and suffix to acquire basic words

Document 1 :blackberry gemini warna merah.
Document 2 :nokia lumia fitur kamera.
Document 3 :blackberry bold garansi tam.
Document 4 :samsung galaxy tipe layar sentuh.
Document 5 :nokia warna putih.

d). Counting *term frequency* step

The step to count number of appearance of the word in a document or content on the database.

Document 1 :blackberry gemini warna merah.
            1      1      1      1
Document 2 :nokia lumia fitur kamera.
            1      1      1      1
Document 3 :blackberry bold garansi tam.
            1      1      1      1
Document 4 :samsung galaxy tipe layar sentuh.
            1      1      1      1      1
Document 5 :nokia warna putih.
            1      1      1

After passed the frequency step calculation (Tf) the next step is calculation of inverse document frequency (idf). The calculation can be seen in the equation:

$$\log(n/df) \quad (4)$$

The next step is weight document term (wdt) calculation. Weight document term is weight of the word in a document or content with the results obtained multiplication between term frequency (Tf ) with *inverse document frequency* ( idf). The calculation can be seen in the equation :

$$wdt = tf.idf \quad (5)$$

The emergence of word calculation processing in the TF-IDF calculation. TF,IDF calculation table can be seen in Table 12 page 9. After that it is processed again with cosine method. Cosine table can be seen in table 13.

The next step is to apply the cosine formula:

$$cosine\ Di = sum(kk\ dot\ Di)/([sqrt(kk2)] \quad (6)\\ * [sqrt(Di2)])$$

cosine D1 = 0.314 / (1.503* 1.135)
          = 0.314 / 1.705
          = 0.184

cosine D2 = 0 / ( 1.503 * 1.272 )
          = 0 / 1.911
          = 0

cosine D3 = 1.131 / (1.503 * 1.272)
          = 1.131 / 1.911
          = 0.591

cosine D4 = 0.487 / ( 1.503 * 1.560)
          = 0.487 / 2.344
          = 0.207

cosine D5 = 0.644 / ( 1.503 * 0.894)
          = 0.644 / 1.343
          = 0.478

From these results it can be seen that the example keyword test, has the highest similarity level with the document 3.

## 4. DATABASE DESIGN AND USER INTERFACE

### 4.1 Database Design

This part explains about database design and relation among the tables used in the semantic search for e-commerce content search. Database is a collection of interconnected data [8]. The database design using 10 tables and 3 tables of them are tables that used during the preprocessing step, include affixes table, basic word table, and stop word table.



*Figure 4 Database design*

Each semantic search engine had a database and table which used to saves data and used for the calculation and data processing in semantic search engine [10].

These are explanations of the function each table in database:

a). tb_affix

This table is used to save a list of affix data. A list of affix used in the stemming process where in the input keyword was founded prefix and suffix word and document content. System will remove prefix and suffix in keywords and document content.

*Table 2. Metadata Of Tb_Suffix*

| Name_field | Data_type | *Length* |
|------------|-----------|--------|
| No | Mediumint | 3 |
| name | Varchar | 10 |
| Id_type | Tinyint | 3 |

b). tb_katadasar

This table is a dictionary of basic words which used to save the Indonesian language basic words that used during a stemming process.

*Table 3. Metadata Of Tb_Katadasar*

| Name_field | Data_type | *Length* |
|------------|-----------|--------|
| Id_katadasar | Integer | 10 |
| katadasar | Varchar | 70 |
| katadasar_type | varchar | 25 |

c.) tb_stopword

This table is used to save the unimportant words during a filtering process, such as "-", "  ", "^", and the other.

*Table 4. Metadata Of Tb_Stopword*

| Name_field | Data_type | Length |
|---|---|---|
| id_stopword | Integer | 5 |
| stopword | Varchar | 25 |

d). sex

This table is used to save gender of the user data.

*Table 5. Metadata Of Sex*

| Name_field | Data_type | Length |
|---|---|---|
| no_sex | Tinyint | 1 |
| sex | varchar | 9 |

e). m_city

This table is used to save a user city living data.

*Table 6. Metadata Of M_City*

| Name_field | Data_type | Length |
|---|---|---|
| no_city | tinyint | 2 |
| City_name | varchar | 20 |

f). tb_articel_grab

This table is used to save article grabbing result and additional content.

*Table 7. Metadata Of Tb_Article_Grab*

| Name_field | Data_type | Length |
|---|---|---|
| Id_content | Bigint | 5 |
| Title | Varchar | 70 |
| content | Text | - |
| Content_date | Datetime | - |
| Link | Text | - |
| Id_category | Mediumint | 9 |
| No_user | Integer | 11 |
| Id_source | Tinyint | 1 |
| Id_brand | Integer | 11 |
| Preprocess_article | Text | - |

g). tb_category

This table is used to save data of content category.

*Table 8. Metadata Of Tb_Category*

| Name_field | Data_type | Length |
|---|---|---|
| id_category | Mediumint | 9 |
| Category_name | Varchar | 30 |

h). tb_user

This table is used to save data of user who wants to upload or add content on the application.

*Table 9. metadata of tb_user*

| Name_field | Data_type | Length |
|---|---|---|
| No_user | Integer | 11 |
| First_name | Varchar | 20 |
| Last_name | Varchar | 40 |
| Pswd | Varchar | 32 |
| Birthdate | Date | - |
| No_sex | Tinyint | 1 |
| Address | Varchar | 50 |
| Telp_no | Varchar | 13 |
| Email | Varchar | 30 |
| No_city | Tinyint | 2 |

i). tb_brand

This table is used to save data of e-commerce brands.

*Table 10. Metadata Of Tb_Brand*

| Name_field | Data_type | Length |
|---|---|---|
| id_brand | mediumint | 10 |
| brand_name | Varchar | 50 |

j). tb_source

This table is used to save data of source of grabbing from e-commerce content.

*Table 11. Metadata Of Tb_Source*

| Name_field | Data_type | Length |
|---|---|---|
| Id_source | Tinyint | 1 |
| Source | Varchar | 10 |
| Logo_source | Varchar | 10 |

### 4.2 User Interface design

These are user interface display on semantic search to find the e-commerce content, where on figure 5 is a display when login.

1) Main Page

This page is a page that contains the menu and the search results sought by user. There are also various links to go to another page.



*Figure 5  User Interface Sistem*

2) Article List

List that has a function to show the history or an article note that has been input by user



*Figure 6  Article List Page*

3) Search Result Page

On this form will displayed search results for entered keyword. In figure 7 entered keyword is blackberry onyx tipe baru warna putih garansi tam and the highest content weight will be displayed on the highest list.



*Figure 7 Search result page*

### 4. CONCLUSION

The conclusion that can be resumed from this research are:

1) Weighted Tree Similarity method can be applied on the search engines. With the application of the method, the output search result be more detail, because of the various calculation processes before finally display on search results.

2) search application with weighted tree similarity method has varying weights of search results with the highest weight is the nearly correct with the search results that the user wants.

3) this application use 900 content that stored in a database.

4) Weakness of this application today is only possible to do the search with indonesian language.
Based on this research, suggestions for further research are:

1) Preferably the applications developed to perform a search by a lot of languages, particularly English

2) This application is designed to search for an input article by the user into the database. For the next research development, the article should be directly connected with the entire existing web, without any input into the application.

### REFRENCES

[1] Riyanarto Sarno, Yeni Anistyasari, and Rahimi Fitri," Semantic Search Pencarian Berdasarkan konten ":.yogyakarta, ANDI, 2012

[2] Ritu Khatri, Kanwalvir Singh Dhindsa, and Vishal Khatri "Investigation and Analysis of New Approach of Intelligent Semantic Web Search Engines" *International Journal of Recent Technology and Engineering (IJRTE)* ISSN: 2277-3878, Volume-1, Issue-1, April 2012, pp. 83-85.

[3] G. Sudeepthi ,G. Anuradha, and  Prof. M. Surendra Prasad Babu "A Survey on  Semantic Web Search Engine" *IJCSI International Journal of Computer Science Issues*, Vol. 9, Issue 2, No 1, March 2012, pp.242

[4] G. Madhu, Dr. A. Govardhan, and Dr. T. V. Rajinikanth "Intelligent Semantic Web Search Engines: A Brief Survey" *International journal of Web & Semantic Technology (IJWesT)* Vol.2, No.1, January 2011, pp. 38-40

[5] Joko Sutrisno, "Strategi Pengembangan Teknologi E-commerce Dengan Metode Swot: Studi Kasus PT.Chingmix Berhan Sejahtera", *Jurnal TELEMATIKA MKOM*, Vol.3 No.2, September 2011, pp.45

[6] Riyanarto Sarno, Faisal Rahutomo, "Penerapan Algoritma Weighted Tree Similarity Untuk PencarianSemantik",Yogyakarta:ANDI, Indonesia,2012.

[7] Paras Nath Gupta,Pawan Singh, Pankaj P Singh ,Punit Kr Singh , Deepak Sinha," A Novel Architecture of Ontology based Semantic Search Engine" *International Journal of Science and Technology* Volume 1 No. 12, December, 2012,pp.651

[8] Kadir, Abdul. "Belajar Database Menggunakan MySQL.":Yogyakarta:CV.ANDI OFFSET.indonesia, 2008.

[9] A. Anil Kumar "Text Data Pre-processing and Dimensionality Reduction Techniques for Document Clustering" *International Journal of Engineering Research & Technology (IJERT)* "Dept of CSE Sri Sivani College of Engineering Srikakulam, India, Vol. 1 Issue 5, July – 2012,pp.2

[10] Salah Sleibi Al-Rawi, Ahmed Tariq Sadiq, Sumaya Abdulla Hamad, "Design and Evaluation of Semantic Guided Search Engine," *International Journal of Web Engineering "*,pp.17-20

[11] Sofi Silvia SP, " Rancang Bangun Search Engine Tafsir AL-Quran Yang Mampu Memproses Teks Bahasa Indonesia Menggunakan Metode Jaccard similarity", malang: UIN maulana malik ibrahim.

*Table 12.Tf-Idf Calculation*

| Token | TF | | | | | | Df | D/df | IDF | W | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | KK | D1 | D2 | D3 | D4 | D5 | | | | KK | D1 | D2 | D3 | D4 | D5 |
| blackberry | 1 | 1 | 0 | 1 | 0 | 0 | 2 | 2.5 | 0.397 | 0.397 | 0.397 | 0 | 0.397 | 0 | 0 |
| gemini | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 5 | 0.698 | 0 | 0.698 | 0 | 0 | 0 | 0 |
| warna | 1 | 1 | 0 | 0 | 0 | 1 | 2 | 2.5 | 0.397 | 0.397 | 0.397 | 0 | 0 | 0 | 0.397 |
| merah | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 5 | 0.698 | 0 | 0.698 | 0 | 0 | 0 | 0 |
| nokia | 0 | 0 | 1 | 0 | 0 | 1 | 2 | 2.5 | 0.397 | 0 | 0 | 0.397 | 0 | 0 | 0.397 |
| lumia | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 5 | 0.698 | 0 | 0 | 0.698 | 0 | 0 | 0 |
| fitur | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 5 | 0.698 | 0 | 0 | 0.698 | 0 | 0 | 0 |
| kamera | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 5 | 0.698 | 0 | 0 | 0.698 | 0 | 0 | 0 |
| bold | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 5 | 0.698 | 0 | 0 | 0 | 0.698 | 0 | 0 |
| garansi | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 5 | 0.698 | 0.698 | 0 | 0 | 0.698 | 0 | 0 |
| tam | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 5 | 0.698 | 0.698 | 0 | 0 | 0.698 | 0 | 0 |
| samsung | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 5 | 0.698 | 0 | 0 | 0 | 0 | 0.698 | 0 |
| galaxy | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 5 | 0.698 | 0 | 0 | 0 | 0 | 0.698 | 0 |
| tipe | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 5 | 0.698 | 0.698 | 0 | 0 | 0 | 0.698 | 0 |
| layar | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 5 | 0.698 | 0 | 0 | 0 | 0 | 0.698 | 0 |
| sentuh | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 5 | 0.698 | 0 | 0 | 0 | 0 | 0.698 | 0 |
| putih | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 5 | 0.698 | 0.698 | 0 | 0 | 0 | 0 | 0.698 |

*Table 13: Vector space model and Cosine Calculation*

| Token | W | | | | | | $W^2$ | | | | | | KK * D1 | KK * D2 | KK * D3 | KK * D4 | KK * D5 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | KK | D1 | D2 | D3 | D4 | D5 | $KK^2$ | $D1^2$ | $D2^2$ | $D3^2$ | $D4^2$ | $D5^2$ | | | | | |
| blackberry | 0.397 | 0.397 | 0 | 0.397 | 0 | 0 | 0.157 | 0.157 | 0 | 0.157 | 0 | 0 | 0.157 | 0 | 0.157 | 0 | 0 |
| gemini | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| warna | 0.397 | 0.397 | 0 | 0 | 0 | 0.397 | 0.157 | 0.157 | 0 | 0 | 0 | 0.157 | 0.157 | 0 | 0 | 0 | 0.157 |
| merah | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| nokia | 0 | 0 | 0.397 | 0 | 0 | 0.397 | 0 | 0 | 0.157 | 0 | 0 | 0.157 | 0 | 0 | 0 | 0 | 0 |
| lumia | 0 | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| fitur | 0 | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| kamera | 0 | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| bold | 0 | 0 | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| garansi | 0.698 | 0 | 0 | 0.698 | 0 | 0 | 0.487 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 |
| tam | 0.698 | 0 | 0 | 0.698 | 0 | 0 | 0.487 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 |
| samsung | 0 | 0 | 0 | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 |
| galaxy | 0 | 0 | 0 | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 |
| tipe | 0.698 | 0 | 0 | 0 | 0.698 | 0 | 0.487 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0.487 | 0 |
| layar | 0 | 0 | 0 | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 |
| sentuh | 0 | 0 | 0 | 0 | 0.698 | 0 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0 | 0 |
| putih | 0.698 | 0 | 0 | 0 | 0 | 0.698 | 0.487 | 0 | 0 | 0 | 0 | 0.487 | 0 | 0 | 0 | 0 | 0.487 |
| CALCULATION | | | | | | | Sum Kk² | Sum W² Di | | | | | Sum Kk Dot Di | | | | |
| | | | | | | | 2.262 | 1.288 | 1.618 | 1.618 | 2.435 | 0.801 | 0.314 | 0 | 1.131 | 0.487 | 0.644 |
| | | | | | | | Sqrt Sum kk² | Sqrt (Sum W² Di) | | | | | | | | | |
| | | | | | | | 1.503 | 1.135 | 1.272 | 1.272 | 1.560 | 0.894 | | | | | |