

ASSOCIATION RULE UNTUK MENGHASILKAN KATA POPULER DARI JUDUL BERITA ONLINE

ASSOCIATION RULES FOR GENERATING POPULAR WORDS FROM NEWS TITLES ON ONLINE NEWS SITES

Muhammad Ihsan Zul¹, Baitul Atiq², Muhammad Fauzan Delfani³

^{1,2,3}Politeknik Caltex Riau

ihsan@pcr.ac.id

ABSTRACT

Online news sites play a vital role in delivering information via digital media. There are numerous online news websites all over the world. This online news portal has generated a large volume of news. These news stories can be used to find significant and popular information. Data mining using frequent itemset mining and the apriori algorithm can provide hot news information. Using the apriori algorithm, the reader can determine the relationship of each word. This research uses frequent itemset mining to generate popular words from news headlines. The apriori algorithm is used to process popular word associations, which are then visualized with the JavaScript library. In this research, the minimal support value is more than 13. This value was decided based on the amount of time it took to process data. This study was effective in generating a method for producing popular terms and their interrelationships in the form of an appropriate visualization.

Keywords : popular words, frequent itemset mining, apriori algorithm, news portals

ABSTRAK

Situs berita online memegang peranan penting dalam penyebaran informasi melalui platform digital. Terdapat banyak situs berita online yang tersebar di dunia. Situs berita online ini telah menghasilkan berita dalam jumlah yang sangat besar. Berita-berita ini dapat dimanfaatkan untuk menggali informasi penting dan populer. Penerapan data mining melalui frequent itemset mining dan algoritma apriori dapat menghasilkan informasi berita yang sedang hangat setiap harinya. Penerapan algoritma apriori memungkinkan pembaca dapat mengetahui keterkaitan setiap kata. Penelitian ini menerapkan frequent itemset mining untuk menghasilkan kata-kata populer yang diambil dari judul berita. Keterkaitan kata populer diolah dengan menggunakan algoritma apriori yang kemudian divisualisasikan dengan menggunakan library javascript. Nilai minimum support yang digunakan dalam penelitian ini adalah besar dari 13. Nilai ini dipilih berdasarkan waktu komputasi yang digunakan saat pemrosesan data. Penelitian ini berhasil mengembangkan suatu pendekatan dalam menghasilkan kata-kata populer dan keterkaitan antar kata tersebut dalam bentuk visualisasi yang menarik.

Kata Kunci: kata-kata populer, frequent itemset mining, algoritma apriori, portal berita

PENDAHULUAN

Media elektronik merupakan salah satu produk hasil perkembangan teknologi informasi. Media Elektronik berperan penting dalam penyampaian informasi kepada masyarakat. Kemudahan akses informasi melalui perkembangan teknologi telah mengubah model bisnis informasi yang dilakukan oleh perusahaan-perusahaan pemilik media. Kemudahan tersebut dimanfaatkan dengan mengembangkan situs berita online. Saat ini telah banyak perusahaan teknologi informasi yang menjalankan bisnisnya dalam pengembangan aplikasi berbasis

web. Sehingga bukan tidak mungkin bagi perusahaan yang tidak memiliki sumber daya dalam mengembangkan website dapat menggunakan jasa perusahaan-perusahaan tersebut. Bagi perusahaan media besar tentu saja mereka memiliki divisi khusus yang menangani bagian tersebut.

Keadaan ini mengakibatkan berkembangnya situs-situs berita online di Indonesia bahkan dunia. Tentu saja tidak semua situs berita online yang dapat dipercaya. Akan tetapi di Indonesia setidaknya terdapat beberapa perusahaan media online yang diakui kredibilitasnya. Setidaknya dapat dijadikan rujukan bagi

masyarakat. Perkembangan media tersebut mengakibatkan terjadinya banjir informasi. Keadaan ini mengakibatkan besarnya jumlah berita harian yang dihasilkan. Jika dikumpulkan dalam waktu satu tahun, media tersebut telah memproduksi data dalam jumlah yang sangat besar. Hal ini merupakan peluang yang dapat dimanfaatkan untuk menggali informasi dari berita-berita tersebut. Pendekatan Data Mining merupakan salah satu cara yang tepat untuk menggali informasi yang dihasilkan oleh situs berita online.

Setiap berita memiliki judul yang dapat digunakan sebagai sumber data. Judul tersebut dapat dimanfaatkan untuk menggali informasi penting yang dihasilkan oleh situs berita nasional. Beberapa hal yang dapat digali adalah: (1) kata-kata terbanyak, (2) keterkaitan antar kata yang digunakan sebagai judul berita, dan; (3) kategori berita. Kata-kata populer dapat dimanfaatkan untuk menggali kata-kata yang menjadi tren berita pada waktu tertentu. Jika kata-kata tersebut berhasil digali maka akan didapatkan keterkaitan antar kata berdasarkan durasi waktu tertentu.

Salah satu pendekatan yang dapat digunakan adalah Frequent Itemset Mining (FIM) dan *association rule* (Han, Pei, & Kamber, 2012). FIM merupakan salah satu teknik yang digunakan untuk menemukan item yang paling sering muncul dalam kelompok data nominal (Borglet, 2012). Teknik ini digunakan untuk dapat digunakan untuk menemukan kata-kata populer yang dipublikasikan oleh situs berita online di Indonesia. Pendekatan *association rule* merupakan salah satu teknik yang digunakan untuk menggali keterkaitan antar item yang terdapat di dalam data set (Han, Pei, & Kamber, 2012). Salah satu algoritma yang termasuk ke dalam *association rule* adalah algoritma apriori.

Berdasarkan kajian tersebut, penelitian ini mengajukan teknik

penggalan kata-kata populer dan keterkaitan antar kata-kata tersebut dengan menggunakan FIM dan algoritma apriori. Data yang digunakan adalah data judul berita yang diambil dari delapan situs berita online yang kredibel di Indonesia. Hasil penelitian ini divisualisasikan dan dapat digunakan untuk menganalisis berita berdasarkan kata-kata yang paling banyak muncul dan keterkaitannya dengan kata lain.

Penelitian yang berhubungan dengan judul berita sudah banyak dilakukan. Penelitian-penelitian tersebut memiliki fokus yang beragam diantaranya pengkategorian judul berita, pengklusteran judul berita dan penelitian lainnya.

Penelitian (Okumura dan Miura, 2016) mengembangkan model yang digunakan untuk memproses judul berita berdasarkan isi berita. Model ini dikembangkan dengan menggunakan Latent Semantics Analysis. Judul berita tidak sepenuhnya diproses secara otomatis, tetapi semiotomatis dengan melibatkan manusia dalam penentuannya. Penelitian (Thu dan Pa, 2020) hampir sama dengan penelitian sebelumnya. Penelitian ini menggunakan pendekatan sequence-to-sequence one-hot encoding dan RNN dalam mendapatkan judul berita. Bahasa yang digunakan dalam penelitian ini adalah bahasa Myanmar.

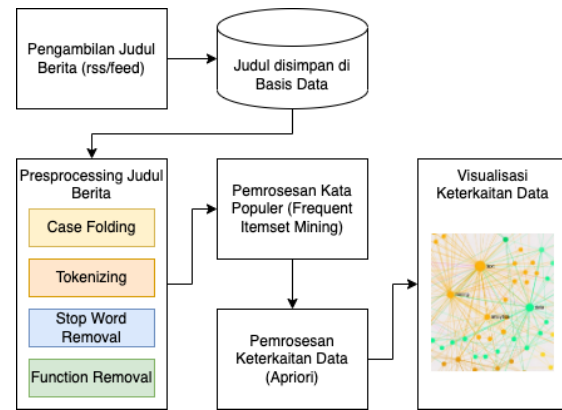
Penelitian (Abyaad dkk., 2020) menggunakan Deep Learning untuk mengkategorikan judul berita. Pendekatan yang digunakan adalah Long-Short Term Memory (LSTM). Pendekatan ini diklaim menghasilkan kategori berita dengan akurasi yang tinggi. Penelitian (Azam dkk., 2018) menggunakan ekstraksi ciri dari teks yang akan diklasifikasikan. Penelitian ini menggunakan algoritma k-nearest-neighbour dalam mengklasifikasikan teks berdasarkan hasil ekstraksi ciri. Penelitian ini juga menggunakan algoritma naïve bayes untuk mengklasifikasikan teks. Penelitian ini menggunakan Rapid Miner dalam penerapannya. Penelitian yang

dilakukan (Fuad dan Syamsuardi, 2019) menggunakan pendekatan natural language processing dalam menentukan trending topik dari text stream di dunia maya. Penelitian ini menggunakan berita sebagai sumber data. Pendekatan clustering digunakan untuk menentukan text yang menjadi trending topic.

Penelitian yang diulas pada paragraph sebelumnya lebih banyak membahas tentang penentuan judul berita berdasarkan konten berita, penentuan kategori berita dan pengklasifikasian berita dengan menggunakan algoritma k-nn dan naïve bayes. Penelitian ini berbeda dengan penelitian yang dipaparkan. Penelitian ini menggunakan judul berita sebagai sumber data yang kemudian judul berita tersebut diolah untuk menghasilkan kata-kata populer yang diambil dari beragam situs berita online di Indonesia. Penelitian ini menggunakan metode *association rule* dengan memanfaatkan frequent itemset mining dan algoritma apriori

METODE

Penelitian ini dilakukan melalui beberapa tahap yang meliputi tahapan-tahapan pengolahan kata. Tahap-tahap tersebut antara lain preprocessing meliputi case folding, tokenization, stop words removal, pemrosesan FIM dan penerapatan algoritma apriori. Hasil pengolahan data ini adalah kata-kata populer yang memiliki keterkaitan dengan kata lainnya. Keterkaitan data ditampilkan dalam bentuk visualisasi menggunakan javascript. Rincian tahapan penelitian dapat dilihat pada Gambar 1.



Gambar 1. Tahapan Penelitian

1. Pengumpulan Data

Data judul berita dikumpulkan dengan merancang aplikasi yang mengambil data dari situs berita *online* melalui rss/feed. RSS merupakan singkatan dari *Really Simple Syndication* yang digunakan untuk memberikan akses kepada pembaca dan pengembang aplikasi pihak ketiga untuk mengambil informasi terkait halaman web yang menyediakan layanan tersebut (Song, 2018). Fasilitas rss/feed ini memberikan akses bagi sistem yang dikembangkan untuk mengambil data judul berita dari media massa tersebut berdasarkan waktu yang ditentukan. Sistem ini bekerja dengan mengambil judul berita dan kemudian menyimpan judul tersebut ke dalam basis data yang disiapkan. Untuk menghindari duplikasi data, dilakukan pengecekan judul dari setiap situs berita dengan membandingkan judul tersebut dengan judul-judul yang telah tersimpan di dalam basis data. Daftar situs berita *online* yang digunakan dalam penelitian ini dapat dilihat pada Tabel 1.

Tabel 1. Daftar Situs Berita Online dan alamat rss/feed

No	Nama Situs Berita	Feed/RSS URL
1	Detik.com	http://rss.detik.com/index.php/detikcom
2	Antara Nasional	http://www.antara.co.id/rss/news.xml
3	Okezone	http://sindikasi.okezone.com/index.php/news/RSS2.0
4	Republika	http://www.republika.co.id/rss

No	Nama Situs Berita	Feed/RSS URL
5	Liputan6.com	http://feed.liputan6.com/rss
6	JPNN	https://www.jpnn.com/index.php?mib=rss
7	Suara.com	http://www.suara.com/rss
8	Merdeka	https://www.merdeka.com/feed/

Semua data judul yang dikumpulkan disimpan di dalam basis data. data yang disimpan tersebut antara lain nama situs berita (nama situs berita), tanggal dan jam penerbitan dan judul berita.

2. Preprocessing Judul Berita

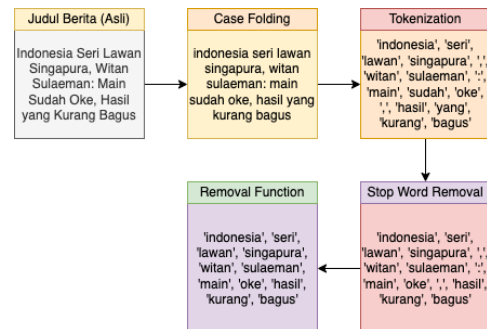
Keterkaitan antar kata yang didapatkan dari judul berita dilakukan dengan melakukan Preprocessing data judul berita. Preprocessing yang diterapkan di dalam penelitian ini adalah (1) *case folding*, (2) *tokenization*, (3) *stop word removal*, dan (4) *removal function*.

Case folding merupakan suatu cara yang dilakukan untuk mengubah semua data kalimat menjadi huruf kecil (Zul, Yulia, & Nurmalasari, 2018) (Chistol & Danubianu, 2021). Setiap judul berita yang disimpan diproses menjadi huruf kecil. Menurut (Chistol & Danubianu, 2021) *tokenization* merupakan suatu tahap yang dilakukan untuk memecah kalimat menjadi kata-kata pembentuk kalimat. Kata tersebut merupakan bentuk asli dari semua kata yang ditulis di dalam judul berita. Setelah kalimat diubah menjadi kata per kata, maka dilakukan proses penyaringan kata (*stop word removal*). Tujuan penyaringan kata adalah untuk membuang semua kata yang tidak diperlukan. Daftar kata yang akan dibuang dalam penelitian ini adalah semua kata sambung (konjungsi) Bahasa Indonesia (Sukarto, 2017).

Tahap berikutnya adalah *stop word removal*. Tahap ini perlu dilakukan untuk efisiensi pemrosesan pada tahap berikutnya. Stop word removal digunakan untuk membuang kata-kata yang termasuk dalam daftar kata yang tidak penting dalam penelitian ini. Tahapan ini akan membuat efisiensi pemrosesan ke tahap berikutnya.

Daftar kata yang akan dibuang adalah kata-kata yang tidak mempengaruhi makna dari penelitian (Ladani & Desai, 2020).

Setelah tahap stop word removal dilakukan, dilanjutkan dengan proses filter tanda baca atau dikenal dengan nama *removal function*. Sehingga kata-kata yang dihasilkan adalah kata terpilih tanpa tanda baca. Ilustrasi pemrosesan data dapat dilihat pada Gambar 2.



Gambar 2. Preprocessing Judul Berita

3. Frequent Itemset Mining (FIM)

Pemrosesan frekuensi kemunculan kata bagian paling penting di dalam penelitian yang berhubungan dengan text processing (Nasreen dkk., 2014). *Frequent itemset* mining digunakan untuk menghitung jumlah kemunculan kata yang ditulis oleh media *online* setiap harinya. Setiap kata hasil preprocessing akan dihitung jumlah kemunculannya (Shah & Patel, 2014). Persamaan yang digunakan untuk menghasilkan kata-kata populer berdasarkan frekuensi kemunculannya adalah sebagai berikut.

$$cover(P) = \{i | T_i \in D \wedge P \subseteq T_i\}$$

Frekuensi kemunculan absolut dilakukan dengan persamaan berikut.

$$freq(P) = |cover(P)|$$

dimana i adalah items (kata), P adalah itemset, T adalah data dan D merupakan dataset.

4. Algoritma Apriori

Algoritma apriori merupakan salah satu algoritma yang digolongkan ke dalam metode *association rule*. Algoritma apriori memiliki peran yang besar dalam menghasilkan data yang frequent dalam penggunaannya. Pada umumnya Apriori algorithm digunakan untuk keperluan *Market Basket Analysis* (MBA) (Ünvan, 2021). Tujuannya adalah untuk menganalisis data transaksi yang dilakukan di sebuah toko/retail. Akan tetapi algoritma apriori juga dapat dimanfaatkan untuk text mining. Salah satunya adalah untuk mengolah keterkaitan antar top words yang dihasilkan dari judul berita media *online*. Tahapan Algoritma apriori dapat dilihat pada *pseudocode* berikut (Han, Pei, & Kamber, 2012).

C_k : Candidate itemset of size k
 L_k : frequent itemset of size k
 $L_1 = \{\text{frequent items}\};$
for ($k = 1; L_k \neq \emptyset; k++$) do begin
 C_{k+1} = candidates generated from L_k ;
for each transaction t in database do
increment the count of all candidates in C_{k+1}
that are contained in t
 L_{k+1} = candidates in C_{k+1} with min_support
end return $\cup_k L_k$;

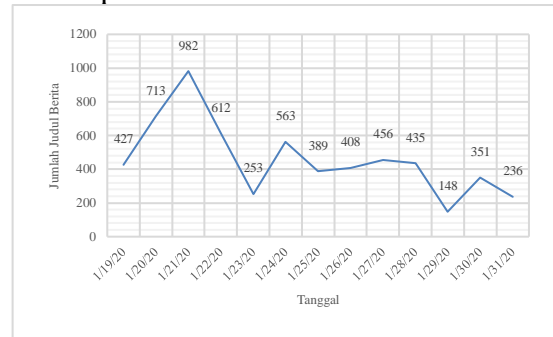
HASIL DAN PEMBAHASAN

1. Data Judul Berita

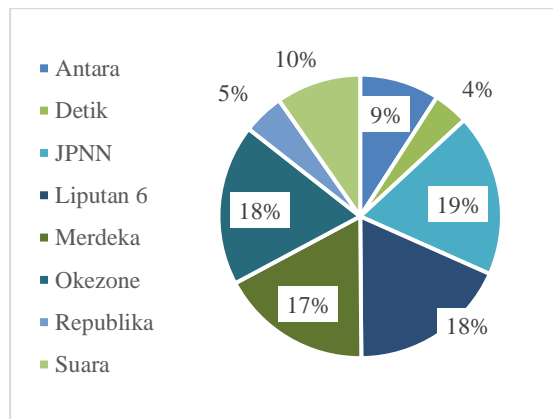
Penelitian ini berhasil mengumpulkan sebanyak 5973 data judul berita yang diambil dalam selama 13 hari antara tanggal 19 s.d. 31 Januari 2020. Pengumpulan data ini dilakukan dengan mengembangkan sistem pengumpul data yang secara otomatis mengambil data dari rss berita yang didaftarkan. Semua data tersimpan di dalam basis data yang dipersiapkan. Jumlah data yang dikumpulkan menurut hari dapat dilihat pada Gambar 3.

Selama pengumpulan data terdapat beberapa kendala yang dihadapi, diantaranya RSS sumber yang tidak bisa diakses, galat/error dari sisi sistem

pengumpul data dan kendala jaringan internet yang tidak stabil. Akan tetapi data tersebut tetap dapat digunakan untuk diolah menggunakan FIM. Jika dilihat dari sumber berita, JPNN, Liputan 6, Merdeka dan Okezone merupakan media massa online yang memproduksi di atas 800 data. Sebaran persentase jumlah data yang dikumpulkan berdasarkan situs berita dapat dilihat pada Gambar 4.



Gambar 3. Produksi Berita Setiap Hari



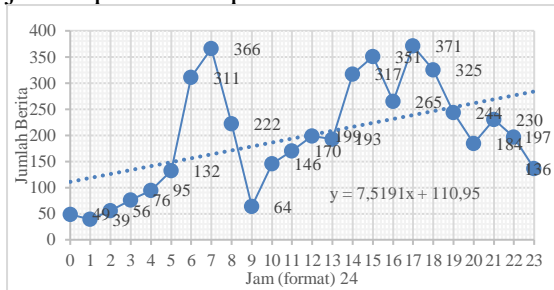
Gambar 4. Persentasi Produksi Berita Setiap Situs Online

2. Pemrosesan dengan FIM

Semua data yang tersimpan di dalam basis data diproses untuk menghasilkan kata-kata yang memiliki frekuensi tinggi berdasarkan nilai minimum support yang ditentukan pada saat pemrosesan melalui algoritma FIM. Semakin kecil nilai minimum support, semakin banyak data yang dihasilkan, dan semakin tinggi waktu komputasi yang dibutuhkan untuk memproses kata.

Sehingga pemrosesan dilakukan berdasarkan hari yang dipilih. Pemilihan

hari ini juga sangat ditentukan oleh jam pemrosesan data. Hal ini berbanding lurus dengan jumlah data yang diproduksi oleh setiap media massa elektronik setiap harinya. Semakin malam, maka akan semakin banyak data yang akan diproses dan akan semakin lama waktu yang diperlukan sistem untuk memproses data. Jumlah data yang dihasilkan berdasarkan jam dapat dilihat pada Gambar 5.



Gambar 5. Produksi Berita Berdasarkan Jam

Jika diperhatikan, produksi data berita tertinggi terjadi pada pagi hari antara pukul 06:00 s.d 08:00, dan sore hingga malam hari antara jam 14:00 s.d 19:00. Sehingga pengujian waktu pemrosesan data ini dilakukan setelah jam 19:00. Karena puncak produksi data pada pagi dan malam hari dapat diakomodir dengan kondisi tersebut. Hasil pengujian dapat dilihat pada Tabel 2.

Tabel 2. Hasil Pengujian Waktu Pemrosesan FIM

No	Nilai Minimum Support	Waktu Pemrosesan FIM (detik)
1	3	5,6
2	5	4,9
3	7	4,75
4	9	4,1
5	11	3,35
6	13	2
7	15	2,1
8	17	1,5
9	19	1,2
10	21	0,86

Berdasarkan data yang disajikan pada tabel xx, terlihat bahwa pemrosesan data dapat dikatakan baik jika nilai *minimum support* yang digunakan besar dari 13. Hasil pemroses FIM adalah kata yang memiliki nilai di atas *minimum support* yang ditentukan. Kata-kata tersebut akan tersaji

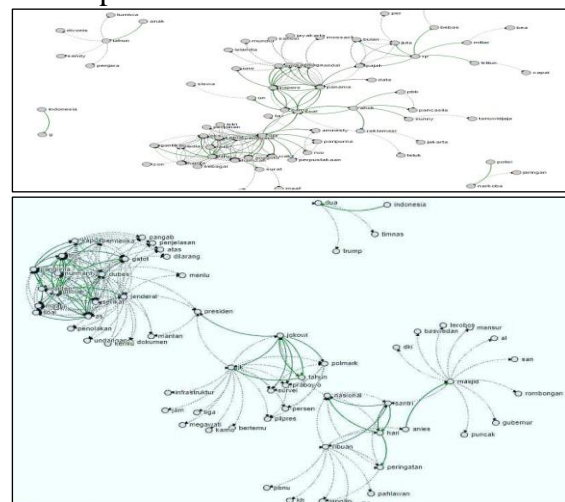
dalam bentuk *top words* seperti pada Tabel 3.

Tabel 3. Potongan kata populer hasil FIM

No	Item 1	Item 2	Item 3	Item 4	Item 5
1	video	update	corona	4	juni
2	update	corona	covid-19	4	juni
3	video	7	penumpa ng	air	asal
4	wni	china	wabah	virus	corona
5	jokowi	ingin	omnibus	law	ketua
6	terus	bertambah	korban	tewas	akibat
7	korban	meninggal	dunia	akibat	virus
8	wabah	corona	bagi	pemer ntah	buat
9	anies	banjir	underpass	kemay oran	pemprov
10	anies	pemprov	dki	bantu	banjir
11	korban	meninggal	china	obat	pasien
12	china	rumah	pasien	virus	corona
13	cegah	virus	corona	anak	anies
14	wabah	virus	korona	tni	siap
15	corona	4	juni	pasien	sembuh
16	update	corona	4	juni	jumlah
...
xxx	update	corona	4	juni	pasien

3. Visualisasi Data Algoritma Apriori

Kata-kata populer yang telah diproses dengan menggunakan FIM kemudian divisualisasikan berdasarkan keterkaitan antar kata dengan menggunakan algoritma apriori. Visualisasi diproses dengan menggunakan library javascript. Hasil visualisasi keterkaitan antar kata dapat dilihat pada Gambar 6.



Gambar 6. Keterkaitan data setelah penerapan algoritma apriori

Berdasarkan Gambar 6, dapat dilihat bahwa beberapa kata menjadi pusat keterkaitan dengan banyak kata lainnya. Selain itu juga terdapat kata yang menjadi penghubung antar kelompok kata. Algoritma ini berhasil memproses keterkaitan antar kata-kata populer yang dihasilkan FIM. Library javascript berperan

penting dalam menampilkan visualisasi kata-kata populer

SIMPULAN

Penelitian telah berhasil mengumpulkan data berita dengan menggunakan rss/feed dari situs berita online yang ditargetkan. Terdapat 5973 judul berita yang dikumpulkan dalam waktu 13 hari. Data tersebut kemudian diolah dengan menggunakan FIM untuk mendapatkan kata-kata populer berdasarkan pemberitaan. Kata-kata populer didapatkan dengan menggunakan nilai minimum support yang besar dari 13. Hal ini dilakukan agar pemrosesan data dapat dilakukan dengan cepat.

Metode *association rule* melalui algoritma apriori berhasil diterapkan untuk menghasilkan keterkaitan antar kata yang diolah dari hasil FIM. Library javascript dimplementasikan untuk menghasilkan visualisasi keterkaitan antar data. Melalui visualisasi ini ditampilkan kata-kata populer yang menjadi pusat dan penghubung antar kelompok kata. Sehingga melalui visualisasi ini, dapat ditarik informasi secara mudah mengenai kata-kata yang sedang tren yang dibahas situs berita online. Penelitian ini dapat dilanjutkan dengan mengidentifikasi konteks berita berdasarkan keterkaitan kata dan mengklasifikasikan judul berita berdasarkan kategori yang diinginkan.

DAFTAR PUSTAKA

Abyaad, R., Kabir, M. R., & Hasan, S. (2020). A Novel Approach to Categorize News Articles From Headlines and Short Text. *2020 IEEE Region 10 Symposium (TENSYMP)*. Dhaka Bangladesh.

Azam, M., Ahmed, T., Sabah, F., & Hussain, M. I. (2018). Feature Extraction based Text Classification using K-Nearest Neighbour Algorithm. *IJCSNS International*

Journal of Computer Science and Network Security, 18(12), 95-101.

- Borglet, C. (2012). Frequent Itemset Mining. *Wires Data Mining and Knowledge Discovery*, 2(6), 437 - 456.
- Chistol, M., & Danubianu, M. (2021). Survey of Text Mining Research Methods and Their Innovative Applicability. *Journal of Danubian Studies and Research*, 11(1), 225-233.
- Fuad, R. G., & Syamsuardi. (2019). *Penentuan Trending Topic Berdasarkan Text Stream Menggunakan Suffix Tree Clustering*. Pelambang: Universitas Sriwijaya.
- Han, J., Kamber, M., & Pei, J. (2012). *Data Mining Concept and Techniques 3rd Edition*. Waltham USA: Morgan Kaufmann Publishers.
- Han, J., Pei, J., & Kamber, M. (2012). *Data Mining: Concepts and Techniques*. Waltham, USA: Morgan Kaufmann.
- Ladani, D. J., & Desai, N. P. (2020). Stopword Identification and Removal Techniques on TC and IR applications: A Survey. *6th International Conference on Advanced Computing and Communication Systems (ICACCS)*. Coimbatore, India.
- Nasreen, S., Azamb, M. A., Shehzada, K., Naeem, U., & Ghazanfara, M. A. (2014). Frequent Pattern Mining Algorithms for Finding Associated Frequent Patterns for Data Streams: A Survey. *The 5th International Conference on Emerging Ubiquitous Systems and Pervasive Networks (EUSPN-2014)* (pp. 109-116). Elsevier B.V.
- Okumura, N., & Miura, T. (2016). Generating Headline Candidates for News Articles . *IEEE 17th International Conference on Information Reuse and Integration*. Las Vegas USA.

- Shah, A., & Patel, P. A. (2014). A Collaborative Approach of Frequent Item Set Mining: A Survey. *International Journal of Computer Applications* , 107(8), 34-36.
- Song, J.-W. (2018). Design and Implementation of Real-Time News App using RSS of the Internet Newspaper. *Journal of Digital Contents Society*, 19(4), 631-637.
- Sukarto, K. A. (2017). Konjungsi Bahasa Indonesia: Suatu Tinjauan. *Pujangga: Jurnal Bahasa dan Sastra*, 3(1), 98-112.
- Thu, Y., & Pa, W. P. (2020). Myanmar News Headline Generation with Sequence-to-sequence Model. *The 23th Conference of The Oriental (COCOSDA)*. Myanmar: -.
- Ünvan, Y. A. (2021). Market Basket Analysis with Association rules. *Communications in Statistics - Theory and Methods* , 50(7), 1615-1628.
- Zul, M. I., Yulia, F., & Nurmalasari, D. (2018). Social Media Sentiment Analysis using K-Means and Naïve Bayes Algorithm. *2nd International Conference on Electrical Engineering and Informatics (ICon EEI)*. Batam, Indonesia.