



Analisis Topic Modelling Persepsi Pengguna Internet Menggunakan Metode Latent Dirichlet Allocation

Angga Reni Dwi Astuti¹, Nuri Cahyono²

anggareni@students.amikom.ac.id, nuricahyono@amikom.ac.id

Universitas Amikom Yogyakarta

Informasi Artikel

Diterima : 4 Feb 2023

Direview : 11 Feb 2023

Disetujui : 26 Feb 2023

Kata Kunci

Text mining, Latent Dirichlet Allocation, Topic Modelling, python

Abstrak

Kemajuan teknologi tidak diragukan lagi memiliki dampak besar pada media informasi. Salah satu dampak kemajuan teknologi adalah keberadaan media berita sebagai sumber informasi publik. Ada juga informasi daerah, baik dalam maupun luar negeri, dan tentunya ada berbagai pembahasan. Data berita dari portal berita online dapat dijadikan sebagai sumber informasi sekaligus sumber penelitian dan analisis. Tentu saja, portal berita mencakup semua jenis berita tentang berbagai topik tertentu. Mengidentifikasi topik yang sering dibahas di portal berita pasti akan memakan banyak waktu. Oleh karena itu, penelitian ini berfokus pada penerapan sistem pemodelan topik untuk mengimplementasikan sistem keputusan topik berita dengan menggunakan metode *Latent Dirichlet Allocation* (LDA). Penelitian ini sukses menerapkan metode *Latent Dirichlet Allocation* (LDA) dalam menentukan topik berita, yang mana terdapat tiga topik kategori yang sering dibahas pada portal berita online detik.com. Topik 1 berisi kejadian bencana alam, topik 2 berisi tokoh dan permasalahan politik, topik 3 berisi berita seputar piala dunia.

Keywords

Text mining, Latent Dirichlet Allocation (LDA), Topic Modelling, python.

Abstrak

Technological advances have undoubtedly had a major impact on information media. One impact of technological progress is the existence of news media as a source of public information. There is also regional information, both domestic and foreign, and of course there are various discussions. News data from online news portals can be used as a source of information as a source of research and analysis. Of course, news portals cover all types of news on various topics. Identifying frequently discussed topics on news portal will definitely take a lot of time. Therefore, this research focuses on applying a topic modelling system to implement a news topic decision system using the Latent Dirichlet Allocation (LDA) method. This research successfully applies the Latent Dirichlet Allocation (LDA) method in determining news topics, of which there are three topic categories that are often discussed on the online news portal detik.com. topic 1 contains natural disaster event, topik 2 contains political figures and issues, topik 3 contains news about the world cup.

A. Pendahuluan

Saat ini adalah era dari kemajuan suatu teknologi, hadirnya teknologi tentunya membawa setiap kehidupan manusia bersinggungan langsung dengannya[1]. sehingga kita dapat melihat begitu banyak kemajuan diberbagai bidang di dalam dunia teknologi. Teknologi yang sangat berkaitan dengan kehidupan sehari – hari kita adalah teknologi informasi[2]. Kemajuan teknologi sendiri tentunya banyak perkembangan di beragam bidang, salah satunya aspek teknologi informasi. Peran dari teknologi informasi yaitu dapat membantu ataupun memudahkan pekerjaan agar berjalan lebih efektif dan efisien[1].

Kemajuan teknologi sendiri tentunya mendukung banyaknya jumlah informasi yang dapat kita terima, hal ini didukung dengan adanya kemajuan teknologi yang pesat sehingga semakin banyak dan semakin deras informasi yang dapat kita terima dan peroleh, dan tentunya akan sangat memudahkan kita untuk dapat mendapatkan dan bertukar suatu informasi didukung dengan teknologi penyedia informasi yang ada. Salah satunya yang akrab kita sebut dengan media sosial, yang mana ikut turut serta membantu kemajuan dari derasnya penyebaran suatu informasi saat ini. Media sosial adalah sebuah aplikasi berbasis internet yang diciptakan untuk memungkinkan penggunaannya untuk melakukan pertukaran konten berupa teks, gambar, video dan lain-lain [4]. Media sosial juga dapat dikatakan sebuah medium di internet yang memungkinkan penggunaannya untuk mengekspresikan diri dan melakukan interaksi, bekerjasama, dan berkomunikasi secara lebih luas, secara daring di era digital [5].

Dengan didukung kemajuan dari teknologi yang semakin berkembang, mengakibatkan banyaknya media sosial yang dapat kita gunakan, yang mana dari semua jenis media sosial tersebut memiliki kelebihan atau keunggulannya masing-masing yang ditawarkan bagi penggunaannya. Ada berbagai macam media sosial yang dapat dengan mudah kita jumpai dan gunakan diantaranya, ada instagram, facebook, line, twitter dan lain – lain.

Twitter menjadi salah satu sosial media yang cukup banyak digunakan dari berbagai jenis media sosial. Twitter termasuk jenis media sosial microblogging yang dapat memberikan penggunaannya untuk aktivitas menulis serta mempublikasi aktivitas serta opini mereka [6]. Twitter juga merupakan media sosial yang cukup banyak digunakan di Indonesia maupun negara lain. Sekitar 73% hingga 87% membaca riview secara online di Twitter [7]. Twitter sendiri tentunya memiliki keunggulan atau kelebihan tersendiri yang ditawarkan kepada penggunaannya, sehingga menarik banyak orang untuk terus menggunakan Twitter. Dengan twitter orang – orang dapat dengan lebih mudah dan leluasa untuk memberikan serta membagikan kehidupan mereka kepada publik, berupa gambar maupun tulisan, yang mana mereka dapat dengan bebas dan leluasa untuk mengekspresikan diri mereka sendiri di twitter.

Kemajuan teknologi sendiri juga ikut mempengaruhi dari cara penyajian berita saat ini. Berita adalah suatu peristiwa yang baru terjadi, dimana di dalam berita terbagi menjadi judul, teras, dan tubuh berita. Berita dibuat guna untuk memberikan informasi mengenai suatu topik permasalahan yang dibahas secara lengkap guna untuk informasi dan pengetahuan [8].

Pada era sebelumnya berita hanya dapat kita nikmati dalam bentuk offline, yaitu dalam bentuk media cetak, lisan, radio, siaran televisi. Namun pada saat ini dengan

didukung dengan kemajuan teknologi kita dapat merasakan dan menikmati berita secara online kapan dan dimana saja, yaitu dengan menggunakan internet seperti portal berita, media sosial dan lain – lain.

Portal berita adalah media yang menawarkan informasi berupa berita online. Keuntungannya ialah informasi yang lebih cepat daripada media lama, seperti koran dan majalah contohnya. Kita dapat mendapatkan berbagai informasi dari suatu portal berita, baik itu berita lokal, daerah, dalam negeri, maupun luar negeri sekalipun dan juga tentunya akan lebih banyak dan beragam jenis dari berita yang menjadi ulasan didalamnya. Hal ini dapat dilakukan dengan melihat postingan gambar ataupun tweet yang dilakukan salah satu media penyedia berita online di salah satu media sosial.

Dari banyaknya kategori dan jenis berita yang dibahas dalam suatu portal berita online, tentunya hal itu memerlukan waktu dalam melakukan identifikasi untuk menentukan topik sebenarnya yang sering kali dibahas dalam suatu akun berita online. Karena itu diharapkan penelitian ini bisa menjadi solusi untuk masalah penentuan topik artikel berita oleh situs berita online di media sosial twitter yaitu akun detik.com.

Penelitian ini menggunakan topik modelling. Topik modelling adalah pembentukan pola tertentu yang berasal dari sekumpulan teks, misal tweet atau dokumen lain. Karena itu kita bisa tahu gambaran umum isi dari topik utama. Ada berbagai macam metode yang dapat digunakan dalam pemodelan topik, salah satunya yaitu *Latent Dirichlet Allocation* (LDA).

Saat ini penelitian ini juga menggunakan metode *Latent Dirichlet Allocation* (LDA) yang merupakan bagian dari metode text mining, dimana ditemukan pola tertentu pada dokumen yang menghasilkan beberapa jenis topik tertentu. Metode *Latent Dirichlet Allocation* (LDA) lebih unggul dibandingkan metode pemodelan topik lainnya dan dapat digunakan untuk mengidentifikasi topik dalam jurnal, memeringkat dan mengklasifikasikan.

B. Metode Penelitian

Tahap ini menjelaskan secara rinci tentang tahapan penelitian, mulai tahap rancangan, pola penelitian, alat yang digunakan, teknik pengumpulan data, analisis, sistem dan hal lain yang terkait dengan strategi pemecahan masalah penelitian.

1. Identifikasi masalah

Pada tahap ini dilakukan guna untuk mengetahui permasalahan yang ada, yaitu dimana dengan derasnya informasi yang diterima,seringkali terjadi kesalah pahaman ataupun kurang tercernanya suatu informasi dari media sosial dengan baik dan karena begitu banyak topik yang dibahas membuat kita sering kali tidak mengetahui berita apa yang sedang dibahas. Sehingga dapat diketahui kebutuhan dan memberikan solusi terhadap masalah yang ada.

2. Mempelajari literatur

Tahap ini dilakukan dengan mempelajari literatur dan ilmu yang berkaitan langsung pada penelitian, misal mempelajari melalui internet, jurnal, buku, dan juga media lainnya yang menjadi sumber pembelajaran. Hal ini dilakukan guna untuk memperoleh informasi dan pengetahuan sebanyak mungkin, sehingga membantu peneliti untuk lebih menguasai dan juga memudahkan peneliti untuk menjalankan penelitian.

3. Pengumpulan data

Pengambilan data ini dilakukan dengan menggunakan library twint, menggunakan git bash, yang mengambil data dari twitter khususnya akun detik.com selama bulan januari – november 2022.

Twint merupakan alat yang dapat digunakan untuk mendapatkan data dari twitter. Twint sendiri dikembangkan dengan menggunakan bahasa pemrograman python dan dapat diinstal menggunakan pip ataupun conda. Bahasa python berisi sejumlah paket yang dapat berguna dan juga cocok untuk analisis [9].

Library twint digunakan karena dapat mengambil data tweet pada twitter tanpa harus memerlukan Twitter API [10].

4. Penentuan algoritma

Algoritma yang digunakan adalah topik modelling menggunakan metode *Latent Dirichlet Allocation* (LDA), yang memberikan hasil maksimal sehingga lebih mudah untuk mencapai hasil yang diinginkan. Topik modelling atau pemodelan topik adalah metode clustering yang termasuk dalam unsupervised learning, artinya tidak melebeli data. Pemodelan topik tersirat dalam pengelompokan, dimana setiap objek pada level tertentu dapat dimiliki oleh lebih dari satu kluster[11]. Pemodelan topik adalah metode analisis topik yang bertujuan untuk mengekstrak topik utama dari sebuah tweet, sehingga kita mendapatkan gambaran tentang topik utama [3].

5. Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) adalah teknik text mining untuk menemukan pola tertentu pada dokumen dengan membuat beberapa jenis topik yang berbeda, sehingga dokumen tidak dikelompokkan secara spesifik berdasarkan topik tertentu[13]. Kinerja *Latent Dirichlet Allocation* (LDA) lebih baik dibandingkan dengan metode pemodelan topik yang lainnya dan dapat diimplementasikan untuk identifikasi topik majalah, klasifikasi, dan clustering [14].

6. Text mining

Text mining adalah istilah untuk melambungkan data dalam bentuk teks dengan tujuan untuk menemukan kata yang dapat mewakili konten dari data sehingga kita dapat mengetahui hubungan antara data - data tersebut. Text mining juga dapat dikatakan sebagai proses mengolah koleksi data teks dari waktu ke waktu yang bertujuan untuk menemukan informasi yang bermanfaat dari sumber data dan mengetahui hubungan antar data [15].

6. Pengkodean

Tahapan ini adalah tahapan memasukan code program dengan bahasa pemrograman python dengan menggunakan google collab.

Tahap ini dilakukan agar dapat menjalankan fungsi dan juga mendapatkan hasil dari penelitian.

C. Hasil dan Pembahasan

Tahap ini memberikan penjelasan tentang menganalisis temuan dari pekerjaan untuk memberikan jawaban atau solusi atas masalah penelitian ini.

a. Pengambilan data

proses pengambilan data ini dilakukan dengan cara memanfaatkan library twint, dengan memanfaatkan git bash, yang mana akan mengambil data berupa tweet dari salah satu akun portal berita online yaitu akun detik.com.

Hasil dari pengambilan data ini disimpan dalam bentuk data csv, yaitu data tersebut berisi tweet yang dilakukan oleh akun detik dari bulan januari hingga bulan november 2022, yang mana data tersebut berisi kurang lebih 5017 data tweet yang dilakukan oleh detik.com selama periode tersebut, dan diberikan nama detikcom.csv pada data tersebut.

Gambar 1. File detikcom.csv

b. Processing data

Langkah selanjutnya adalah preprocessing data. Langkah ini dilakukan untuk menghapus beberapa simbol dalam proses analisis topik. Langkah ini penting saat membuat teks tidak berstruktur dan menyimpan kata kunci yang dapat membantu dalam menyajikan topik. Berikut adalah beberapa langkah preprocessing teks yang digunakan dalam penelitian ini :

1) Stop word removal

Tahapan ini dilakukan penghapusan kepada kata yang tidak mengandung atau memiliki makna arti lebih.

2) Tokenizing

Tokenizing merupakan tahapan untuk menghapus kata yang berlebihan serta memenggal kata yang menyusun suatu kalimat.

3) Case folding

Tahap ini berguna untuk mengubah seluruh character pada tweet menjadi huruf kecil.

4) Removal Punctuation

Menghilangkan karakter, url, sebutan, tagar, nomor, spasi ganda, dan tanda baca yang tidak diperlukan dalam proses penguraian.

5) Stemming

Yaitu untuk menghapus sufiks di dalam suatu kata.

6) Sastrawi

Proses pengurangan kata sehingga sesuai dengan standar kamus bahasa indonesia, dengan memanfaatkan library python.

Berikut adalah hasil dari proses processing data yang dilakukan.

E	F	G	H
Keywords	Text		
via, korban, gemj	[hotel, 'bidakara', 'jenazah', 'ferry', 'mursyid', 'temu', 'mobil', 'parkir', 'area', 'drop', 'off', 'birawa', 'assambly', 'hall', 'hotel', 'bidakara]		
via, korban, gemj	[video, 'viral', 'tonton', 'kali', 'wanita', 'hasil', 'kulit', 'sampo', 'anti', 'ketombe', 'aman', 'nggak]		
jokowi, presiden,	[rachel, 'vennya', 'keringat', 'olahraga', 'hatihati]		
dunia, piala, grup	[gagal, 'penalti', 'lionel', 'messi', 'argentina', 'vs', 'polandia', 'debat', 'absah', 'fifa', 'dukung', 'penuh', 'putus', 'wasit', 'danny', 'makiele]		
dunia, piala, grup	[teliti, 'perancis', 'jerman', 'rusia', 'hidup', 'jenis', 'virus', 'es', 'tanah', 'siberia', 'permafrost]		
via, korban, gemj	['korban', 'tinggal', 'dunia', 'akibat', 'gempa', 'cianjur', 'orang', 'data', 'korban', 'korban', 'bencana', 'rawat', 'rumah', 'sakit', 'via', 'jabar]		
dunia, piala, grup	[jerman, 'bahan', 'olokolok', 'piala', 'dunia', 'situasi', 'beda', 'alami', 'mesut', 'oezil', 'kes', 'sayang', 'helat', 'akbar]		
dunia, piala, grup	[bmw, 'group', 'resmi', 'luncur', 'bmw', 'seri', 'sedan', 'premium', 'alami', 'ubah', 'eksterior', 'kendara]		
dunia, piala, grup	[kisah', 'lengkap', 'kisah', 'cinta', 'pasang', 'beda', 'negara', 'sorot', 'warganet', 'bule', 'amrik', 'libruan', 'jepara', 'kepincut', 'wanita', 'semaran]		
jokowi, presiden,	[heboh', 'warganet', 'media', 'sosial', 'wanita', 'batal', 'nikah', 'calon', 'suami', 'meni', 'wanita', 'calon', 'suami', 'kaget]		
via, korban, gemj	['kabaops', 'polres', 'jakpus', 'akbp', 'saufi', 'salamun', 'luka', 'duga', 'akibat', 'lempar', 'batu', 'aman', 'demo', 'mahasiswa', 'papua', 'aliansi',]		
via, korban, gemj	[nik', 'hadir', 'saksi', 'hitc', 'jakarta', 'desember', 'tiket', 'batas]		
jokowi, presiden,	[insentif', 'nonfiskal', 'perintah', 'salah', 'mobil', 'listrik', 'bebas', 'ganjil', 'genap', 'bijak', 'nilai', 'diskriminatif]		

Gambar 2. Hasil Preprocessing

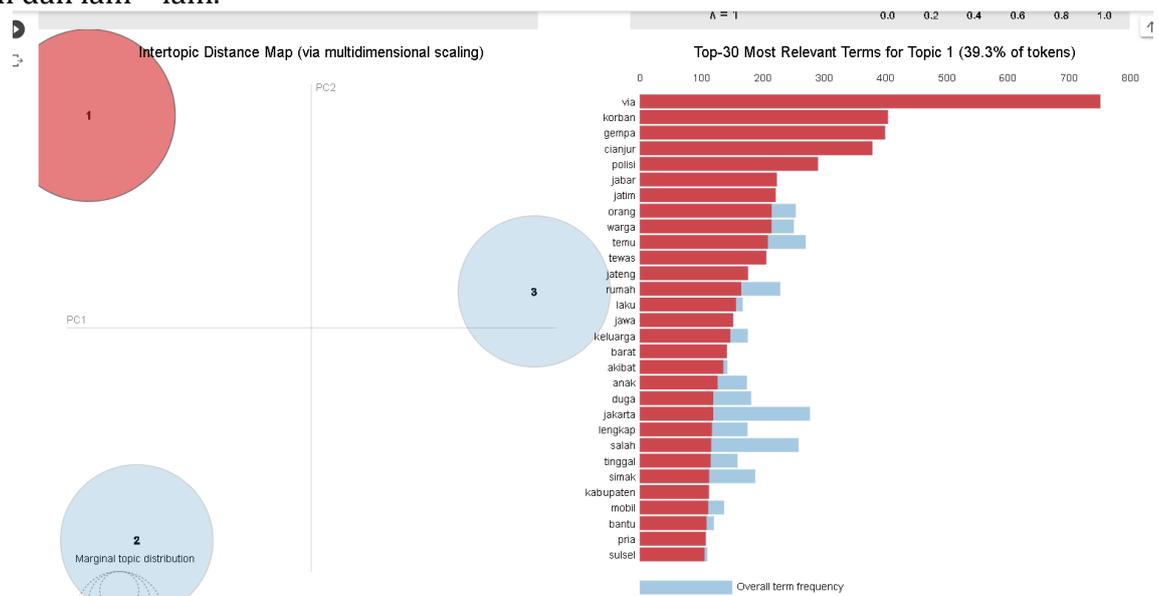
c. Hasil pengolahan

Pengolahan data pada penelitian ini dilakukan dengan preprocessing menggunakan metode *Latent Dirichlet Allocation* (LDA). Tujuan dari metode ini adalah untuk menentukan kata topik yang mewakili setiap topik dan diimplementasikan menggunakan library python. Hasil pengolahan data berupa diagram yang menggambarkan isi berita yang mencakup 3 topik yang sering menjadi pembahasan di akun media sosial twitter detik.com. hasil pemodelan topik dapat divisualisasikan menggunakan library pydavis dan hasil kajian ini dapat dilihat pada gambar dibawah ini.

1) Topik 1

Ini akan dibagi menjadi 3 topik yang dimana dalam 1 topik utama terdapat 30 pembahasan dalam satu topik.

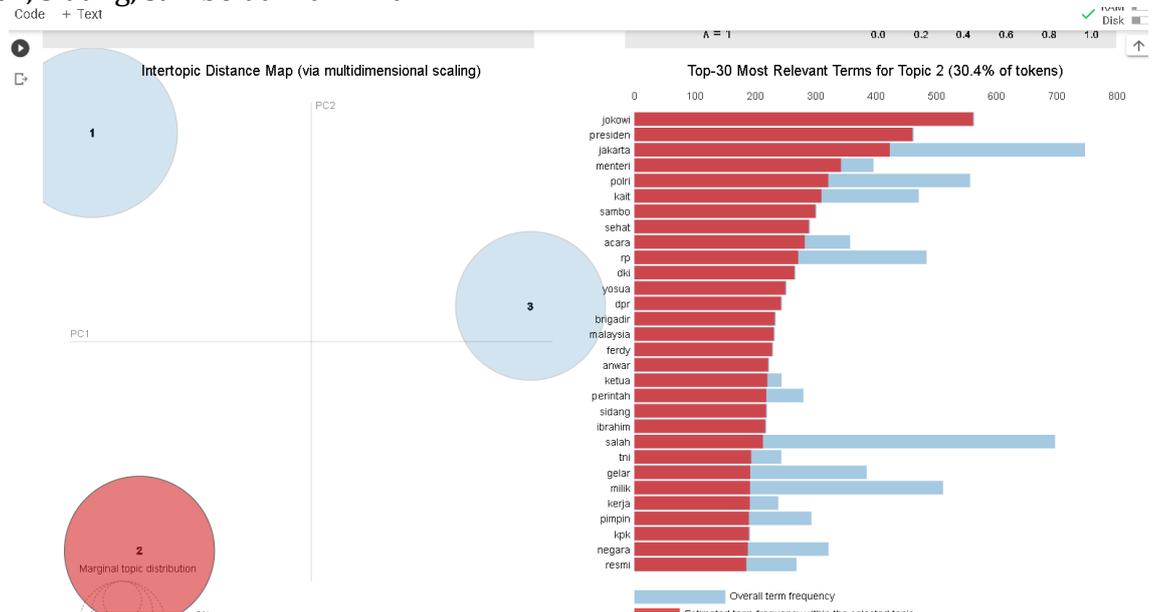
Topik pertama lebih dominan membahas mengenai bencana alam yang terjadi di indonesia, seperti gempa yang terjadi di cianjur, yang membahas gempa, lokasi, korban dan lain – lain.



Gambar 3. Diagram Dominan Topik 1

2) Topik 2

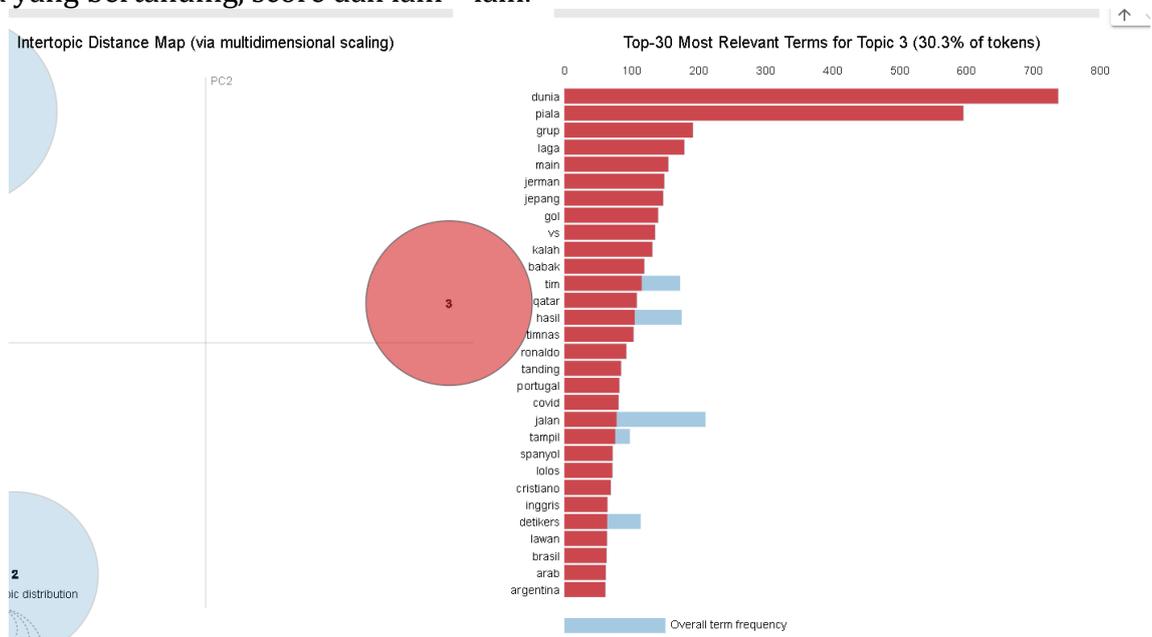
pada topik kedua ini hal yang sering atau dominan dibahas adalah mengenai tokoh politik dan kasus – kasus dalam pemerintahan indonesia, yaitu presiden, menteri, dpr, kpk, sidang, sambo dan lain – lain.



Gambar 4. Diagram Dominan Topik 2

3) Topik 3

Topik 3 menjadi topik akhir dari penelitian ini, yang berisikan informasi mengenai seputaran piala dunia di Qatar pada tahun 2022 lalu, yaitu berisikan grup, laga, negara yang bertanding, score dan lain – lain.



Gambar 5. Diagram Dominan Topik 3

4) Word Cloud

Word cloud juga dikenal sebagai cloud tag, adalah representasi visual dari teks tertulis yang diatur berdasarkan frekuensi. Teknik ini sering digunakan untuk

- Telkom Jl Terusan Buah Batu, U., Dayeuhkolot, K., & Barat, J. (2022). Analisis Persepsi Masyarakat Terhadap Komunikasi Kebijakan Menggunakan Topic Modelling (Kebijakan Protokol Kesehatan Covid-19 Dalam Penggunaan Masker). *Jurnal Sains Komputer & Informatika (J-SAKTI)*, 6(1), 253–266.
- [5] Indah Nurhafida, S., & Sembiring, F. (2021). Analisis Text Clustering Masyarakat Di Twiter Mengenai Mcdonald’Sxbts Menggunakan Orange Data Mining. *SISMATIK (Seminar Nasional Sistem Informasi Dan Manajemen Informatika)*, 28–35.
- [6] Kannitha, D. Z. T., Mustafid, M., & Kartikasari, P. (2022). Pemodelan Topik Pada Keluhan Pelanggan Menggunakan Algoritma Latent Dirichlet Allocation Dalam Media Sosial Twitter. *Jurnal Gaussian*, 11(2), 266–277. <https://doi.org/10.14710/j.gauss.v11i2.35474>
- [7] Kencana, W. H., Situmeang, I. V. O., & Januar, K. (n.d.). 1509-Article Text-2448-1-10-20211028. 6(2), 136–145.
- [8] Habibie, M. I., Widiaputra, T., & Yulianingsani, Y. (2022). Web Scraping of Disease Information From Social Media Twitter. *Jurnal Teknoinfo*, 16(2), 246. <https://doi.org/10.33365/jti.v16i2.1871>
- [9] Syaifuddin, A., & Muslimin, M. (2022). Analisis Sentimen Pada Sosial Media Tentang Implementasi Kebijakan Pse Kominfo Menggunakan Algoritme Lexicon Based. *Seminar Nasional Fakultas Teknik*, 1(1), 7–14. <https://doi.org/10.36815/semastek.v1i1.2>
- [10] Firdaus, M. R., Rizki, F. M., Gaus, F. M., & Susanto, I. K. (2020). Analisis Sentimen Dan Topic Modelling Dalam Aplikasi Ruangguru. *J-SAKTI (Jurnal Sains Komputer Dan Informatika)*, 4(1), 66. <https://doi.org/10.30645/j-sakti.v4i1.188>
- [11] Patmawati, P., & Yusuf, M. (2021). Analisis Topik Modelling Terhadap Penggunaan Sosial Media Twitter oleh Pejabat Negara. *Building of Informatics, Technology and Science (BITS)*, 3(3), 122–129. <https://doi.org/10.47065/bits.v3i3.1012>
- [12] Nurlyayli, A., & Nasichuddin, M. A. (2019). Topik Modeling Penelitian Dosen Jptei Uny Pada Google Scholar Menggunakan Latent Dirichlet Allocation. *Elinvo (Electronics, Informatics, and Vocational Education)*, 4(2), 154–161. <https://doi.org/10.21831/elinvo.v4i2.28254>
- [13] Khatib, J., Dalam, S., Satria, B., Sidauruk, A., Wardhana, R., Akbar, A. Al, Ihsan, A., Gama, A. M., Yogyakarta, U. A., Bengkulu, U. D., Selatan, P. A., & Kunci, K. (2022). *Indonesian Journal of Computer Science*. 11(1), 566–576.
- [14] Informasi, E. (2021). *Pawiyatan Universitas IVET* <http://e-journal.ikip-veteran.ac.id/index.php/pawiyatan>. 3, 1–9.
- [15] Suparyanto dan Rosad (2015). (2020). METODE LATENT DIRICHLET ALLOCATION UNTUK MENENTUKAN TOPIK TEKS SUATU BERITA. *Suparyanto Dan Rosad (2015)*, 5(3), 248–253.