

APLIKASI *CLUSTERING* BERITA DENGAN METODE *K MEANS* DAN PERINGKAS BERITA DENGAN METODE *MAXIMUM MARGINAL RELEVANCE*

Edy Susanto ¹⁾, Viny Christanti Mawardi ²⁾, Manatap Dolok Lauro ³⁾

¹⁾ Teknik Informatika, FTI, Universitas Tarumanagara
Jl. Letjen S Parman no 1, Jakarta 11440 Indonesia

edy.535170077@stu.untar.ac.id ¹⁾, viny@fti.untar.ac.id ²⁾, manataps@fti.untar.ac.id ³⁾

ABSTRACT

News is information about facts or opinions that are interesting to know. News can be obtained from various media such as newspapers and the internet. As is well known, news has various topics, such as politics, sports and others. There is also the same story written with the addition of a little information. This causes it to take more time to get the headline of the news. Therefore we need a system for news clustering using the K-Means method and news summarizing using the Maximum Marginal Relevance (MMR) method in order to obtain information from news more easily and efficiently. News that is processed in the form of a collection of files (multi document) with the extension txt. The summarization process goes through the text preprocessing stage, which consists of sentence segmentation, case folding, tokenizing, filtering, stemming. The next step is TF-IDF calculation to calculate word weight then Cosine Similarity to calculate the similarity between documents. After that, enter the K-Means stage for clustering division and proceed with determining the summary with MMR. Based on the results testing that has been done, this application is running well, the results of clustering and summarizing news can make it easier for users to get news summaries from some similar news.

Key words

Berita, Clustering, K-Means, Maximum Marginal Relevance, Peringkasan teks.

1. Pendahuluan

1.1 Latar Belakang

Berita adalah sebuah informasi mengenai fakta atau opini yang menarik orang untuk diketahui. Berita bisa didapatkan dari berbagai media seperti surat kabar dan internet. [1]

Seperti yang diketahui, berita memiliki berbagai macam topik, seperti tentang kejahatan, politik, olahraga dan lain-lain. Dengan demikian, untuk mempermudah pembaca berita diperlukan pengelompokan berdasarkan topik yang dibicarakan untuk menemukan informasi yang dibutuhkan. Akan tetapi terkadang ada juga berita yang sama dituliskan dengan penambahan sedikit informasi. Hal ini menyebabkan untuk mendapatkan informasi utama dari berita tersebut menggunakan waktu yang lebih banyak. Oleh karena itu diperlukan peringkasan kumpulan berita ini agar dapat memperoleh informasi dari berita lebih mudah dan efisien.

1.2 Rumusan Rancangan

Rancangan aplikasi yang digunakan untuk clustering berita dengan menggunakan metode K-Means dan peringkasan berita dengan menggunakan metode Maximum Marginal Relevance (MMR) ini berbasis website (browser) sehingga dapat dengan mudah diakses dan digunakan. Aplikasi ini memberikan ringkasan berita yang merupakan hasil dari peringkasan beberapa

dokumen berita yang sudah dikelompokkan dan ringkasan tersebut dalam bentuk ekstraksi.

Tampilan dari aplikasi ini dirancang dengan menggunakan bahasa pemrograman ASP.Net dan sedangkan untuk perhitungan clustering berita dan peringkasan berita menggunakan bahasa pemrograman Python.

1.3 Tujuan Rancangan

Tujuan dari perancangan aplikasi clustering dan peringkasan berita ini sebagai berikut :

1. Membuat sebuah aplikasi berbasis web untuk clustering berita dengan metode K-Means.
2. Membuat sebuah aplikasi berbasis web untuk peringkasan berita dengan metode Maximum Marginal Relevance.

2. Landasan Teori

2.1 Peringkasan Teks Otomatis

Peringkasan teks Otomatis adalah proses dimana akan diambil informasi dan point-point penting dari sebuah teks sehingga menghasilkan kesimpulan dari sumber teks tersebut. Sistem akan dimasukkan berupa teks dan menghasilkan sebuah ringkasan dari teks aslinya.

2.2 Clustering

Clustering adalah proses untuk mengelompokkan kelas dengan kesamaan objek-objek tertentu. Clustering dan klasifikasi adalah hal yang berbeda, dimana pada clustering tidak adanya target dalam melakukan pengelompokkan. Clustering biasanya dilakukan untuk proses awal dalam data mining untuk melakukan suatu analisis. Untuk melakukan clustering, banyak algoritma yang bisa digunakan untuk clustering, salah satunya adalah K-Means. Walaupun algoritmanya berbeda-beda namun tetap mempunyai prinsip yang sama, yaitu mengelompokkan data dan mengukur kemiripan antar data dalam satu kelompok. [2]

2.3 Text Mining

Text mining atau penambangan teks adalah proses ekstraksi atau pemisah pola yang berupa pengetahuan dan informasi-informasi yang penting dari sumber dokumen teks seperti dokumen PDF, dokumen Word dan sejenisnya. Proses yang sering dilakukan oleh teks mining di antaranya adalah meringkas otomatis dan mendeteksi plagiarisme.

2.4 Text Preprocessing

Text preprocessing merupakan langkah awal dalam melakukan persiapan data teks yang tidak terstruktur

menjadi lebih terstruktur untuk diolah lebih lanjut .[3] Adapun tahapan-tahapan dari proses Text Preprocessing, yaitu :

1. Segmentasi Kalimat

Pada proses segmentasi kalimat, teks berita akan dipisahkan berdasarkan tanda pemisah dalam suatu kalimat seperti tanda titik, tanda seru, dan tanda tanya yang nantinya akan digunakan untuk proses yang lebih lanjut.

2. Case Folding

Case Folding adalah proses dimana semua huruf dalam teks dokumen diubah menjadi huruf kecil dan semua tanda baca, nomor dan simbol akan dihilangkan.

3. Tokenizing

Tokenizing adalah proses dimana mengubah dari bentuk kalimat menjadi bentuk kata-kata tunggal dengan memisahkan berdasarkan spasi, hal ini bertujuan agar dapat melakukan proses *stemming*.

4. Filtering

Filtering yaitu proses pemilihan kata-kata penting yang bisa mewakili isi dari sebuah dokumen tersebut. Proses ini dilakukan dengan menghilangkan stopwords, yang bisa mempengaruhi hasil dari ringkasan dikarenakan bobot katanya yang besar. Stopword bisa berupa kata ganti, kata penghubung dan lain-lain. Seperti dia, yang, dari, dan lain-lain.

5. Stemming

Stemming merupakan proses mengubah kata menjadi kata dasar dengan menghilangkan kata-kata imbuhan seperti prefix, sufiks, infiks, dan konfiks pada setiap kata.

2.5 Term Weighting (Pembobotan)

Term Frequency atau TF merupakan menghitung bobot kata dengan menjumlahkan kata yang muncul pada dokumen tersebut. Sedangkan Inverse Document Frequency atau IDF merupakan jumlah kemunculan suatu kata pada semua dokumen yang ada. [4]

Untuk rumus TF dapat dilihat pada Persamaan 1.

$$Tf_{t,a} = \begin{cases} 1 + \log(tf_{t,a}) & \text{if } tf_{t,a} > 0 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Dimana tf adalah jumlah dari kemunculan sebuah kata dalam satu dokumen. Sedangkan untuk rumus IDF dapat dilihat pada Persamaan 2.

$$IDF_t = \log\left(\frac{N}{df_t}\right) \quad (2)$$

Dimana N adalah jumlah dari semua dokumen dalam data set dan df_t adalah jumlah dari dokumen yang mengandung Term t di dalamnya.

TF-IDF adalah hasil perkalian antara TF dengan IDF pada setiap katanya. Untuk rumus TF-IDF dapat dilihat pada Persamaan 3.

$$W_{t,d} = Tf_{t,d} * IDF_t = Tf_{t,d} * \log\left(\frac{N}{df_t}\right) \quad (3)$$

Nilai dari TF-IDF juga dapat dilakukan normalisasi untuk rumusnya dapat dilihat pada Persamaan 4.

$$W_{t,d \text{ norm}} = \frac{W_{t,d}}{\sqrt{\sum_{t=1}^n W_{t,d}^2}} \quad (4)$$

Dimana W adalah bobot Term t pada dokumen d, t adalah kata ke-t dari Term, d adalah dokumen ke-d, N adalah total dokumen, tf adalah jumlah kemunculan Term t pada dokumen d, df adalah jumlah dokumen yang memiliki Term t.

2.6 Cosine Similarity

Cosine similarity adalah tahapan dimana untuk menghitung kemiripan dari antar dokumen satu dengan dokumen lainnya. Nilai dari cosine similarity digunakan untuk mewakili kesamaan antara dokumen dalam melakukan peringkasan informasi, untuk rumus cosine similarity dapat dilihat pada Persamaan 5 dan Persamaan 6. [4]

1. Untuk Cosine Similarity tanpa normalisasi TF-IDF:

$$CosSim(d_j, q) = \frac{\sum_{i=1}^t (W_{i,j} \cdot W_{i,q})}{\sqrt{\sum_{i=1}^t W_{i,j}^2 \cdot \sum_{i=1}^t W_{i,q}^2}} \quad (5)$$

2. Untuk Cosine Similarity dengan normalisasi TF-IDF:

$$CosSim(d_j, q) = \sum_{i=1}^t (W_{i,j} \text{norm} \times W_{i,q} \text{norm}) \quad (6)$$

Dimana i adalah Term dalam Kalimat, $W_{i,j}$ adalah bobot Term i dalam blok j, $W_{i,q}$ adalah bobot Term i dalam blok q.

2.7 Algoritma Clustering K-Means

Algoritma K-Means cukup populer dikarenakan cukup mudah untuk diimplementasikan. Adapun kelemahan dari algoritma K-Means yaitu sensitif terhadap inisialisasi cluster. Untuk algoritma dari K-Means Clustering adalah sebagai berikut: [5]

1. Inisialisasi cluster k. Biasanya dilakukan dengan cara random. Namun pada penelitian ini untuk menginisialisasi jumlah k dapat menggunakan rumus pada Persamaan 7. Dimana n adalah jumlah total dokumen. [6]

$$k = \sqrt{n/2} \quad (7)$$

2. Masukkan dokumen ke cluster berdasarkan ukuran kedekatan dengan centroid. Dimana centroid merupakan vektor yang menjadi pusat dari cluster. Untuk rumus kedekatannya dapat dilihat pada persamaan 8.

$$d(P, Q) = \sqrt{\sum_{j=1}^p (x_j(P) - x_j(Q))^2} \quad (8)$$

3. Setelah seluruh dokumen masuk ke cluster, maka hitung ulang centroidnya berdasarkan dokumen yang ada di dalam cluster tersebut dengan mencari nilai rata-rata centroid dari tiap cluster.
4. Jika nilai dari centroidnya berubah maka kembali ke langkah 2. Jika nilai dari centroidnya tidak berubah, maka berhenti.

2.8 Maximum Marginal Relevance (MMR)

Maximum Marginal Relevance atau MMR adalah salah satu metode yang digunakan untuk meringkas dokumen tunggal maupun banyak dokumen. Maximum Marginal Relevance merupakan Teknik peringkasan yang tidak mengandung redundansi dan mengambil informasi yang relevan. Maximum Marginal Relevance meringkas dokumen dengan memperhitungkan kesamaan antara query dengan kalimat isi dokumen untuk mendapatkan nilai yang nantinya akan dijadikan ringkasan. Untuk rumus dari Maximum Marginal Relevance dapat dilihat pada Persamaan 9.[4]

$$MMR = \operatorname{argmax}[\lambda * Sim1(S_i, Q) - (1 - \lambda) * \max Sim2(S_i, S')] \quad (9)$$

Dimana λ adalah parameter yang mempengaruhi tingkat relevansi, S_i adalah vektor bobot kata yang menjadi kandidat, S' adalah vektor bobot kata lainnya selain kandidat, Q adalah vektor bobot kata dari query, $Sim1(S_i, Q)$ adalah nilai similarity antar kalimat ke-i dengan query dan $Sim2(S_i, S')$ adalah nilai similarity antar kalimat ke-i dengan kalimat hasil ekstraksi.

Nilai dari parameter λ merupakan 1 atau 0 atau antara ($0 < \lambda < 1$). Jika pada saat parameter λ sama dengan 1 maka nilai Maximum Marginal Relevance yang diperoleh akan cenderung relevan terhadap dokumen aslinya. Sedangkan ketika λ sama dengan 0 maka nilai dari Maximum Marginal Relevance yang diperoleh akan cenderung relevan dengan kalimat yang sudah diekstrak sebelumnya. Oleh sebab itu, nilai dari parameter λ perlu dioptimalkan agar mendapatkan hasil ringkasan yang baik dengan nilai λ berada antara 0 sampai 1. Pada peringkasan dengan dokumen yang sedikit, seperti artikel akan menghasilkan hasil ringkasan dengan baik jika nilai dari parameter λ adalah 0,7.[4]

2.9 Silhouette Coefficient

Silhouette Coefficient adalah salah satu metode yang digunakan untuk menguji kualitas dari sebuah cluster. Metode ini merupakan kombinasi dari metode *Cohesion* dan metode *Separation*. Metode *Cohesion* adalah metode yang digunakan untuk mengukur kedekatan relasi antar objek dalam satu cluster yang sama, sedangkan metode *Separation* adalah metode yang digunakan untuk mengukur kejauhan sebuah cluster terpisah dengan cluster yang lainnya. *Silhouette Coefficient* mempunyai tiga tahap dalam perhitungannya, Berikut adalah tahapan untuk perhitungan *Silhouette Coefficient*. [3]

1. Menghitung nilai rata-rata dari jarak objek dengan semua dokumen yang ada di dalam satu cluster dengan menggunakan rumus pada persamaan 10.

$$a(i) = \frac{1}{|A|-1} \sum_{j \in A, j \neq i} d(i, j) \tag{10}$$

Dimana a(i) merupakan perbedaan rata-rata objek (i) dengan semua objek lain pada A, d(i,j) adalah jarak antara data i dengan j, dan A adalah cluster.

2. Lalu menghitung jarak dari objek dengan semua dokumen antar cluster dengan menggunakan rumus pada persamaan 11.

$$d(i, C) = \frac{1}{|C|} \sum_{j \in C} d(i, j) \tag{11}$$

Dimana d(i,C) adalah perbedaan rata-rata objek (i) dengan semua objek lain pada C (cluster lain selain cluster A atau cluster C berbeda dengan A). Setelah menghitung d(i, C) untuk semua C, cari nilai minimum dengan menggunakan rumus pada persamaan 12.

$$b(i) = \min_{C \neq A} d(i, C) \tag{12}$$

Dimana cluster B yang mencapai minimum (yaitu, d(i,B) = b(i)) disebut tetangga dari objek (i). Ini merupakan cluster terbaik kedua untuk objek (i).

3. Kemudian hitung nilai *Silhouette Coefficient* dengan menggunakan rumus pada persamaan 13.

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))} \tag{13}$$

Nilai hasil *Silhouette Coefficient* mempunyai bervariasi antara -1 hingga 1. Jika nilai dari *Silhouette Coefficient* mendekati nilai 1, maka semakin baik pula pengelompokan data dalam satu cluster. Sebaliknya jika nilai dari *Silhouette Coefficient* nilai -1, maka semakin buruk pula pengelompokan data dalam satu cluster. Berikut adalah tabel tingkat kualitas struktur dari *Silhouette Coefficient*. [3]

Tabel 1 Nilai *Silhouette Coefficient*

Nilai <i>Silhouette Coefficient</i>	Struktur
0.7 < <i>Silhouette Coefficient</i> <= 1	Kuat
0.5 < <i>Silhouette Coefficient</i> <= 0.7	Sedang
0.25 < <i>Silhouette Coefficient</i> <= 0.5	Lemah
<i>Silhouette Coefficient</i> <= 0.25	Tidak Terstruktur

3. Hasil Pengujian

3.1 Pengujian Clustering

Hasil Pengujian akurasi dari clustering menggunakan Metode *Silhouette Coefficient* yang nantinya akan menampilkan seberapa baik kualitas dari hasil clustering tersebut.

Dari berita 1 sampai 5 tersebut didapatkan hasil clustering seperti gambar 1.

No	Judul	Cluster
1	4 Fakta Man United Vs Arsenal Autameyang Gerakan Selam Merah di Old Trafford.txt	1
2	4 Rancangan Bill Gates untuk Tahun 2021, dan Pandemi hingga Perubahan Kim.txt	1
3	101 Dokter Gugur karena COVID-19 (D) Waspada Lonjakan Selepas Libur Panjang.txt	1
4	412.784 Kasus Covid-19 di Indonesia dan Persentase Kemarian yang Masih di Atas Angka Duma.txt	2
5	AI Bisa Deteksi COVID-19 Lewat Suara Batah.txt	2

Gambar 1 Hasil Clustering Berita 1-5

Disini dijelaskan bahwa berita 1,2, dan 3 termasuk ke cluster 1. Dan berita 4 dan 5 termasuk ke cluster 2. Untuk hasil kualitas dari clusteringnya bisa dilihat pada gambar 2.

Dokumen	a(i)	b(i)	c(i)
D1	1.2265861780351	4.4778186544851	0.72560255840832
D2	1.15729416228448	3.91832414296307	0.7046532778856
D3	1.20251224809488	2.81404915218012	0.54426638070845
D4	0.468848046857075	4.1447747118748	0.8888213698415
D5	0.468848046857075	5.11327028933237	0.800307366304797

Silhouette Coefficient : 0.753954214262983

Gambar 2 Hasil *Silhouette Coefficient* Berita 1-5

Pada gambar diatas, didapatkan hasil *Silhouette Coefficient* (SC) adalah 0.75. Maka kualitas dari clustering berita 1 sampai 5 diatas adalah kuat, karena nilai dari *Silhouette Coefficient* lebih besar dari 0.7. Dari hasil clustering 5 berita diatas, jika dinilai secara manual terdapat kekurangan pada saat pembagian

cluster, dikarenakan berita 2 dan berita 3 jika dilihat dari judul seharusnya masuk ke cluster ke 2.

Sedangkan untuk berita 6 sampai 10 didapatkan hasil clustering seperti pada Gambar 3.

No	Judul	Cluster
1	Anggota DPRD Sumbang Meninggal Dunia Setelah Positif Corona.txt	1
2	Arsenal Dungham Manchester United, Takik Arreta Disebut Mirip Mourinho.txt	1
3	Bis WHO Tindak Teles COVID-19 Usak Kontak Cekat, In Awasannya.txt	1
4	Cara Benar Pakai dan Lepas Sarung Tangan di Masa Pandemi Covid-19.txt	1
5	Cegah Covid-19, Hindari Menyentuh Barang-barang In di Mal.txt	2

Gambar 3 Hasil Clustering Berita 6 - 10

Disini dijelaskan bahwa berita 6 sampai 9 (1-4) termasuk ke cluster 1. Dan sisanya termasuk ke cluster 2. Untuk hasil kualitas dari clusteringnya dapat dilihat pada gambar 4.

Hasil			
Dokumen	a(i)	b(i)	c(i)
D1	1.41313857910768	1.44750114307886	0.0233786358422089
D2	0.78213846684044	1.5326284087594	0.483151986857024
D3	0.832094117389903	1.33078070571257	0.57473611383535
D4	1.24204106380836	1.78552786445197	0.286504411818824
D5	0	0.05488157220897	1

Silhouette Coefficient : 0.435626429320278

Gambar 4 Hasil Silhouette Coefficient Berita 6 – 10

Pada gambar diatas, didapatkan hasil Silhouette Coefficient (SC) adalah 0.43. Maka kualitas dari clustering berita 6 sampai 10 di atas adalah lemah, karena nilai dari Silhouette Coefficient terletak diantara 0.25 dan 0.5. Dari hasil clustering 5 berita diatas, jika dinilai secara manual terdapat kekurangan pada saat pembagian cluster, dikarenakan pada centroid awal terletak pada berita pertama dan berita terakhir, sehingga untuk berita 2 masuk ke cluster 1.

3.2 Pengujian Peringkasan

Hasil pengujian akurasi dari hasil peringkasan menggunakan Question and Answering dikarenakan pada peringkasan tidak terdapat definisi ringkasan ideal. Maka dari itu diperlukan beberapa responden untuk menilai hasil ringkasan berita tersebut. Untuk nilai MMR dari berita 1 sampai 5 adalah sebagai berikut.

Iterasi	D1	D2	D3	D4	D5	D6	D7	D8
1	5.8410687068861	14.5388043543955	12.7058633804371	4.9449334110889	0	0	7.22459645107746	
2	0.973594132400415	7.11945217719777	3.63024668288202	1.23483948579358	-3.3396861218058	0	3.61229822553873	
3	0.64906275483381	4.74830145148518	2.42018445532135	0.82322832882386	-2.2264574145372	0	2.40819881702582	
4	0.486797086200207	3.55972608858888	1.81512334149101	0.617410742886789	-1.6888430090029	0	1.80614911278936	
5	0.389437852980186	2.84778097807911	1.45209867319281	0.403935794317432	-1.33587444872232	0	1.44491929021549	
6	0.324531737468805	2.37315072573259	1.21008222768067	0.41161316183193	-1.1132287072886	0	1.20409940851291	
7	0.278189752114404	2.03412918348508	1.03721333799486	0.352811281855308	-0.954198034801657	0	1.03208520729678	
8	0.243388533100104	1.77986304428944	0.907561670745505	0.308709871448395	-0.8402153045145	0	0.903074556384682	
9	0.216354251644537	1.58210048382173	0.808721485107115	0.274408774820795	-0.7421524715124	0	0.802732039080007	
10	0.194718826480083	1.42389043543955	0.726049338596404	0.248687897158716	-0.6878723431616	0	0.722459645107746	
11	0.177017114981894	1.28444585039859	0.680044851451276	0.224516270144287	-0.60721585851046	0	0.656781498552496	
12	0.162285888733402	1.18657536288629	0.605041113830337	0.205808580865596	-0.5596143538343	0	0.602049704256455	
13	0.149783712878887	1.0953003348935	0.558498488889541	0.188975305500704	-0.5137978648802	0	0.55573818854442	
14	0.13804878057202	1.01708459674254	0.518606888897431	0.178405940827854	-0.477098017400829	0	0.51604260384839	

Gambar 5 MMR Berita 1-5

Sedangkan untuk nilai MMR dari Berita 6-10 yaitu sebagai berikut.

Iterasi	D1	D2	D3	D4	D5	D6	D7	D8	D9	D10
1	0	0	8.15778387018422	58.7296288945919	0	22.2917351580519	0	0	0	0
2	0	0	-0.02882128204371	29.214813447286	0	11.1488678702829	0	-0.0481087911725	-0.5018300086338	-3.091847088795
3	0	0	-0.81774738029147	18.4785422881973	0	7.4057838801729	0	-2.8888738834115	-1.86788872730882	-2.06128888598833
4	0	0	-0.48331083102186	14.807486723848	0	5.57283378951287	0	-2.01746548755883	-1.25091504548189	-1.54586735448875
5	0	0	-0.37848584817488	11.8889253789184	0	4.45834703161037	0	-1.6133843800469	-1.00073203888535	-1.2387738835918
6	0	0	-0.308873754014574	9.7382714808886	0	3.7152881830884	0	-1.344838889170575	-0.83384338885446	-1.0304480298317
7	0	0	-0.26474882012482	8.3470885637028	0	3.184533584900741	0	-1.15288313574779	-0.714888897418189	-0.883488819851284
8	0	0	-0.23168531551093	7.30370381822089	0	2.7884688475848	0	-1.00870274377931	-0.82545752140845	-0.77288887724874
9	0	0	-0.205815888008716	6.4821807888577	0	2.48888888888888	0	-0.888888888888888	-0.688888888888888	-0.688888888888888
10	0	0	-0.185324282408744	5.84288888888888	0	2.22817351580519	0	-0.808888888888888	-0.500888888888888	-0.418388888888888
11	0	0	-0.168478888888888	5.31178426314472	0	2.02852137800472	0	-0.733888888888888	-0.454888888888888	-0.502168888888888
12	0	0	-0.1544388877007287	4.88913557454833	0	1.837844588888888	0	-0.672488888888888	-0.418888888888888	-0.515332451488883
13	0	0	-0.142557117237498	4.48458888888888	0	1.714748888888888	0	-0.620740150318839	-0.384888888888888	-0.475888888888888
14	0	0	-0.132874488888888	4.17584477818514	0	1.592288888888888	0	-0.576401888888888	-0.357404288888888	-0.441704888888888
15	0	0	-0.123548888888888	3.885308458888888	0	1.48811587720348	0	-0.537874788888888	-0.33357345481784	-0.412257881937288
16	0	0	-0.115827887755485	3.651851888888888	0	1.38832344737824	0	-0.504351371888888	-0.31228781370423	-0.388481838822437

Gambar 6 MMR Berita 6-10

Dimana dari nilai MMR maksimum tersebut akan dijadikan ringkasan. Untuk hasil ringkasan dari berita 1 sampai 5. Didapatkan 2 buah ringkasan. Dikarenakan terdapat 2 cluster berita yaitu berita 1 sampai 3 termasuk cluster 1 dan berita 4 sampai 5 termasuk cluster 2. Berikut ini adalah hasil ringkasan berita per cluster.

Tabel 2 Hasil Ringkasan Berita 1 - 5

Hasil Ringkasan Berita 1 sampai 5	
Cluster	Ringkasan
Cluster 1	KOMPAS Liga Man United vs Arsenal pada pekan ketujuh Liga Inggris 2020-2021 yang digelar di Stadion Old Trafford, Minggu (1/11/2020) berakhir dengan kemenangan untuk tim tamu Arsenal sukses menaklukkan Manchester United dengan skor tipis 1-0 Selain menjadi kekalahan ketiga Setan Merah di Liga Inggris musim ini, laga Manchester United vs Arsenal juga menyajikan sejumlah fakta menarik lainnya Berikut empat fakta laga Man United vs Arsenal yang dihimpun dari Opta
Cluster 2	JAKARTA, KOMPAS696 kasus baru Covid-19 dalam 24 jam terakhir Penambahan itu menyebabkan total kasus Covid-19 di Indonesia kini berjumlah 412899414 orang yang sudah diperiksa Dengan demikian, total pasien Covid-19 yang sudah sembuh dari virus corona terhitung sejak awal pandemi berjumlah 341

Sedangkan untuk hasil ringkasan dari berita 6 sampai 10 juga didapatkan 2 buah ringkasan. Karena juga terdapat 2 cluster berita yaitu berita 1 sampai 4 termasuk cluster 1 dan berita 5 termasuk cluster 2. Berikut ini adalah hasil ringkasan berita per cluster.

Tabel 3 Hasil Ringkasan Berita 6 - 10

Hasil Ringkasan Berita 6 sampai 10	
Cluster	Ringkasan
Cluster 1	PADANG, KOMPAS Politisi Partai Amanat Nasional (PAN) itu meninggal dunia, Sabtu (31/10/2020) pukul 23 Supardi menyebut, SF merupakan salah seorang anggota DPRD Sumbar yang aktif dan vokal dalam menyuarakan aspirasi konstituennya SF sebelum menjabat sebagai anggota DPRD Sumbar, terlebih dahulu menjadi anggota DPRD Dharmasraya Neville menilai kekuatan lini tengah menjadi kunci kemenangan Arsenal atas Manchester United "Jangan lupa, mencuci tangan dilakukan dengan air mengalir dan sabun selama 20 detik serta dikeringkan agar penggunaan sarung tangan tidak meningkatkan kelembaban dan menyebabkan jamur," kata Kesh "Saat melepaskan sarung tangan, jangan sampai merobeknya
Cluster 2	TEMPO Anggota Tim Pakar Universitas Lambung Mangkurat (ULM) untuk Percepatan Penanganan COVID-19, Prof Dr dr Syamsul Arifin MPd, mengingatkan masyarakat untuk mewaspadai transmisi COVID-19 di mal yang kerap diabaikan pengunjung "Hindari bersentuhan dengan beberapa tempat yang perlu diwaspadai jadi transmisi COVID-19 dan keharusan membawa sendiri penyaniitasi tangan selama kunjungan ke mal," katanya Kemudian, pegangan kereta barang atau keranjang merupakan area yang sering disentuh oleh para pembeli pada saat membawa barang belanjaan Oleh karena itu, pada saat menyentuhnya bisa menggunakan tisu dan membersihkan tangan dengan hand sanitizer setelah selesai digunakan Terakhir untuk tombol ATM, peneliti di Cina menemukan masing-masing tombol mengandung rata-rata 1

Dari data Survey yang disebarkan menggunakan Question and Answering yang berupa kuesioner. Didapatkan total responden sebanyak 52 orang. Responden diminta untuk menjawab 5 pertanyaan yang dari tiap ringkasan sehingga total yang harus dijawab sebanyak 20 pertanyaan.

Tabel 4 Akurasi MMR

Ya Ringkasan 1	Ya Ringkasan 2	Ya Ringkasan 3	Ya Ringkasan 4
156	192	196	205
Avg Ringkasan 1	Avg Ringkasan 2	Avg Ringkasan 3	Avg Ringkasan 4

60	73.84615385	75.38461538	78.84615385
AVG Total MMR	72.01923077		

Dari data tabel 4 dapat dilihat total yang menjawab ya pada ringkasan pertama sebanyak 156 dan rata-rata 60%, pada ringkasan kedua yang menjawab ya sebanyak 192 dan rata-rata 73.84%, pada ringkasan ketiga yang menjawab ya sebanyak 196 dan rata-rata 75.38%, dan pada ringkasan keempat yang menjawab ya sebanyak 205 dan rata-rata 78.84%. Dari keempat ringkasan tersebut dapat dilihat bahwa nilai akurasi dari MMR adalah 72.01%.

4. Kesimpulan

Berdasarkan pengujian program “Aplikasi Clustering Berita Dengan Metode K-Means Dan Peringkasan Berita Dengan Metode Maximum Marginal Relevance”, maka didapatkan kesimpulan sebagai berikut :

1. Pada pengujian Clustering berita 1 sampai 5 mendapatkan nilai Silhouette Coefficient yaitu 0.75, Maka kualitas dari clustering berita 1 sampai 5 diatas adalah kuat atau baik, karena nilai dari Silhouette Coefficient lebih besar dari 0.7.
2. Pada pengujian Clustering berita 6 sampai 10 mendapatkan nilai Silhouette Coefficient yaitu 0.43. Maka kualitas dari clustering berita 6 sampai 10 di atas adalah lemah, karena nilai dari Silhouette Coefficient terletak diantara 0.25 dan 0.5.
3. Pada pengujian peringkasan dari keempat ringkasan tersebut dapat dilihat bahwa nilai akurasi dari MMR adalah 72.01%.
4. Pada saat ada lebih dari 2 topik berita dari 3 sampai 5 berita, pembagian cluster tetap menjadi 2 kelompok atau cluster.
5. Batasan dari peringkasan ini ditentukan dari perolehan nilai MMR ketika nilai dari MMR maksimum dari tiap kalimat tersebut lebih besar dari 0 maka akan dijadikan ringkasan.
6. Banyaknya file dan besarnya ukuran file akan mempengaruhi lamanya proses komputasi.

REFERENSI

[1] Rofiqi, Ach. Yasir. “CLUSTERING BERITA OLAHRAGA BERBAHASA INDONESIA MENGGUNAKAN METODE K-MEDOID BERSYARAT”. Jurnal Simantec. Vol. VI, Nomor 1. Jawa Timur: Universitas Trunojoyo Madura, Juni 2017.

[2] Sukma Sindi; R. O. N. Weni; Sihombing, Irma Agustika; P.P.P.A.N.W. Fikrul Ilmi R.H.Zer; dan Hartama, Dedy. “ANALISIS ALGORITMA K-MEDOIDS CLUSTERING DALAM PENGELOMPOKAN PENYEBARAN COVID-19 DI INDONESIA”. Jurnal Teknologi Informasi. Vol. IV, Nomor 1. Pematangsiantar: STIKOM Tunas Bangsa Pematangsiantar, 2020.

- [3] Hudin, Muhammad Sholeh; Fauzi ,M Ali; dan Adinugroho, Sigit. “Implementasi Metode Text Mining dan K-Means Clustering untuk Pengelompokan Dokumen Skripsi (Studi Kasus: Universitas Brawijaya)”. *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*. Vol. II, Nomor 11. November 2018.
- [4] Saraswati, Nirmala Fa’izah; Indriati; dan Perdana, Rizal Setya; “Peringkasan Teks Otomatis Menggunakan Metode *Maximum Marginal Relevance* Pada Hasil Pencarian Sistem Temu Kembali Informasi Untuk Artikel Berbahasa Indonesia”. *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*. Vol. II, Nomor 11. Jawa Timur: Fakultas Ilmu Komputer Universitas Brawijaya. November 2018.
- [5] Husni; D. P. N. Yudha; dan M. Syarif. “Clusterisasi Dokumen Web (Berita) Bahasa Indonesia Menggunakan Algoritma K-Means”. *Jurnal SimanteC*. Vol. IV, Nomor 3. Madura: Fakultas Teknik Universitas Trunojoyo, Juni 2015.
- [6] Kodinariya, Trupti M. and Dr. Prashant R. Makwana. “Review on deTermining number of Cluster in K-Means Clustering”. *International Journal of Advance Research in Computer Science and Management Studies*. Vol. I, Nomor 6. November 2013.

Edy Susanto, Seorang mahasiswa pada program studi Fakultas Teknologi Informasi di Universitas Tarumanagara.

Viny Christanti, Seorang dosen tetap Program Studi Teknik Informatika Fakultas Teknologi Informasi Universitas Tarumanagara, Jakarta.

Manatap Dolok Lauro, Seorang dosen tetap Program Studi Teknik Informatika Fakultas Teknologi Informasi Universitas Tarumanagara, Jakarta.