

**IMPLEMENTASI *MODIFIED K-MEANS CLUSTERING* PADA
PENGELOMPOKAN DATA KOMENTAR SENTIPOL**

Skripsi



oleh

RUDDY CAHYANTO

71130097

**PROGRAM STUDI INFORMATIKA FAKULTAS TEKNOLOGI INFORMASI
UNIVERSITAS KRISTEN DUTA WACANA**

2020

**IMPLEMENTASI *MODIFIED K-MEANS CLUSTERING* PADA
PENGELOMPOKAN DATA KOMENTAR SENTIPOL**

Skripsi



Diajukan kepada Program Studi Informatika Fakultas Teknologi Informasi
Universitas Kristen Duta Wacana
Sebagai Salah Satu Syarat dalam Memperoleh Gelar
Sarjana Komputer

Disusun oleh

RUDDY CAHYANTO

71130097

PROGRAM STUDI INFORMATIKA FAKULTAS TEKNOLOGI INFORMASI
UNIVERSITAS KRISTEN DUTA WACANA

2020

PERNYATAAN KEASLIAN SKRIPSI

Saya menyatakan dengan sesungguhnya bahwa skripsi dengan judul:

IMPLEMENTASI MODIFIED K-MEANS CLUSTERING PADA PENGELOMPOKAN DATA KOMENTAR SENTIPOL

yang saya kerjakan untuk melengkapi sebagian persyaratan menjadi Sarjana Komputer pada pendidikan Sarjana Program Studi Informatika Fakultas Teknologi Informasi Universitas Kristen Duta Wacana, bukan merupakan tiruan atau duplikasi dari skripsi kesarjanaan di lingkungan Universitas Kristen Duta Wacana maupun di Perguruan Tinggi atau instansi manapun, kecuali bagian yang sumber informasinya dicantumkan sebagaimana mestinya.

Jika dikemudian hari didapati bahwa hasil skripsi ini adalah hasil plagiasi atau tiruan dari skripsi lain, saya bersedia dikenai sanksi yakni pencabutan gelar kesarjanaan saya.

Yogyakarta, 8 Januari 2020



RUDDY CAHYANTO

71130097

**HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS
AKHIR UNTUK KEPENTINGAN AKADEMIS SECARA ONLINE
UNIVERSITAS KRISTEN DUTA WACANA YOGYAKARTA**

Saya yang bertanda tangan di bawah ini:

NIM : 71130097
Nama : Ruddy Cahyanto
Prodi / Fakultas : Informatika / Fakultas Teknologi Informasi
Judul Tugas Akhir : Implementasi Modified K-Means Clustering pada
Pengelompokan Data Komentar Sentipol

bersedia menyerahkan Tugas Akhir kepada Universitas melalui Perpustakaan untuk keperluan akademis dan memberikan **Hak Bebas Royalti Non Eksklusif (*Non-exclusive Royalty-free Right*)** serta bersedia Tugas Akhirnya dipublikasikan secara online dan dapat diakses secara lengkap (full access).

Dengan Hak Bebas Royalti Noneklusif ini Perpustakaan Universitas Kristen Duta Wacana berhak menyimpan, mengalihmedia/formatkan, mengelola dalam bentuk database, merawat, dan memublikasikan Tugas Akhir saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta. Demikian pernyataan ini saya buat dengan sebenar-benarnya.

Yogyakarta, 22 Januari 2020



(71130097 – Ruddy Cahyanto)

HALAMAN PERSETUJUAN

Judul Skripsi : IMPLEMENTASI MODIFIED K-MEANS
CLUSTERING PADA PENGELOMPOKAN DATA
KOMENTAR SENTIPOL

Nama Mahasiswa : RUDDY CAHYANTO

N I M : 71130097

Matakuliah : Skripsi (Tugas Akhir)

Kode : TIW276

Semester : Gasal

Tahun Akademik : 2019/2020

Telah diperiksa dan disetujui di
Yogyakarta,
Pada tanggal 8 Januari 2020

Dosen Pembimbing I



Antonius Rachmat C., S.Kom., M.Cs.

Dosen Pembimbing II



Danny Sebastian, S.Kom., M.M., M.T.

HALAMAN PENGESAHAN

IMPLEMENTASI MODIFIED K-MEANS CLUSTERING PADA PENGELOMPOKAN DATA KOMENTAR SENTIPOL

Oleh: RUDDY CAHYANTO / 71130097

Dipertahankan di depan Dewan Penguji Skripsi
Program Studi Informatika Fakultas Teknologi Informasi
Universitas Kristen Duta Wacana - Yogyakarta
Dan dinyatakan diterima untuk memenuhi salah satu syarat memperoleh gelar
Sarjana Komputer
pada tanggal 17 Desember 2019

Yogyakarta, 8 Januari 2020
Mengesahkan,

Dewan Penguji:

1. Antonius Rachmat C., S.Kom., M.Cs.
2. Danny Sebastian, S.Kom., M.M., M.T.
3. Joko Purwadi, M.Kom
4. Willy Sudiarto Raharjo, S.Kom., M.Cs.

DUTA WACANA

Dekan

Ketua Program Studi


(Restyandito, S.Kom., MSIS., Ph.D.)


(Gloria Virginia, Ph.D.)

UCAPAN TERIMA KASIH

Dalam menyelesaikan penelitian ini, penulis mengucapkan terima kasih kepada Bapak Antonius Rachmat dan Bapak Danny Sebastian yang telah membimbing, mendampingi serta banyak membantu dalam proses penelitian. Penulis juga mengucapkan terima kasih kepada rekan-rekan dan keluarga yang selalu mendoakan dan memberi dukungan dalam menyelesaikan penelitian ini.

©UKDW

DAFTAR ISI

HALAMAN JUDUL.....	
UCAPAN TERIMA KASIH.....	vii
INTISARI.....	viii
DAFTAR ISI.....	ix
DAFTAR TABEL.....	xiii
DAFTAR GAMBAR.....	xv
BAB I.....	1
PENDAHULUAN.....	1
1.1. Latar Belakang Masalah.....	1
1.2. Perumusan Masalah.....	2
1.3. Batasan Masalah.....	2
1.4. Tujuan Penelitian.....	3
1.5. Manfaat Penelitian.....	3
1.6. Metodologi Penelitian.....	3
1.7. Sistematika Penulisan.....	4
BAB II.....	6
TINJAUAN PUSTAKA DAN LANDASAN TEORI.....	6
2.1. Tinjauan Pustaka.....	6
2.2. Landasan Teori.....	8
2.2.1. Text Mining.....	8
2.2.2. Term Weighting.....	10
2.2.3. Clustering.....	11
2.2.4. K-means Clustering.....	12
2.2.5. Modified K-means.....	13

2.2.6	<i>Purity</i>	17
BAB III	20
ANALISIS DAN PERANCANGAN SISTEM	20
3.1	Kebutuhan Penelitian	20
3.1.1	<i>Software</i>	20
3.1.2	<i>Hardware</i>	20
3.1.3	Data Penelitian	20
3.1.4	<i>Library</i>	20
3.2	Perancangan Alur Kerja Sistem	21
3.2.1	Data Input	21
3.2.2	Normalisasi	22
3.2.3	Tokenisasi	23
3.2.4	<i>Stopword Removal</i>	24
3.2.5	<i>Feature Generation</i>	26
3.2.6	<i>Feature Selection</i>	26
3.2.7	<i>Clustering</i>	27
3.2.8	Evaluation	31
3.3	Perancangan Basis Data	32
3.3.1	Tabel Komentar	32
3.3.2	Tabel TF-IDF <i>Features</i>	33
3.3.3	Tabel Kata	33
3.3.4	Tabel Hasil <i>Clustering</i>	34
3.4	Use Case Diagram	34
3.5	Use Case Description	35
3.5.1	Admin	35

3.5.2	<i>User</i>	37
3.6	Perancangan <i>User Interface</i>	38
3.6.1	Admin	38
3.6.2	<i>User</i>	45
3.7	Rancangan Pengujian	46
BAB IV		49
HASIL DAN PEMBAHASAN.....		49
4.1	Tampilan Form dan Pembahasan Program	49
4.1.1	Admin	49
4.1.2	<i>User</i>	58
4.2	Pengujian.....	61
4.2.1	Pengujian Metode <i>K-means</i>	61
4.2.2	Pengujian Metode <i>Modified K-means</i>	63
4.3	Analisis Perbandingan <i>Stemming</i> dan Tanpa <i>Stemming</i>	64
4.4	Analisis Perbandingan <i>K-means</i> dan <i>Modified K-means</i>	68
4.4.2	<i>Purity</i>	68
4.4.3	Pengaruh Faktor <i>Random Centroid</i> Awal <i>K-Means</i>	72
4.5	Analisis Nilai <i>Purity</i>	74
4.5.1	Pengujian pada Perbaikan Label.....	76
4.5.2	Pengujian pada Penambahan Jumlah <i>Cluster</i>	78
BAB V.....		82
KESIMPULAN DAN SARAN.....		82
5.1	Kesimpulan	82
5.2	Saran.....	83
DAFTAR PUSTAKA		84

LAMPIRAN.....	86
1. Kartu Konsultasi	86
2. Berita Acara Pendadaran	88
3. Formulir Perbaikan (Revisi) Skripsi	89
4. Source Code Program	90
▪ Input File Dataset	90
▪ Normalisasi	91
▪ Tokenisasi	95
▪ Stemming	96
▪ Stopword Removal.....	96
▪ Term Weighting	96
▪ Feature Selection.....	97
▪ Matrix Vector.....	97
▪ Get Centroid Awal K-Means	98
▪ Get Centroid Awal Modified K-Means	99
▪ Get Final Cluster	101
▪ Analisis Cluster	103

DAFTAR TABEL

Tabel 2. 1 Hasil pengujian waktu pemrosesan <i>clustering</i>	7
Tabel 2. 2 Tabel daftar penelitian	7
Tabel 2. 3 Contoh tabel tf-idf.....	14
Tabel 2. 4 Contoh hasil penghitungan jarak pada tiap dokumen	15
Tabel 3. 1 Contoh Komentar Sebelum dan Setelah Normalisasi.....	23
Tabel 3. 2 Contoh kata-kata dalam daftar <i>stopword</i> bahasa Indonesia.....	25
Tabel 3. 3 Tabel komentar	33
Tabel 3. 4 Tabel tfidf_ <i>features</i>	33
Tabel 3. 5 Tabel kata.....	33
Tabel 3. 6 Tabel hasil <i>clustering</i>	34
Tabel 3. 7 Rancangan pengujian tanpa <i>stemming</i>	47
Tabel 3. 8 Rancangan pengujian pakai <i>stemming</i>	47
Tabel 3. 9 Rancangan pengujian faktor <i>random</i>	48
Tabel 4. 1 Hasil Pengujian <i>K-Means</i> tanpa <i>Stemming</i>	61
Tabel 4. 2 Hasil Pengujian <i>K-Means</i> Menggunakan <i>Stemming</i>	62
Tabel 4. 3 Hasil pengujian <i>modified k-means</i> tanpa <i>stemming</i>	63
Tabel 4. 4 Hasil pengujian <i>modified k-means</i> menggunakan <i>stemming</i>	64
Tabel 4. 5 Persentase perbandingan rata-rata nilai <i>purity stemming</i> vs non- <i>stemming</i>	67
Tabel 4. 6 Persentase Perbandingan rata-rata Nilai <i>Purity K-Means</i> vs <i>Modified K-Means</i>	71
Tabel 4. 7 Hasil pengujian faktor <i>random</i> pada metode <i>k-means</i>	73
Tabel 4. 8 Hasil pengujian pada metode <i>modified k-means</i>	73
Tabel 4. 9 Hasil pengujian <i>modified k-means</i> menggunakan <i>stemming</i> dengan perbaikan label	76
Tabel 4. 10 Hasil pengujian <i>modified k-means</i> menggunakan <i>stemming</i> dengan perbaikan label dan jumlah <i>cluster</i> 4	78
Tabel 4. 11 Hasil pengujian <i>modified k-means</i> menggunakan <i>stemming</i> dengan perbaikan label dan jumlah <i>cluster</i> 5	79

Tabel 4. 12 Hasil pengujian <i>modified k-means</i> menggunakan <i>stemming</i> dengan perbaikan label dan jumlah <i>cluster</i> 6	79
Tabel 4. 13 Hasil pengujian <i>modified k-means</i> menggunakan <i>stemming</i> dengan perbaikan label dan jumlah <i>cluster</i> 7	80

©UKDW

DAFTAR GAMBAR

Gambar 2. 1 Proses knowledge discovery pada data mining.....	9
Gambar 2. 2 Contoh Penghitungan <i>Purity</i>	18
Gambar 3. 1 Perancangan Alur Kerja Sistem	21
Gambar 3. 2 Flowchart Input File	22
Gambar 3. 3 Flowchart Tokenisasi	24
Gambar 3. 4 Flowchart <i>Stopword Removal</i>	25
Gambar 3. 5 Flowchart Pembobotan Kata	26
Gambar 3. 6 Flowchart <i>Feature Selection</i>	27
Gambar 3. 7 Flowchart <i>K-means</i> Gambar 3. 8 Flowchart <i>Modified K-means</i>	28
Gambar 3. 9 Flowchart <i>Modified Centroid Selection</i>	30
Gambar 3. 10 Rancangan ER Diagram Sistem <i>Clustering</i>	32
Gambar 3. 11 Use Case Diagram <i>Clustering</i>	35
Gambar 3. 12 Rancangan Halaman Input File <i>Dataset</i>	38
Gambar 3. 13 Rancangan Halaman <i>Feature Generation</i>	39
Gambar 3. 14 Rancangan Halaman <i>Feature Selection</i>	40
Gambar 3. 15 Rancangan Halaman <i>Clustering</i>	41
Gambar 3. 16 Rancangan Halaman Informasi <i>Centroid</i> Awal Terpilih	42
Gambar 3. 17 Rancangan Halaman Hasil <i>Clustering</i>	43
Gambar 3. 18 Rancangan Halaman Analisis <i>Cluster</i>	44
Gambar 3. 19 Rancangan Halaman Input Komentar	45
Gambar 3. 20 Rancangan Halaman Hasil <i>Clustering</i>	46
Gambar 4. 1 Implementasi Tampilan Form Input <i>Dataset</i>	49
Gambar 4. 2 Implementasi Tampilan Form <i>Feature Generation</i>	50
Gambar 4. 3 Implementasi Tampilan Form <i>Feature Selection</i>	52
Gambar 4. 4 Implementasi Tampilan Form <i>Clustering</i>	53
Gambar 4. 5 Implementasi Tampilan Form <i>Centroid</i>	54
Gambar 4. 6 Implementasi Tampilan Form Hasil <i>Clustering</i>	56
Gambar 4. 7 Implementasi Tampilan Form Analisis <i>Cluster</i>	57

Gambar 4. 8 Implementasi Tampilan Form Input Komentar Baru - <i>User</i>	58
Gambar 4. 9 Implementasi Tampilan Form Hasil <i>Clustering</i> Komentar - <i>User</i> ...	59
Gambar 4. 10 Implementasi Tampilan Form Analisis <i>Cluster</i> - <i>User</i>	60
Gambar 4. 11 Grafik Perbandingan <i>Stemming</i> vs No <i>Stemming</i> pada <i>K-means</i> ...	65
Gambar 4. 12 Grafik Perbandingan Rata-Rata <i>Purity Stemming</i> vs No <i>Stemming</i> pada <i>K-means</i>	65
Gambar 4. 13 Grafik Perbandingan <i>Stemming</i> vs No <i>Stemming</i> pada <i>Modified K-</i> <i>means</i>	66
Gambar 4. 14 Grafik Perbandingan Rata-Rata <i>Purity Stemming</i> vs No <i>Stemming</i> pada <i>Modified K-means</i>	67
Gambar 4. 15 Grafik Perbandingan <i>K-means</i> vs <i>Modified K-means</i> Tanpa <i>Stemming</i>	68
Gambar 4. 16 Grafik Perbandingan Rata-Rata <i>Purity K-means</i> vs <i>Modified K-means</i> Tanpa <i>Stemming</i>	69
Gambar 4. 17 Grafik Perbandingan <i>K-means</i> vs <i>Modified K-means</i> Menggunakan <i>Stemming</i>	70
Gambar 4. 18 Grafik Perbandingan Rata-Rata <i>Purity K-means</i> vs <i>Modified K-means</i> Menggunakan <i>Stemming</i>	71
Gambar 4. 19 Contoh Cluster yang Memiliki Komentar Identik dengan Label Berbeda	75
Gambar 4. 20 Contoh <i>cluster</i> yang bisa dikelompokkan lagi	75
Gambar 4. 21 Grafik Perbandingan Nilai <i>Purity</i> pada <i>Clustering</i> Menggunakan Data Label Sebelum dan Setelah Perbaikan.....	77
Gambar 4. 22 Grafik Perbandingan Rata-Rata <i>Purity</i> pada Jumlah <i>Cluster</i> yang Berbeda	81

BAB I

PENDAHULUAN

1.1. Latar Belakang Masalah

Pada saat ini, media sosial tidak hanya menjadi sarana untuk menjalin pertemanan dengan banyak orang tanpa dibatasi ruang dan waktu. Media sosial juga banyak digunakan untuk hal-hal lain, misalnya sebagai sarana untuk mempromosikan produk hingga menjadi sarana untuk kampanye politik, seperti yang terjadi pada pemilihan presiden Republik Indonesia periode 2014-2019 yang mempertemukan Jokowi dan Prabowo sebagai lawan politik.

Rachmat dan Lukito (2016) pada jurnal penelitiannya yang berjudul “Sentipol: *Dataset* Sentimen Komentar pada Kampanye Pemilu Presiden 2014 dari Facebook Page” melakukan pembangunan *dataset* dengan mengumpulkan status dan komentar terhadap calon presiden Indonesia pada masa kampanye pemilu tahun 2014 dari facebook *page*. Penelitian tersebut menghasilkan sebuah *dataset* yang berisi 3400 komentar dari 68 status. Masing-masing komentar diberi label positif, negatif dan netral.

Pada jurnal penelitian tersebut, penulis juga telah mencoba melakukan klasifikasi pada *dataset* yang telah dihasilkan. Klasifikasi dilakukan dengan menggunakan metode *Naive Bayes* dan *Support Vector Machine* dimana keakuratan hasil klasifikasi tersebut masing-masing 83,32% dan 84,82%. Namun pada penelitian tersebut, Rachmat dan Lukito (2016) belum mencoba melakukan *clustering* pada *dataset* sentipol yang dibangun. Oleh karena itu, pada penelitian ini penulis akan mencoba menerapkan teknik *clustering* pada *dataset* sentipol.

Clustering merupakan salah satu teknik dalam data *mining* yang mengelompokkan suatu himpunan data ke dalam kelompok-kelompok (*clusters*) data yang serupa (Han & Kamber, 2006). Salah satu algoritma yang umum digunakan dalam *clustering* adalah algoritma *K-means*. Algoritma *K-means* sering digunakan dalam *clustering* karena kemudahan penggunaannya. Namun, kualitas

hasil akhir *clustering* dengan menggunakan algoritma *K-means* sangat bergantung pada pemilihan *centroid* (titik pusat *cluster*) awal (Alrabea, Senthilkumar, Al-Shalabi, & Bader, 2013). Pada algoritma *K-means*, pada umumnya pemilihan *centroid* awal dilakukan secara *random*, sehingga memungkinkan hasil akhir *clustering* berbeda meskipun menggunakan data uji yang sama (Fabregas, Gerardo, & Tanguilig III, 2017). Selain itu, jika *centroid* awal yang dipilih buruk, maka *cluster* yang dihasilkan bisa sangat tidak optimal.

Oleh karena itu, berdasarkan pada salah satu kelemahan algoritma *K-means* tersebut, Sujatha dan Sona (2013) melakukan penelitian mengenai algoritma *K-means* yang dimodifikasi, dengan menambahkan algoritma pada penentuan *centroid* awal. Hasilnya pun cukup baik, dalam proses *clustering*, algoritma yang digunakan membutuhkan jumlah iterasi dan waktu yang lebih sedikit dari algoritma *K-means* biasa. Berdasarkan pada penelitian yang dilakukan oleh Sujatha dan Sona (2013), penulis akan menggunakan algoritma *Modified K-means* pada sistem *clustering* yang akan dibangun, sehingga diharapkan mampu menghasilkan hasil *clustering* yang sama pada data uji yang sama pula. Selain itu, sistem yang dibangun juga diharapkan bisa menghasilkan *cluster* yang memiliki akurasi yang lebih baik dibandingkan algoritma *K-means* konvensional.

1.2. Perumusan Masalah

Rumusan masalah pada penelitian ini adalah sebagai berikut:

1. Bagaimana penerapan *Modified K-means Clustering* dalam mengklusterisasi data komentar pada *dataset* sentipol?
2. Bagaimana hasil analisis *clustering* data komentar pada *dataset* sentipol berdasarkan nilai *purity* yang dihasilkan?

1.3. Batasan Masalah

Sistem *clustering* sentipol yang akan dibangun memiliki beberapa batasan sebagai berikut:

1. *Dataset* sentipol diambil dari penelitian yang telah dilakukan oleh Rachmat dan Lukito (2016) melalui <http://ti.ukdw.ac.id/~crowd/dataset.php>.
2. *Dataset* yang akan diteliti menggunakan bahasa Indonesia.

3. Fokus penelitian hanya pada data komentar.
4. Daftar *stopword* diambil dari http://static.hikaruyuuki.com/wp-content/uploads/stopword_list_tala.txt.
5. *Stemming* dilakukan menggunakan library JSastrawi untuk bahasa pemrograman java.

1.4. Tujuan Penelitian

Adapun tujuan dari penelitian ini adalah sebagai berikut:

1. Menerapkan metode *Modified K-Means Clustering* dalam mengklasterisasi data komentar pada *dataset* sentipol.
2. Menganalisis hasil *clustering* data komentar pada *dataset* sentipol berdasarkan nilai *purity* yang dihasilkan.

1.5. Manfaat Penelitian

Manfaat penelitian yang akan dilakukan adalah sebagai berikut:

1. Mengetahui seberapa baik *dataset* sentipol yang telah dibangun untuk proses pembelajaran *clustering* berdasarkan analisis *cluster*.
2. Dengan menggunakan algoritma *Modified K-Means*, diharapkan hasil *clustering* bisa memiliki nilai *purity* yang lebih besar, serta menghasilkan *cluster* yang lebih konsisten dibanding algoritma *K-Means* konvensional.

1.6 Metodologi Penelitian

Adapun metodologi penelitian yang akan dilakukan adalah sebagai berikut:

1. Studi Pustaka dan Pengumpulan Data

Tahap pertama pada penelitian ini adalah studi pustaka. Pada tahap ini, penulis mencari dasar-dasar teori terkait topik penelitian dari berbagai sumber referensi, seperti buku, jurnal, internet dan lain-lain. Selain itu, pada tahap ini juga dilakukan pengambilan *dataset* sentipol yang akan digunakan sebagai *input* dari sistem yang akan dibangun melalui link berikut <http://ti.ukdw.ac.id/~crowd/dataset.php>.

2. Perancangan Sistem

Pada tahap ini, dilakukan perancangan sistem *clustering* yang meliputi perancangan alur kerja sistem, perancangan *database*, dan perancangan antarmuka sistem yang akan dibangun.

3. Implementasi

Tahap ini merupakan tahap implementasi atau pembuatan sistem dengan mengacu pada rancangan yang telah dibuat pada tahap sebelumnya.

4. Pengujian Sistem

Pada tahap ini dilakukan pengujian sistem, apakah sistem yang dibangun sudah berjalan dengan baik dan sesuai dengan rancangan awal. Selain itu juga akan dilakukan evaluasi hasil *clustering* dengan menggunakan metode *purity*.

5. Penulisan Laporan

Proses terakhir dalam penelitian ini adalah penulisan laporan Tugas Akhir. Laporan berisi Bab 1 sampai dengan Bab 5 sesuai dengan format yang telah ditetapkan oleh pihak Universitas Kristen Duta Wacana.

1.7 Sistematika Penulisan

Sistematika penulisan laporan tugas akhir ini terdiri dari 5 bab. Bab 1 merupakan pendahuluan, yang berisi penjabaran mengenai latar belakang, rumusan masalah, batasan masalah, dan tujuan dilakukannya penelitian. Selain itu, dalam bab ini juga menjabarkan mengenai metode penelitian yang akan dilakukan serta sistematika penulisan laporan tugas akhir.

Bab 2 berisi tinjauan pustaka dan landasan teori. Tinjauan pustaka berisi pemaparan mengenai penelitian-penelitian yang pernah dilakukan sebelumnya yang terkait dengan penelitian yang akan dilakukan oleh penulis. Sedangkan landasan teori berisi penjelasan mengenai teori-teori yang terkait dengan penelitian yang akan dilakukan.

Bab 3 berisi tentang metodologi penelitian yang akan dilakukan. Bab ini memaparkan mengenai kebutuhan sistem, alat-alat yang dibutuhkan dalam pembuatan sistem, dan perancangan sistem yang meliputi perancangan alur kerja

sistem, perancangan *database* serta perancangan antarmuka sistem. Selain itu, pada bab ini juga akan dijelaskan mengenai rancangan pengujian yang akan dilakukan.

Bab 4 berisi hasil dan pembahasan. Bab ini memaparkan hasil dan pembahasan sistem yang telah dibangun sesuai dengan rancangan sistem pada bab 3. Selain itu, pada bab ini juga memaparkan hasil pengujian yang telah dilakukan. Sedangkan pada bab 5 berisi kesimpulan terkait penelitian yang telah dilakukan dan juga saran untuk pengembangan penelitian selanjutnya.

©UKDW

BAB V

KESIMPULAN DAN SARAN

5.1 Kesimpulan

Berdasarkan hasil pengujian dan analisis *cluster* yang telah dilakukan, maka penulis dapat mengambil beberapa kesimpulan dari penelitian mengenai Implementasi *Modified K-means Clustering* pada Pengelompokan Data Komentar Sentipol sebagai berikut:

1. Penggunaan *stemming* mampu meningkatkan kualitas *cluster* yang dihasilkan. Pada metode *K-means*, penggunaan *stemming* meningkatkan kualitas *cluster* sebesar 1,7%. Sedangkan pada metode *Modified K-means*, penggunaan *stemming* meningkatkan kualitas *cluster* sebesar 1,9% berdasarkan pada rata-rata nilai *purity* yang dihasilkan.
2. Metode *Modified K-means* menghasilkan kualitas *cluster* atau nilai *purity* yang lebih baik dibandingkan metode *K-means* biasa. Pada *clustering* yang menggunakan *stemming*, *Modified K-Means* mampu menghasilkan persentase rata-rata nilai *purity* sebesar 42% sedangkan *K-Means* menghasilkan persentase rata-rata nilai *purity* sebesar 39,1%. Sementara pada *clustering* yang tidak menggunakan *stemming*, *Modified K-Means* mampu menghasilkan persentase rata-rata nilai *purity* sebesar 40,1% sedangkan persentase rata-rata nilai *purity* yang dihasilkan pada *K-Means* sebesar 37,4%.
3. Faktor *random* pada penentuan *centroid* awal juga mempengaruhi kualitas dan konsistensi hasil *cluster* yang ada. Dari 5 kali pengujian dengan data yang sama, metode *Modified K-Means* yang dalam penerapannya menggunakan algoritma *modified centroid selection* dalam menentukan *centroid* awal tidak mengalami perubahan sama sekali, baik dari segi *purity* maupun *cluster* yang dihasilkan karena tidak adanya faktor *random*. Sementara pada metode *K-Means* yang mana terdapat faktor *random* pada pemilihan *centroid* awal, selalu terjadi perubahan *cluster* pada setiap pengujian yang dilakukan sehingga mempengaruhi nilai *purity* yang dihasilkan.
4. Rata-rata persentase nilai *purity* yang dihasilkan dibawah 50% karena masalah pelabelan dan kurangnya jumlah *cluster* yang ditetapkan. Dari hasil pengujian

dengan label yang diperbaiki secara manual, terjadi peningkatan nilai *purity* sebesar 0,5%. Sedangkan pada pengujian penambahan jumlah *cluster*, dari jumlah *cluster* sebesar 3 sampai dengan 7, selalu terjadi peningkatan nilai *purity* setiap kali jumlah *cluster* ditambah.

5.2 Saran

Adapun saran untuk pengembangan yang lebih lanjut mengenai penelitian ini adalah sebagai berikut:

1. Perlu adanya penelitian lebih lanjut menggunakan *dataset* lain dengan topik yang berbeda sebagai perbandingan.
2. Melakukan pendekatan pengujian yang berbeda, misalnya dengan melakukan penghitungan *purity* secara manual atau tidak menggunakan label yang ada. Selain itu juga bisa dilakukan penelitian untuk mencari jumlah *cluster* yang terbaik.

© UKDW

DAFTAR PUSTAKA

- Alrabea, A., Senthilkumar, A. V., Al-Shalabi, H., & Bader, A. (2013, June). Enhancing K-Means Algorithm with Initial Cluster Centers Derived from Data Partitioning along the Data Axis with PCA. *Journal of Advances in Computer Networks*, 1(2), 137-142. doi:10.7763/JACN.2013.V1.28
- Fabregas, A. C., Gerardo, B. D., & Tanguilig III, B. T. (2017). Enhanced Initial Centroids for K-Means Algorithm. *International Journal of Information Technology and Computer Science*, 9(1), 26-33. doi:10.5815/ijitcs.2017.01.04
- Feldman, R., & Sanger, J. (2007). *The Text Mining Handbook: Advanced Approaches in Analyzing Unstructured Data*. United States of America: Cambridge University Press.
- Han, J., & Kamber, M. (2006). *Data Mining Concepts and Techniques* (2 ed.). San Francisco: Elsevier Inc.
- Kadhim, A. I., Cheah, Y. N., & Ahamed, N. H. (2014). Text Document Preprocessing and Dimension Reduction Techniques for Text Document Clustering. *International Conference on Artificial Intelligence with Applications in Engineering and Technology*, (hal. 69-73). Kinabalu.
- Maimon, O., & Rokach, L. (2010). *Data Mining and Knowledge Discovery Handbook* (2 ed.). New York, USA: Springer Science+Business Media. doi:10.1007/978-0-387-09823-4
- Manning, C. D., Raghavan, P., & Schütze, H. (2009). *An Introduction to Information Retrieval*. Cambridge, England: Cambridge University Press.
- Rachmat, A., & Lukito, Y. (2016). Sentipol: Dataset Sentimen Komentar pada Kampanye Pemilu Presiden Indonesia 2014 dari Facebook Page. *Konferensi Nasional Teknologi Informasi dan Komunikasi*, 218-228. Retrieved from <https://knastik.ukdw.ac.id/2016/makalah/artikel/e7-j1.pdf>
- Rudy. (2009). *Perbandingan Metode K-Means dan Hierarchical Agglomerative Clustering untuk Pengelompokan Dokumen Teks*. Duta Wacana Christian University. Special Region of Yogyakarta: Duta Wacana Christian University. Retrieved from <http://sinta.ukdw.ac.id>

Sujatha, S., & Sona, A. S. (2013, February). New Fast K-Means Clustering Algorithm using Modified Centroid Selection Method. *International Journal of Engineering Research & Technology*, 2(2), 1-9. Retrieved from <https://www.ijert.org/download/2353/new-fast-k-means-clustering-algorithm-using-modified-centroid-selection-method>

Tan, P.-N., Steinbach, M., Karpatne, A., & Kumar, V. (2013). *Introduction to Data Mining* (2 ed.). New York: Pearson.

©UKDW