

**PENERAPAN ALGORITMA *COMPLETE LINKAGE*
UNTUK *CLUSTERING* DOKUMEN TEKS**

SKRIPSI



Oleh:

LINA KUSUMA DEWI

22 08 4407



PROGRAM STUDI TEKNIK INFORMATIKA
FAKULTAS TEKNOLOGI INFORMASI
UNIVERSITAS KRISTEN DUTA WACANA
YOGYAKARTA

2012

**PENERAPAN ALGORITMA *COMPLETE LINKAGE*
UNTUK *CLUSTERING* DOKUMEN TEKS**

SKRIPSI



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



Disusun oleh:

LINA KUSUMA DEWI

22 08 4407

PROGRAM STUDI TEKNIK INFORMATIKA
FAKULTAS TEKNOLOGI INFORMASI
UNIVERSITAS KRISTEN DUTA WACANA
YOGYAKARTA

2012

PERNYATAAN KEASLIAN TUGAS AKHIR

Saya menyatakan dengan sesungguhnya bahwa tugas akhir dengan judul:

PENERAPAN ALGORITMA *COMPLETE LINKAGE* UNTUK *CLUSTERING*
DOKUMEN TEKS

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

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

Yogyakarta, 18 Oktober 2012



LINA KUSUMA DEWI

22 08 4407



HALAMAN PERSETUJUAN

Judul : Penerapan Algoritma Complete Linkage untuk Clustering Dokumen Teks
Nama : Lina Kusuma Dewi
NIM : 22 08 4407
Mata Kuliah : Skripsi (Tugas Akhir)
Kode : TIW276
Semester : Ganjil
Tahun Akademik : 2012/2013

Telah diperiksa dan disetujui
di Yogyakarta,
Pada Tanggal 1-Nov-2012

Dosen Pembimbing I



Budi Susanto, S.Kom., M.T.

Dosen Pembimbing II



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

HALAMAN PENGESAHAN

SKRIPSI

PENERAPAN ALGORITMA *COMPLETE LINKAGE* UNTUK *CLUSTERING*
DOKUMEN TEKS

Oleh: LINA KUSUMA DEWI / 22 08 4407

Dipertahankan di depan Dewan Penguji Skripsi
Program Studi Teknik Informatika Fakultas Teknologi Informasi
Universitas Kristen Duta Wacana – Yogyakarta

Dan dinyatakan diterima untuk memenuhi salah satu syarat
untuk memperoleh gelar

Sarjana Komputer

pada tanggal

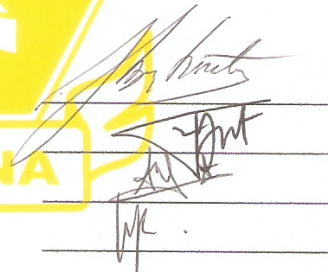
30 November 2012

Yogyakarta, 14/12/12

Mengesahkan,

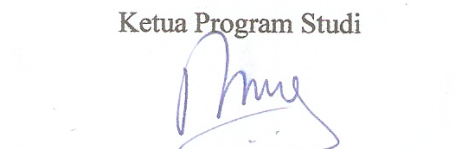
Dewan Penguji:

1. Budi Susanto, S.Kom., M.T.
2. Antonius Rachmat C., S.Kom., M.Cs.
3. Aditya Wikan Mahastama, S.Kom
4. Rosa Delima, S.Kom, M.Kom.



Dekan

(Drs. Wimmie Handiwiidjojo, MIT)

Ketua Program Studi

(Nugroho Agus Haryono, S.Si.,M.Si.)

UCAPAN TERIMA KASIH

Puji dan syukur penulis panjatkan kepada Tuhan kehadiran Tuhan Yang Maha Esa yang telah melimpahkan rahmat dan anugerah sehingga penulis dapat menyelesaikan Tugas Akhir dengan judul Penerapan Algoritma Complete Linkage untuk Clustering Dokumen Teks.

Penulisan laporan ini merupakan kelengkapan dan pemenuhan dari salah satu syarat dalam memperoleh gelar Sarjana Komputer. Selain itu bertujuan melatih mahasiswa untuk dapat menghasilkan suatu karya yang dapat dipertanggungjawabkan secara ilmiah, sehingga dapat bermanfaat bagi penggunaannya.

Dalam menyelesaikan pembuatan program dan laporan Tugas Akhir ini, penulis telah banyak menerima bimbingan, saran, dan masukan dari berbagai pihak, baik secara langsung maupun secara tidak langsung. Untuk itu dengan segala kerendahan hati, pada kesempatan ini penulis menyampaikan ucapan terima-kasih kepada :

1. Bpk. Budi Susanto, S.Kom., M.T. selaku dosen pembimbing I yang telah memberikan bimbingan dengan sabar dan baik kepada penulis, juga kepada
2. Bpk. Antonius Rachmat C., S.Kom, M.Cs selaku dosen pembimbing II atas bimbingan, petunjuk, dan saran yang diberikan selama pengerjaan Tugas Akhir ini sejak awal hingga akhir, juga kepada
3. Dosen-dosen Universitas Kristen Duta Wacana yang telah membantu memberikan pengarahan dan masukan kepada penulis.
4. Keluarga tercinta yang setia memberikan dukungan, doa, dan semangat.
5. Yanuar Budi Prasetyo atas saran, dukungan, dan semangat yang diberikan.
6. Pihak lain yang tidak dapat penulis sebutkan satu per satu, sehingga Tugas Akhir ini dapat terselesaikan dengan baik

Penulis menyadari bahwa program dan laporan Tugas Akhir ini masih jauh dari sempurna. Oleh karena itu, penulis sangat mengharapkan kritik dan

saran yang membangun dari pembaca sekalian. Sehingga suatu saat penulis dapat memberikan karya yang lebih baik lagi.

Akhir kata penulis ingin meminta maaf bila ada kesalahan baik dalam penyusunan laporan maupun yang pernah penulis lakukan sewaktu membuat program Tugas Akhir. Sekali lagi penulis memohon maaf yang sebesar-besarnya. Dan semoga dapat berguna bagi kita semua.

Yogyakarta, 18 Oktober 2012

Penulis



© UKDWN

INTISARI

Penerapan Algoritma *Complete Linkage* untuk *Clustering* Dokumen Teks

Perkembangan teknologi yang semakin cepat cenderung membawa dampak positif bagi kehidupan manusia. Manusia dapat memperoleh informasi dengan mudah melalui berbagai macam media. Semakin banyak informasi yang ada, semakin mudah informasi tersebut didapatkan, namun semakin sulit untuk mendapatkan informasi yang relevan. Proses seleksi atau pengelompokan informasi dibutuhkan untuk mendapatkan hasil yang relevan.

Melihat masalah di atas, penulis membangun sebuah aplikasi berbasis algoritma *Complete Linkage* untuk menemukan *cluster-cluster* pada dokumen teks. Proses untuk menemukan *cluster-cluster* tersebut menggunakan *Vector Space Model* untuk pembobotan dokumen, *Cosine Similarity* untuk menghitung nilai *similarity* antara dua dokumen, dan *thesaurus WordNet*. *WordNet* merupakan sebuah *database* yang menyimpan sinonim kata berbahasa Inggris. Dengan penggunaan *WordNet* diharapkan dapat meningkatkan kualitas *cluster* yang terbentuk.

Sistem yang dibangun telah mampu menemukan *cluster-cluster* dengan nilai *purity* sebesar 0,8308 untuk penggunaan *WordNet* dengan *cutting point* 0,02. Dengan uji coba tanpa menggunakan *WordNet* untuk *cutting point* yang sama, sistem menghasilkan *cluster-cluster* dengan nilai *purity* sebesar 0,9077.

Kata Kunci: *cluster*, *Complete Linkage*, *Cosine Similarity*, *cutting point*, dokumen teks, *purity*, *Vector Space Model*, *WordNet*

DAFTAR ISI

HALAMAN JUDUL	
PERNYATAAN KEASLIAN TUGAS AKHIR.....	i
HALAMAN PERSETUJUAN	ii
HALAMAN PENGESAHAN	iii
UCAPAN TERIMA KASIH	iv
INTISARI	vi
DAFTAR ISI	vii
DAFTAR TABEL	xi
DAFTAR GAMBAR DAN GRAFIK	xii
BAB I PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	2
1.3 Batasan Masalah	2
1.4 Hipotesis	2
1.5 Tujuan Penelitian	3
1.6 Metode Penelitian	3
1.7 Sistematika Penulisan	3
BAB II TINJAUAN PUSTAKA	5
2.1 Tinjauan Pustaka.....	5
2.2 Landasan Teori	6
2.2.1 <i>Clustering</i>	6
2.2.2 <i>Feature Selection</i>	7
2.2.3 <i>Frequency-based Feature Selection</i>	7
2.2.4 <i>Hierarchical Agglomerative Clustering (HAC)</i>	8
2.2.5 <i>Complete Linkage Algorithm</i>	9
2.2.6 <i>Vector Space Model</i>	12
2.2.7 <i>Purity</i>	14
BAB III ANALISIS DAN PERANCANGAN SISTEM	16
3.1 Kebutuhan Perangkat Keras dan Perangkat Lunak	16

3.2	Kamus Data	17
3.2.1	Tabel Documents	17
3.2.2	Tabel Files	17
3.2.3	Tabel Stopword	18
3.2.4	Tabel Training	18
3.2.5	Tabel New_Files	19
3.2.6	Tabel Selection	19
3.2.7	Tabel Similarity	20
3.3	Skema Database	20
3.4	Diagram Use Case	21
3.5	Diagram Alir Sistem.....	22
3.5.1	Tokenisasi dan Menghilangkan <i>Stopword</i>	23
3.5.2	Menghitung <i>Tf-Idf Weighting</i>	23
3.5.3	Menghitung <i>Similarity</i> Antar Dokumen	26
3.5.4	<i>Clustering</i> Dokumen Teks	27
3.6	Rancangan Antarmuka Sistem	28
3.6.1	Rancangan Antarmuka <i>Form</i> Main Menu	28
3.6.1.1	Rancangan Antarmuka <i>Form</i> Upload	29
3.6.1.2	Rancangan Antarmuka <i>Form</i> Stopword	30
3.6.2	Rancangan Antarmuka <i>Form</i> Process	30
3.6.2.1	Rancangan Antarmuka <i>Form</i> Replaced Word	30
3.6.2.2	Rancangan Antarmuka <i>Form</i> Weight	31
3.6.2.3	Rancangan Antarmuka <i>Form</i> Result.....	32
3.6.3	Rancangan Antarmuka <i>Form</i> Hasil Cluster	33
3.7	Contoh Perhitungan Manual Sistem	34
3.7.1	Perhitungan <i>Tf-Idf Weighting</i>	34
3.7.2	Perhitungan <i>Similarity</i> Dokumen.....	39
3.7.3	<i>Clustering</i> Dokumen Teks Menggunakan <i>Complete Linkage</i> .	40
3.8	Rancangan Pengujian Sistem	48
3.7.1	Pengujian terhadap Keefektifan Penggunaan <i>WordNet</i> untuk <i>Clustering</i> Dokumen Teks	42

3.7.2 Pengujian terhadap Hasil <i>Clustering</i> Dokumen Teks Tanpa Menggunakan <i>WordNet</i>	42
BAB IV IMPLEMENTASI DAN ANALISIS SISTEM	44
4.1 Implementasi Sistem	44
4.1.1 <i>Form</i> Main Menu	44
4.1.2 <i>Form</i> Process	46
4.1.3 <i>Form</i> Matrix	49
4.1.4 <i>Form</i> Hasil Cluster	50
4.2 Implementasi Proses	50
4.2.1 Menghitung Nilai <i>Tf</i> Baru	50
4.2.2 Menghitung <i>Df</i> , <i>Idf</i> , dan <i>Weight Term</i>	52
4.2.3 Menghitung <i>Magnitude</i> Masing-Masing Dokumen	53
4.2.4 Menghitung <i>Dot Product</i> , <i>Magnitude</i> , dan <i>Similarity</i> Antar Dokumen	53
4.2.5 Mengelompokkan Dokumen Teks dengan Menggunakan <i>Complete Linkage Algorithm</i>	54
4.3 Analisis Sistem	55
4.3.1 Pengujian dengan <i>WordNet</i>	55
4.3.1.1 Menggunakan <i>Cutting Point</i> Sebesar 0,30	57
4.3.1.2 Menggunakan <i>Cutting Point</i> Sebesar 0,25	58
4.3.1.3 Menggunakan <i>Cutting Point</i> Sebesar 0,20	58
4.3.1.4 Menggunakan <i>Cutting Point</i> Sebesar 0,02	59
4.3.2 Pengujian Tanpa <i>WordNet</i>	61
4.3.2.1 Menggunakan <i>Cutting Point</i> Sebesar 0,30	62
4.3.2.2 Menggunakan <i>Cutting Point</i> Sebesar 0,25	63
4.3.2.3 Menggunakan <i>Cutting Point</i> Sebesar 0,20	64
4.3.2.4 Menggunakan <i>Cutting Point</i> Sebesar 0,02	65
4.3.3 Pengujian Pengaruh Feature Selection Terhadap Nilai Purity dengan <i>Cutting Point</i> 0,02	67
BAB V KESIMPULAN DAN SARAN	69
5.1 Kesimpulan	69

5.2 Saran	69
DAFTAR PUSTAKA	68

© UKDW

DAFTAR TABEL

Tabel 3.1	Tabel Documents	17
Tabel 3.2	Tabel Files	18
Tabel 3.3	Tabel Stopword	18
Tabel 3.4	Tabel Training	18
Tabel 3.5	Tabel New_Files	19
Tabel 3.6	Tabel Selection	19
Tabel 3.7	Tabel Similarity	20
Tabel 3.8	Tabel New_Files Setelah Perhitungan.....	37
Tabel 3.9	Tabel Selection Setelah Proses <i>Feature Selection</i>	38
Tabel 3.10	Tabel Documents Setelah Proses Perhitungan	39
Tabel 3.11	Tabel Similarity Setelah Proses Perhitungan.....	40
Tabel 3.12	Tabel Matriks Cluster	40
Tabel 3.13	Tabel Matriks Baru	41
Tabel 3.14	Tabel Hasil Cluster.....	41
Tabel 4.1	Percobaan untuk Menentukan Nilai <i>Minimum Feature Selection</i> dengan <i>WordNet</i>	56
Tabel 4.2	<i>Cluster</i> Baru dengan <i>Cutting Point</i> 0,30	57
Tabel 4.3	<i>Cluster</i> Baru dengan <i>Cutting Point</i> 0,25	58
Tabel 4.4	<i>Cluster</i> Baru dengan <i>Cutting Point</i> 0,20	59
Tabel 4.5	<i>Cluster</i> Baru dengan <i>Cutting Point</i> 0,02	60
Tabel 4.6	Percobaan untuk Menentukan Nilai <i>Minimum Feature Selection</i> Tanpa <i>WordNet</i>	61
Tabel 4.7	<i>Cluster</i> Baru dengan <i>Cutting Point</i> 0,30	63
Tabel 4.8	<i>Cluster</i> Baru dengan <i>Cutting Point</i> 0,25	63
Tabel 4.9	<i>Cluster</i> Baru dengan <i>Cutting Point</i> 0,20	64
Tabel 4.10	<i>Cluster</i> Baru dengan <i>Cutting Point</i> 0,02	65
Tabel 4.11	<i>Purity</i> Hasil Pengujian.....	66
Tabel 4.12	Pengujian Pengaruh <i>Feature Selection</i> Terhadap Nilai <i>Purity</i> dengan <i>Cutting Point</i> 0,02.....	67

DAFTAR GAMBAR DAN GRAFIK

Gambar 2.1	Perbandingan <i>Feature Selection</i>	8
Gambar 2.2	Contoh <i>Dendogram</i>	9
Gambar 2.3	Contoh <i>Complete Linkage Clustering</i>	10
Gambar 2.4	Contoh <i>Single Linkage Clustering</i>	11
Gambar 2.5	Perbedaan <i>Complete Linkage</i> dengan Algoritma Lain	11
Gambar 2.6	<i>Complete Linkage</i> dengan <i>Outlier</i>	12
Gambar 2.7	Ilustrasi <i>Cosine Similarity</i>	13
Gambar 2.8	<i>Pseudocode</i> untuk <i>Cosine Similarity</i>	14
Gambar 2.9	Contoh Perhitungan <i>Purity</i>	15
Gambar 3.1	Skema <i>Database</i>	20
Gambar 3.2	Diagram <i>Use Case</i>	21
Gambar 3.3	Diagram Alir Sistem	22
Gambar 3.4	Diagram Alir Menghitung <i>Tf</i> Awal Dokumen	23
Gambar 3.5	Diagram Alir Menghitung <i>Tf</i> Baru	24
Gambar 3.6	Diagram Alir Menghitung <i>Df</i> , <i>Idf</i> , <i>Weight Term</i>	25
Gambar 3.7	Diagram Alir Menghitung <i>Magnitude</i> Masing-Masing Dokumen	26
Gambar 3.8	Diagram Alir Menghitung <i>Dot Product</i> , <i>Magnitude</i> , dan <i>Similarity</i> Antar Dokumen	27
Gambar 3.9	<i>Clustering</i> Dokumen Teks dengan <i>Complete Linkage</i>	28
Gambar 3.10	Rancangan Antarmuka <i>Form Upload</i>	29
Gambar 3.11	Rancangan Antarmuka <i>Form Stopword</i>	30
Gambar 3.12	Rancangan Antarmuka <i>Form Replaced Word</i>	31
Gambar 3.13	Rancangan Antarmuka <i>Form Weight</i>	32
Gambar 3.14	Rancangan Antarmuka <i>Form Result</i>	32
Gambar 3.15	Rancangan Antarmuka <i>Form Hasil Cluster</i>	33
Gambar 3.16	Ilustrasi Perhitungan <i>Tf</i> Awal	34
Gambar 3.17	Ilustrasi Pergantian Kata.....	35
Gambar 3.18	Ilustrasi <i>Feature Selection</i>	38

Gambar 3.19 Hierarki <i>Tree</i>	41
Gambar 4.1 <i>Form</i> Main Menu 1	44
Gambar 4.2 <i>Form</i> Main Menu 2.....	45
Gambar 4.3 <i>Form</i> Main Menu 3.....	46
Gambar 4.4 <i>Form</i> Process 1	47
Gambar 4.5 <i>Form</i> Process 2	48
Gambar 4.6 <i>Form</i> Process 3	48
Gambar 4.7 <i>Form</i> Matrix	49
Gambar 4.8 <i>Form</i> Hasil Cluster	50
Gambar 4.9 <i>Pseudocode</i> Menghitung Nilai <i>Tf</i> Baru	51
Gambar 4.10 <i>Pseudocode</i> Menghitung <i>Df</i> dan <i>Idf</i>	52
Gambar 4.11 <i>Pseudocode</i> Menghitung <i>Weight Term</i>	52
Gambar 4.12 <i>Pseudocode</i> Menghitung <i>Magnitude</i> Dokumen	53
Gambar 4.13 <i>Pseudocode</i> Menghitung <i>Dot Product</i> , <i>Magnitude</i> , dan <i>Similarity</i> Antar Dokumen	54
Gambar 4.14 <i>Pseudocode</i> Mengelompokkan Dokumen Teks dengan <i>Complete</i> <i>Linkage Algorithm</i>	55
Grafik 4.1 Grafik Percobaan <i>WordNet</i> dengan Nilai <i>Max Similarity</i> Sebagai Pembanding	56
Grafik 4.2 Grafik Percobaan <i>Non-WordNet</i> dengan Nilai <i>Max Similarity</i> Sebagai Pembanding.....	62
Grafik 4.3 Grafik <i>Purity</i> Hasil Pengujian	66
Grafik 4.4 Grafik Pengaruh <i>Feature Selection</i> dengan Menggunakan <i>WordNet</i> Terhadap Nilai <i>Purity</i>	68
Grafik 4.5 Grafik Pengaruh <i>Feature Selection</i> Tanpa Menggunakan <i>WordNet</i> Terhadap Nilai <i>Purity</i>	68

INTISARI

Penerapan Algoritma *Complete Linkage* untuk *Clustering* Dokumen Teks

Perkembangan teknologi yang semakin cepat cenderung membawa dampak positif bagi kehidupan manusia. Manusia dapat memperoleh informasi dengan mudah melalui berbagai macam media. Semakin banyak informasi yang ada, semakin mudah informasi tersebut didapatkan, namun semakin sulit untuk mendapatkan informasi yang relevan. Proses seleksi atau pengelompokan informasi dibutuhkan untuk mendapatkan hasil yang relevan.

Melihat masalah di atas, penulis membangun sebuah aplikasi berbasis algoritma *Complete Linkage* untuk menemukan *cluster-cluster* pada dokumen teks. Proses untuk menemukan *cluster-cluster* tersebut menggunakan *Vector Space Model* untuk pembobotan dokumen, *Cosine Similarity* untuk menghitung nilai *similarity* antara dua dokumen, dan *thesaurus WordNet*. *WordNet* merupakan sebuah *database* yang menyimpan sinonim kata berbahasa Inggris. Dengan penggunaan *WordNet* diharapkan dapat meningkatkan kualitas *cluster* yang terbentuk.

Sistem yang dibangun telah mampu menemukan *cluster-cluster* dengan nilai *purity* sebesar 0,8308 untuk penggunaan *WordNet* dengan *cutting point* 0,02. Dengan uji coba tanpa menggunakan *WordNet* untuk *cutting point* yang sama, sistem menghasilkan *cluster-cluster* dengan nilai *purity* sebesar 0,9077.

Kata Kunci: *cluster*, *Complete Linkage*, *Cosine Similarity*, *cutting point*, dokumen teks, *purity*, *Vector Space Model*, *WordNet*

BAB I

PENDAHULUAN

1.1 Latar Belakang Masalah

Perkembangan teknologi yang semakin cepat cenderung membawa dampak positif bagi kehidupan manusia. Manusia dapat memperoleh informasi dengan mudah melalui berbagai macam media. Dengan kemudahan yang ditawarkan, proses untuk mendapatkan informasi atau data tidak menghabiskan banyak waktu. Namun, informasi dan data-data tersebut belum tentu memenuhi kriteria yang diinginkan. Semakin banyak informasi yang ada, semakin mudah informasi tersebut didapatkan, namun semakin sulit untuk mendapatkan informasi yang benar-benar dibutuhkan.

Untuk memperoleh informasi yang benar-benar sesuai, diperlukan proses seleksi terhadap semua informasi yang tersedia. Proses seleksi ini dapat dilakukan secara *manual* dengan membandingkan semua informasi tersebut. Akan tetapi, hal ini tidak disarankan karena proses seleksi tidak optimal dan bila informasi yang ada mencapai ribuan dokumen, waktu yang dibutuhkan untuk mengelompokkan dokumen akan sangat lama dan tidak menutup kemungkinan terjadi kesalahan dalam proses seleksi. Proses seleksi dapat dilakukan secara otomatis untuk mempermudah *user* dalam mencari dokumen yang berkualitas.

Proses seleksi atau *clustering* adalah suatu proses untuk menentukan kategori. Dokumen yang memiliki isi yang relatif sama dapat dikelompokkan dalam satu kelompok yang sama. Dokumen yang dimaksud dapat berupa berbagai macam dokumen, misalnya dokumen perpajakan, jurnal, dokumen keuangan, dsb. Dalam implementasinya di komputer, dokumen yang dimaksud adalah dokumen *plain text* yang artinya dokumen yang tidak memiliki format apapun. Penelitian ini memadukan metode *Vector Space Model* dengan algoritma *Complete Linkage*

guna memperoleh *cluster* dokumen yang memiliki kesamaan topik. Selain itu, penelitian ini menggunakan *thesaurus* dari *WordNet* guna meningkatkan kualitas *cluster* yang terbentuk.

1.2 Perumusan Masalah

Penelitian ini berfokus tentang seberapa akurat *clustering* dokumen teks dengan mengimplementasikan algoritma *Complete Linkage* dipadukan dengan kamus kata *WordNet*.

1.3 Batasan Masalah

Program *clustering* yang akan dibuat ini memiliki beberapa batasan masalah seperti berikut ini :

- Sistem dibuat dalam bentuk aplikasi berbasis *desktop* dengan dokumen teks berupa *file plain text* berekstensi *.txt*
- Korpus dokumen merupakan sekumpulan berita berbahasa Inggris yang diambil secara acak dari <http://www.voanews.com>
- Pembobotan kata diperoleh dengan menggunakan *tf-idf weighting* dan *similarity measure* diperoleh dengan *Cosine Similarity*
- Sistem menggunakan *thesaurus WordNet* yang diambil dari www.androidtech.com/downloads/wordnet20-from-prolog-all-3.zip untuk meningkatkan kualitas *cluster*

1.4 Hipotesis

Clustering dokumen teks dengan menggunakan *thesaurus* dari *WordNet* mampu memberikan hasil *cluster* dengan tingkat akurasi sebesar 75%.

1.5 Tujuan Penelitian

Penelitian ini memiliki tujuan untuk mengelompokkan dokumen teks dengan berbagai macam topik menjadi sekumpulan dokumen yang memiliki topik yang sama. Fokus penelitian ini yakni untuk meneliti tingkat keakuratan *cluster* yang terbentuk dengan menggunakan *thesaurus WordNet*.

1.6 Metode / Pendekatan

Metode yang digunakan dalam penelitian ini antara lain sebagai berikut :

a. Metode Pengumpulan Data

Data-data yang digunakan dalam penelitian ini didapatkan penulis dengan cara melakukan studi pustaka melalui literatur-literatur yang mendukung penyelesaian penelitian ini. Selain itu, penulis juga mengumpulkan beberapa dokumen teks untuk dijadikan sebagai korpus dokumen. Data dari dokumen teks diambil dari sebuah situs berita berbahasa Inggris (www.voanews.com).

b. Metode Pengembangan Sistem

Sistem akan dikembangkan dengan menggunakan algoritma *complete linkage*. Algoritma ini merupakan salah satu algoritma dalam *hierarchical clustering*.

c. Metode Evaluasi

Metode evaluasi yang digunakan adalah dengan menghitung nilai *purity cluster* untuk beberapa *cutting point* tertentu .

1.7 Sistematika Penulisan

Untuk memudahkan dalam mendapatkan gambaran yang lengkap dan jelas mengenai penelitian yang akan dilakukan, penulis membagi laporan ini menjadi 5

(lima) bab yaitu Bab I Pendahuluan, Bab II Tinjauan Pustaka, Bab III Analisis dan Perancangan Sistem, Bab IV Implementasi dan Analisis Sistem, dan Bab V Kesimpulan dan Saran.

Dalam bab pertama menguraikan hal-hal seperti latar belakang masalah, perumusan masalah, batasan masalah, hipotesis, tujuan penelitian, metode/pendekatan yang digunakan serta sistematika penulisan laporan Tugas Akhir.

Bab kedua berisi tentang tinjauan pustaka serta landasan teori yang diperlukan untuk memecahkan masalah dan merumuskan hipotesis.

Bab ketiga berisikan tentang analisis teori yang digunakan dalam penelitian, uraian tentang variabel dan data yang akan dikumpulkan, flowchart dan arsitektur sistem, cara perancangan dan simulasi yang dilakukan .

Bab keempat berisi tentang hasil penelitian/implementasi serta pembahasan/analisis dari penelitian yang telah dilakukan dan dijelaskan secara terpadu.

Bab terakhir berisikan kesimpulan dari sistem yang telah dibuat dan saran yang akan berguna untuk pengembangan sistem dalam ma. Dengan adanya saran, diharapkan nantinya dapat membantu perusahaan untuk lebih meningkatkan kualitas sistem yang sudah ada sebelumnya.



BAB V

KESIMPULAN DAN SARAN

5.1 Kesimpulan

Berdasarkan analisis dan implementasi sistem, maka diperoleh kesimpulan sebagai berikut:

- Dengan menggunakan algoritma *Complete Linkage* sistem dapat membentuk *cluster-cluster* dengan nilai *purity* tanpa penggunaan *WordNet* lebih tinggi daripada penggunaan *WordNet* meski selisihnya tidak mencapai 0,1.
- Dengan uji coba menggunakan *WordNet* didapatkan nilai *purity* sebesar 0,8308 untuk *cutting point* 0,02. Uji coba tanpa menggunakan *WordNet* menghasilkan nilai *purity* sebesar sebesar 0,9077 untuk *cutting point* yang sama.
- Pengujian dengan *WordNet* menghasilkan *term* yang lebih sedikit bila dibandingkan tanpa penggunaan *WordNet*, hal ini dikarenakan ada beberapa *term* yang dihapus dan digantikan oleh *term* baru dari *WordNet*. Pergantian *term* ini akan meningkatkan bobot *term* sehingga meningkatkan kemungkinan lolos *feature selection*.

5.2 Saran

Saran untuk pengembangan dan perbaikan sistem adalah:

- Perlunya perbaikan dalam menggunakan *feature selection*, misalnya dengan mengganti metode *feature selection* dengan metode yang lebih baik dari *Frequency-based Feature Selection*.
- Dapat ditambahkan dengan penggunaan *stemming* dan *lemmatization* untuk *pre-processing*.
- Dapat menambahkan tabel-tabel yang digunakan dalam *WordNet*.

DAFTAR PUSTAKA

- Ackermann, M.R., Blomer, J., Kuntze, D., & Sohler, S. (2010). *Analysis of Agglomerative Clustering*. Diakses 25 Februari 2012 dari <http://drops.dagstuhl.de/opus/volltexte/2011/2994/pdf/30.pdf>.
- Agirre, E., & Lopez de Lacalle, O. *Clustering WordNet Word Senses*. (n.d.). Di akses 29 Februari 2012 dari <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.80.527&rep=rep1&type=pdf>.
- Frakes, W.B., & Baeza-Yates, R. (Eds). (1992). *Information Retrieval : Data Structure & Algorithms*. United States of America : Prentice-Hall Inc.
- Grossman, D.A., & Frieder, O. (2004). *Information Retrieval Algorithms and Heuristics 2nd Ed*. Netherlands : Springer.
- Han, J., & Kamber, M. (2006). *Data Mining Concepts & Techniques*. San Fransisco : Morgan Kaufmann.
- Hartini, E. *Metode Clustering Hirarki*. (n.d.). Diakses 1 Februari 2012 dari http://www.batan.go.id/ppin/lokakarya/LKSTN_15/Entin.pdf
- Manning, C.D., Raghavan, P., & Schütze, H. (2009). *An Introduction to Information Retrieval*. England : Cambridge University Press.
- Moens, M., (2006). *Information Extraction: Algorithms and Prospects in a Retrieval Context*. Netherlands : Springer.
- Passos, A., & Wainer, J. *WordNet-Based Metrics Do Not Seem to Help Document Clustering*. Di akses 29 Februari 2012 dari <http://www.ic.unicamp.br/~tachard/docs/wnccluster.pdf>.
- Sambamoorthi, N. (n.d.). *Hierarchical Cluster Analysis*. Di akses 4 September 2012 dari http://www.crmportals.com/hierarchical_cluster_analysis.pdf.

Soraya, Y. (2011). *Perbandingan Kinerja Metode Single Linkage, Metode Complete Linkage dan Metode K-Means dalam Analisis Cluster*. Di akses 8 Desember 2011 dari <http://lib.unnes.ac.id/6803/1/8503.pdf>.

Tan, P., Steinbach, M., & Kumar, V. (2006). *Introduction to Data Mining*. Boston : Addison-Wesley.

© UKDW