

**IMPLEMENTASI ROCCHIO'S CLASSIFICATION DALAM
MENGKATEGORIKAN RENUNGAN HARIAN KRISTEN**

Skripsi



oleh

ELISABETH ADELIA WIDJOJO

22094680

PROGRAM STUDI TEKNIK INFORMATIKA FAKULTAS TEKNOLOGI INFORMASI
UNIVERSITAS KRISTEN DUTA WACANA

2013

IMPLEMENTASI ROCCHIO'S CLASSIFICATION DALAM MENGKATEGORIKAN RENUNGAN HARIAN KRISTEN

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

ELISABETH ADELIA WIDJOJO
22094680

PROGRAM STUDI TEKNIK INFORMATIKA FAKULTAS TEKNOLOGI INFORMASI
UNIVERSITAS KRISTEN DUTA WACANA
2013

PERNYATAAN KEASLIAN SKRIPSI

Saya menyatakan dengan sesungguhnya bahwa skripsi dengan judul:

IMPLEMENTASI ROCCHIO'S CLASSIFICATION DALAM MENGKATEGORIKAN RENUNGAN HARIAN KRISTEN

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 keserjanaan 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 keserjanaan saya.

Yogyakarta, 15 Januari 2013



ELISABETH ADELIA WIDJOJO
22094680

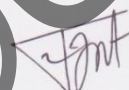
©UKDW

HALAMAN PERSETUJUAN


Judul Skripsi : IMPLEMENTASI ROCCHIO'S CLASSIFICATION
DALAM MENKATEGORIKAN RENUNGAN
HARIAN KRISTEN
Nama Mahasiswa : ELISABETH ADELIA WIDJOJO
N I M : 22094680
Matakuliah : Skripsi (Tugas Akhir)
Kode : TIW276
Semester : Gasal
Tahun Akademik : 2012/2013

Telah diperiksa dan disetujui di
Yogyakarta,
Pada tanggal 15 Januari 2013

Dosen Pembimbing I


Antonius Rachmat C., SKom.,M.Cs

Dosen Pembimbing II


Drs. R. Gunawan Santosa, M.Si.

HALAMAN PENGESAHAN

**IMPLEMENTASI ROCCHIO'S CLASSIFICATION DALAM
MENGKATEGORIKAN RENUNGAN HARIAN KRISTEN**

Oleh: ELISABETH ADELIA WIDJOJO / 22094680


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 memperoleh gelar
Sarjana Komputer
pada tanggal 10 Januari 2013

Yogyakarta, 15 Januari 2013
Mengesahkan,

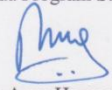
Dewan Penguji:

1. Antonius Rachmat C., SKom.,M.Cs
2. Drs. R. Gunawan Santosa, M.Si.
3. Joko Purwadi, M.Kom
4. Lukas Chrisantyo, M.Eng.

Dekan


(Drs. Wimmie Handiwidjojo, MIT.)

Ketua Program Studi


(Nugroho Agus Haryono, M.Si.)

UCAPAN TERIMA KASIH

Puji syukur penulis panjatkan ke hadirat Tuhan Yesus Kristus yang telah melimpahkan rahmat dan anugerah, sehingga penulis dapat menyelesaikan skripsi yang berjudul “Implementasi Rocchio's Classification dalam Mengkategorikan Renungan Harian Kristen” ini dengan tepat waktu.

Dalam menyelesaikan penelitian ini, penulis menyadari banyak menerima masukan dan saran dari berbagai pihak, baik secara langsung maupun secara tidak langsung. Oleh karena itu, pada kesempatan ini penulis ingin menyampaikan terima kasih kepada:

1. Antonius Rachmat C, S.Kom., M.Cs. dan Drs. R. Gunawan Santosa, M.Si., selaku dosen pembimbing yang telah banyak membimbing penulis dalam menyelesaikan penelitian ini.
2. Budi Susanto, S.Kom., M.T. selaku koordinator Tugas Akhir.
3. Keluarga terkasih yang selalu mendukung dan memberikan motivasi kepada penulis selama ini.
5. Semua teman dan pihak lain yang telah mendukung penulis dalam penyelesaian penelitian ini yang tidak dapat penulis sebutkan satu per satu.

Akhir kata penulis ingin meminta maaf bila ada kesalahan dalam penyusunan laporan ini. Terima kasih.

Yogyakarta, 15 Januari 2013

Penulis

INTISARI

IMPLEMENTASI ROCCHIO'S CLASSIFICATION DALAM MENGKATEGORIKAN RENUNGAN HARIAN KRISTEN

Pada masa sekarang, kebanyakan gereja atau lembaga-lembaga Kristen mulai menggunakan media digital untuk penyimpanan data, seperti teks, gambar, suara, maupun video rohani. Bahkan renungan harian Kristen yang biasanya berupa buku pun sekarang sudah banyak dipublikasikan dalam media digital melalui internet. Dengan semakin banyaknya renungan Kristen yang dipublikasikan di internet, akan sangat sulit untuk menemukan renungan yang sesuai dengan topik yang diinginkan karena data-data renungan yang ada masih belum dikelompokkan.

Untuk mempermudah pengelompokkan tersebut, pada penelitian ini penulis menggunakan *Rocchio's classification*, yang menekankan penggunaan *tf-idf weighting* untuk pembobotan token yang digunakan untuk proses klasifikasi, serta perhitungan *prototype vector (centroid)* setiap kategori. Setiap dokumen uji akan diuji kemiripannya dengan setiap *centroid*. Dalam penelitian ini penulis melakukan beberapa pengujian untuk mengetahui pada kondisi apa keakuratan sistem dicapai.

Sistem klasifikasi yang dibangun dapat memberikan akurasi cukup tinggi pada saat *feature selection* 20% yaitu sebesar 73,33%. Nilai *precision* tertinggi jatuh pada kategori hikmat dengan nilai *precision* 1 dalam semua *feature selection*. Sedangkan nilai *recall* tertinggi jatuh pada kategori motivator dengan nilai *recall* 1 dalam semua *feature selection*.

Kata kunci : klasifikasi, pengkategorian, renungan, Rocchio, *centroid*, *similarity*

DAFTAR ISI

HALAMAN JUDUL.....	
PERNYATAAN KEASLIAN SKRIPSI.....	iii
HALAMAN PERSETUJUAN.....	iv
HALAMAN PENGESAHAN.....	v
UCAPAN TERIMA KASIH.....	vi
INTISARI.....	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	xi
DAFTAR GAMBAR.....	xiii
DAFTAR LISTING.....	xiv
DAFTAR GRAFIK.....	xv
BAB 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 Metode Penelitian.....	3
1.6 Sistematika Penulisan.....	4
BAB 2 TINJAUAN PUSTAKA.....	5
2.1 Tinjauan Pustaka.....	5
2.2 Landasan Teori.....	6
2.2.1 Text Mining.....	6
2.2.2 Klasifikasi.....	7
2.2.2.1 Preprocessing Data.....	7
2.2.2.1.1 Tokenisasi.....	8
2.2.2.1.2 Penghapusan Stopword.....	8
2.2.2.2 Pembobotan TF/IDF.....	8

2.2.2.3 <i>Frequency-based Feature Selection</i>	10
2.2.2.4 Algoritma Rocchio untuk Klasifikasi Teks.....	10
2.2.3 Evaluasi Sistem.....	12
BAB 3 ANALISIS DAN PERANCANGAN SISTEM	14
3.1 Spesifikasi Kebutuhan Perangkat Keras dan Perangkat Lunak.....	14
3.1.1 Spesifikasi Kebutuhan Perangkat Keras.....	14
3.1.2 Spesifikasi Kebutuhan Perangkat Lunak.....	14
3.2 Spesifikasi Sistem.....	15
3.3 Diagram Use Case.....	16
3.4 Flowchart.....	17
3.5 Kamus Data.....	24
3.6 Skema Diagram Basis Data.....	27
3.7 Rancangan Antarmuka.....	30
3.8 Rancangan Pengujian dan Evaluasi Sistem.....	33
3.8.1 Rancangan Pengujian.....	33
3.8.2 Rancangan Evaluasi Sistem.....	33
3.9 Contoh Kasus Klasifikasi.....	35
BAB 4 IMPLEMENTASI DAN ANALISIS SISTEM	40
4.1 Implementasi Sistem.....	40
4.1.1 Antarmuka Sistem.....	40
4.1.2 Pengumpulan Dokumen.....	42
4.1.3 Pseudocode Sistem.....	43
4.1.3.1 Tahap Preprocessing.....	43
4.1.3.2 Tahap Perhitungan Centroid.....	45
4.1.3.3 Tahap Klasifikasi.....	46
4.2 Evaluasi Sistem.....	47
4.2.1 Feature Selection 10%.....	49
4.2.2 Feature Selection 20%.....	53
4.2.3 Feature Selection 30%.....	58
4.2.4 Feature Selection 40%.....	62
4.2.5 Feature Selection 50%.....	67

4.2.6 Evaluasi Precision Recall Menurut Sumber Data.....	71
4.2.7. Grafik Hasil Pengujian.....	75
4.2.7.1 Grafik Keakuratan Sistem.....	75
4.2.7.2 Grafik Precision Recall Kategori Berkat.....	75
4.2.7.3 Grafik Precision Recall Kategori Motivator.....	76
4.2.7.4 Grafik Precision Recall Kategori Iman.....	76
4.2.7.5 Grafik Precision Recall Kategori Hikmat.....	77
4.2.7.6 Grafik Precision Recall <i>Feature Selection</i> 10%.....	77
4.2.7.7 Grafik Precision Recall <i>Feature Selection</i> 20%.....	78
4.2.7.8 Grafik Precision Recall <i>Feature Selection</i> 30%.....	79
4.2.7.9 Grafik Precision Recall <i>Feature Selection</i> 40%.....	80
4.2.7.10 Grafik Precision Recall <i>Feature Selection</i> 50%.....	80
4.2.7.11 Grafik Precision Recall Sumber Data AH.....	81
4.2.7.12 Grafik Precision Recall Sumber Data BT.....	82
4.2.7.13 Grafik Precision Recall Sumber Data SP.....	83
4.2.7.14 Grafik Precision Recall Sumber Data OL.....	83
4.2.7.15 Grafik Precision Recall Sumber Data GK.....	84
BAB 5 KESIMPULAN DAN SARAN.....	85
5.1 Kesimpulan.....	85
5.2 Saran.....	85
DAFTAR PUSTAKA.....	86
LAMPIRAN	

DAFTAR TABEL

Tabel 2.1	Frekuensi Kemunculan Token.....	9
Tabel 3.1	Struktur Tabel Kategori.....	24
Tabel 3.2	Struktur Tabel Dokumen.....	24
Tabel 3.3	Struktur Tabel Token.....	25
Tabel 3.4	Struktur Tabel Stopword.....	25
Tabel 3.5	Struktur Tabel Dokumen Uji.....	25
Tabel 3.6	Struktur Tabel Token Uji.....	26
Tabel 3.7	Struktur Tabel Dokumen Token.....	26
Tabel 3.8	Struktur View_1.....	27
Tabel 3.9	Struktur View_2.....	27
Tabel 3.10	Struktur View_3.....	28
Tabel 3.11	Struktur View_4.....	28
Tabel 3.12	TF/IDF Dokumen Pelatihan.....	36
Tabel 3.13	TF/IDF Dokumen Uji.....	36
Tabel 3.14	Sorting Bobot Dokumen Pelatihan Kategori Iman.....	37
Tabel 3.15	Sorting Bobot Dokumen Pelatihan Kategori Motivator.....	37
Tabel 4.1	Hasil Pengujian dengan FS 10%.....	49
Tabel 4.2	Confusion Matrix dengan FS 10%.....	50
Tabel 4.3	Hasil Pengujian dengan FS 20%.....	54
Tabel 4.4	Confusion Matrix dengan FS 20%.....	54
Tabel 4.5	Hasil Pengujian dengan FS 30%.....	58
Tabel 4.6	Confusion Matrix dengan FS 30%.....	59
Tabel 4.7	Hasil Pengujian dengan FS 40%.....	63
Tabel 4.8	Confusion Matrix dengan FS 40%.....	63
Tabel 4.7	Hasil Pengujian dengan FS 50%.....	67
Tabel 4.10	Confusion Matrix dengan FS 50%.....	68
Tabel 4.11	Inisial Sumber Data.....	71
Tabel 4.12	Sumber Data Dokumen Uji.....	72
Tabel 4.13	Confusion matrix Sumber Data AH.....	72

Tabel 4.14	Confusion matrix Sumber Data BT.....	73
Tabel 4.15	Confusion matrix Sumber Data SP.....	73
Tabel 4.16	Confusion matrix Sumber Data OL.....	74
Tabel 4.17	Confusion matrix Sumber Data GK.....	74

©UKDW

DAFTAR GAMBAR

Gambar 2.1	Contoh Perhitungan TF/IDF.....	10
Gambar 2.2	Contoh Confusion Matrix.....	13
Gambar 3.1	Use Case Diagram.....	16
Gambar 3.2	Flowchart Sistem.....	17
Gambar 3.3	Flowchart Preprocessing Dokumen.....	18
Gambar 3.4	Flowchart Tokenisasi.....	19
Gambar 3.5	Flowchart Tambah Dokumen Pelatihan.....	20
Gambar 3.6	Flowchart Pelatihan Dokumen.....	21
Gambar 3.7	Flowchart Uji Dokumen dan Pemilihan Kategori.....	22
Gambar 3.8	Flowchart Menghitung Cosine Similarity.....	23
Gambar 3.9	Skema Diagram Basis Data.....	29
Gambar 3.10	Form Login.....	30
Gambar 3.11	Form Lihat.....	30
Gambar 3.12	Form Tambah.....	31
Gambar 3.13	Form Klasifikasi.....	32
Gambar 4.1	Form Login.....	40
Gambar 4.2	Menu Lihat Renungan (admin dan user).....	41
Gambar 4.3	Menu Tambah Renungan (admin).....	41
Gambar 4.4	Menu Klasifikasi (admin dan user).....	42
Gambar 4.5	Contoh pengujian dengan FS 10% (benar).....	48
Gambar 4.6	Contoh pengujian dengan FS 10% (salah).....	48

DAFTAR LISTING

Listing 4.1	Pseudocode Tokenisasi dan Penghapusan Stopword.....	43
Listing 4.2	Pseudocode Perhitungan TF/IDF.....	44
Listing 4.3	Pseudocode Perhitungan Centroid dan Pemilihan FS.....	45
Listing 4.4	Pseudocode Proses Klasifikasi.....	46
Listing 4.5	Pseudocode Pemilihan Kategori.....	47

©UKDW

DAFTAR GRAFIK

Grafik 4.1	Grafik Keakuratan Sistem.....	75
Grafik 4.2	Grafik Precision Recall Kategori Berkat.....	75
Grafik 4.3	Grafik Precision Recall Kategori Motivator.....	76
Grafik 4.4	Grafik Precision Recall Kategori Iman.....	76
Grafik 4.5	Grafik Precision Recall Kategori Hikmat.....	77
Grafik 4.6	Grafik Precision Recall Feature Selection 10%.....	77
Grafik 4.7	Grafik Precision Recall Feature Selection 20%.....	78
Grafik 4.8	Grafik Precision Recall Feature Selection 30%.....	79
Grafik 4.9	Grafik Precision Recall Feature Selection 40%.....	80
Grafik 4.10	Grafik Precision Recall Feature Selection 50%.....	80
Grafik 4.11	Grafik Precision Recall Sumber Data AH.....	81
Grafik 4.12	Grafik Precision Recall Sumber Data BT.....	82
Grafik 4.13	Grafik Precision Recall Sumber Data SP.....	83
Grafik 4.14	Grafik Precision Recall Sumber Data OL.....	83
Grafik 4.15	Grafik Precision Recall Sumber Data GK.....	84

BAB 1

PENDAHULUAN

1.1 Latar Belakang Masalah

Pada masa sekarang, kebanyakan gereja atau lembaga-lembaga Kristen mulai menggunakan media digital untuk penyimpanan data, seperti teks, gambar, suara, maupun video rohani. Bahkan renungan harian Kristen yang biasanya berupa buku pun sekarang sudah banyak dipublikasikan dalam media digital melalui internet. Orang seringkali mencari renungan dengan topik tertentu sebagai bahan khotbah atau mungkin sekedar menguatkan hati dalam menghadapi pergumulan hidup. Dengan semakin banyaknya renungan Kristen yang dipublikasikan di internet, akan sangat sulit untuk menemukan renungan yang sesuai dengan topik yang diinginkan karena data-data renungan yang ada masih belum dikelompokkan.

Berbicara mengenai pengelompokan data, dalam hal ini berupa dokumen teks, ada banyak metode klasifikasi yang bisa diterapkan antara lain Naïve Bayes, k-Nearest Neighbor, Decision Tree, Support Vector Machines, dan lain-lain. Salah satu teknik yang umum digunakan dalam pemrosesan dokumen teks adalah dengan pembobotan kata-kata yang dianggap penting dalam dokumen, yaitu dengan menghitung *term-frequency* (tf). Akan tetapi, jumlah dokumen di mana sebuah token/kata muncul (df) juga harus dipertimbangkan untuk melihat seberapa penting suatu token dalam dokumen tersebut (*scarcity of tokens*). Dalam hal ini akan dihitung *inverse document frequency* untuk tiap token (idf). Dari tf dan idf, kita akan mendapatkan *tf-idf weighting* (Robertson, 2005), yaitu suatu formula untuk menghitung bobot hubungan suatu token di dalam dokumen.

Oleh karena itu, pada penelitian ini penulis menggunakan *Rocchio's classification*, yang menekankan penggunaan *tf-idf weighting* untuk pembobotan token yang digunakan untuk klasifikasi dokumen. Klasifikasi dilakukan sesuai dengan algoritma Rocchio yang menggunakan *prototype vector (centroid)* untuk

membantu klasifikasi/pengelompokkan renungan harian. Dengan adanya penelitian ini, maka dapat membantu siapa saja yang ingin mencari renungan harian sesuai dengan topik yang diinginkan.

1.2 Perumusan Masalah

Rumusan masalah yang dapat dibuat dalam tugas akhir ini adalah:

- Seberapa akuratkah *Rocchio's classification* dalam mengkategorikan renungan harian Kristen?
- Kategori manakah yang memiliki nilai precision dan recall yang paling tinggi?
- Apakah pemilihan *feature selection* mempengaruhi nilai precision dan recall pada masing-masing kategori?
- Sumber data manakah yang memberikan nilai precision dan recall yang paling tinggi?

1.3 Batasan Masalah

Batasan-batasan masalah dalam pembuatan aplikasi ini adalah:

- Sistem yang dibangun adalah berupa aplikasi desktop
- Sumber data dari renungan harian Kristen diambil dari :
 - <http://www.pelitahidup.com/>
 - <http://www.renungan-spirit.com/renungan-kristen.html>
 - <http://gkysydney.org/renungan-gema-2012/index.php>
 - <http://www.sabdaharian.com>
 - <http://www.bethanygraha.org>
 - <http://renungan-harian-online.com>
 - Renungan Harian Air Hidup
 - <http://renungan-harian-online.com>
 - <http://www.gkpi.or.id/renungan/>

- Renungan harian Kristen akan dibagi dalam 4 kategori yaitu : berkat, iman, hikmat, dan motivator.
- Tidak ada tahap *stemming* dalam tahap *preprocessing* dokumen.
- *Stoplist* yang digunakan adalah data *stopwords* yang bersumber dari <http://www.ilc.uva.nl/Research/Reports/MoL-2003-02.text.pdf>, kata kunci *database* dan tanda baca.
- Dokumen yang digunakan berekstensi .txt dimana :
 - nama dokumen mewakili kode kategori dan judul dari renungan harian Kristen
 - isi dokumen merupakan isi dari satu renungan harian Kristen
- Menggunakan pembobotan TF/IDF untuk merepresentasikan dokumen teks

1.4 Tujuan Penelitian

Tujuan dari penelitian ini adalah melihat keakuratan *Rocchio's classification* dalam mengkategorikan renungan harian Kristen, khususnya dalam topik berkat, iman, hikmat, dan motivator.

1.5 Metode Penelitian

Metodologi yang akan digunakan dalam pembuatan aplikasi ini adalah dengan studi pustaka dan pengumpulan data. Studi pustaka dilakukan dengan mempelajari teori-teori melalui buku, artikel, jurnal, dan bahan lain yang mendukung yang berhubungan dengan klasifikasi dokumen. Pengumpulan data dilakukan dengan mencari referensi resmi (sumber data) untuk setiap renungan harian yang akan digunakan.

1.6 Sistematika Penulisan

Penulisan laporan tugas akhir ini dibagi menjadi lima bab yaitu:

Bab 1, Pendahuluan, yang memberikan gambaran umum mengenai apa yang diteliti dalam tugas akhir ini. Pendahuluan berisi latar belakang masalah, perumusan masalah, batasan masalah, tujuan penelitian, metode penelitian, dan sistematika penulisan.

Bab 2, Tinjauan Pustaka, yang terdiri dari dua bagian utama, yakni tinjauan pustaka dan landasan teori. Tinjauan pustaka menguraikan berbagai teori mengenai klasifikasi dengan metode *Rocchio's classification* yang didapatkan dari berbagai sumber pustaka yang digunakan untuk penyusunan tugas akhir. Landasan teori memuat penjelasan tentang konsep dan prinsip utama yang diperlukan untuk memecahkan masalah dalam penelitian. Hanya penjelasan yang berhubungan dengan penelitian yang dilakukan yang akan dicantumkan di sini.

Bab 3, Analisis dan Perancangan Sistem, yang mencakup perancangan sistem yang akan dibuat, yakni mengenai kebutuhan *hardware* dan *software*, spesifikasi sistem, arsitektur sistem, *use case diagram*, algoritma yang digunakan, skema *database*, dan rancangan antarmuka sistem.

Bab 4, Implementasi dan Analisis Sistem, yang memuat hasil implementasi dan pembahasan mengenai pengujian sistem yang dibuat berdasarkan bab 3, beserta hasil dari sistem yang dijalankan dan analisis dari sistem yang dibuat.

Bab 5, Kesimpulan dan Saran, berisi kesimpulan dari hasil penelitian yang didapatkan dan saran untuk memberikan analisis dan pengembangan yang lebih baik lagi pada penelitian ke depannya dalam topik yang serupa.

BAB 5

KESIMPULAN DAN SARAN

5.1 Kesimpulan

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

1. Sistem klasifikasi Rocchio memberikan akurasi cukup tinggi untuk *feature selection* 20% yaitu sebesar 73,33%.
2. Nilai precision tertinggi jatuh pada kategori hikmat dengan nilai precision 1 dalam semua *feature selection*. Sedangkan nilai recall tertinggi jatuh pada kategori motivator dengan nilai recall 1 dalam semua *feature selection*.
3. Peningkatan *feature selection* tidak terlalu mempengaruhi nilai precision dan recall pada setiap kategori.
4. Sumber data Renungan Harian Spirit sangat cocok untuk kategori motivator karena memiliki nilai precision dan recall 1 dari total 13 renungan harian yang diambil.

5.2 Saran

Saran yang diajukan penulis untuk perbaikan dan pengembangan sistem adalah sebagai berikut :

1. Diperlukan penggunaan *store procedure* pada VB.NET untuk mempercepat *preprocessing* data dan mengurangi penggunaan memori.
2. Dapat ditambahkan proses *stemming* dalam bahasa Indonesia untuk lebih meningkatkan akurasi sistem.

DAFTAR PUSTAKA

- Auvil, Loretta, et all. (2007). VAST to Knowledge : Combining Tools for Exploration and Mining. *IEEE VAST*, 197-198. Diakses pada tanggal 1 September 2012 dari <http://webcache.googleusercontent.com/search?q=cache:http://vac.nist.gov/2007/summaries/auvil.pdf>
- Kurniawan, Erick. (2011). *Cepat Mahir Visual Basic 2010*. Yogyakarta : Penerbit Andi.
- Koster, Cornelis H. A. (2003). *Chapter 2 : Document Classification*. Diakses pada tanggal 17 Mei 2012 dari <http://www.cs.ru.nl/~kees/ir2/papers/h03.pdf>
- Feldman, R., dan Sanger, J. (2007). *The Text Mining Handbook: Advanced Approaches in Analyzing Unstructured Data*. Cambridge : Cambridge University Press.
- Grossman, David A. dan Frieder, Ophir. (2004). *Information Retrieval Algorithms dan Heuristics, 2nd edition*. New York : Springer.
- Han, J. & Kamber, M. (2006). *Data Mining : Concepts and Technique 2nd Edition*. San Fransisco : Morgan Kauffman Publishers.
- Intan, Rolly & Defeng, Andrew. (2006). *HARD : Subject Based Search Engine Menggunakan TF-IDF dan Jaccard's Coeffisient*. Diakses pada tanggal 25 Agustus 2012 dari <http://puslit.petra.ac.id/files/published/journals/IND/IND060801/IND06080106.pdf>
- Joachims, T. (1997). *A Probabilistic Analysis of the Rocchio Algorithm with TF/IDF for Text Categorization*. Diakses pada tanggal 20 Oktober 2012 dari http://www.cs.cornell.edu/People/tj/publications/joachims_97a.ps.gz
- Manning, Christopher D., et all. (2008). *Introduction to Information Retrieval*. New York : Cambridge University Press.
- Salton, G, et all. (1975). A Vector Space Model for Automatic Indexing. *Communications of the ACM*, 12, 613-620. Diakses pada tanggal 1 September 2012 dari http://www.cs.uiuc.edu/class/fa05/cs511/Spring05/other_papers/p613-salton.pdf

Uden, Mark Van. (n.d). *Rocchio : Relevance Feedback in Learning Classification Algorithms*. Diakses pada tanggal 23 Agustus 2012 dari <http://www.cs.kun.nl/nscs/artikelen/markuden.ps.Z>.

Weiss, Sholom M., et all. (2005). *Text Mining : Predictive Methods for Analyzing Unstructured Information*. New York : Springer

©UKDW