

**PENERAPAN KLASIFIKASI *TWEETS* PADA BERITA TWITTER
MENGUNAKAN METODE *K-NEAREST NEIGHBOR* DAN
QUERY EXPANSION BERBASIS *DISTRIBUTIONAL SEMANTIC***

SKRIPSI

Untuk memenuhi sebagian persyaratan memperoleh
gelar Sarjana Komputer

Disusun oleh:
Galih Nuring Bagaskoro
NIM: 135150218113027



PROGRAM STUDI TEKNIK INFORMATIKA
JURUSAN TEKNIK INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS BRAWIJAYA
MALANG
2018

PENGESAHAN

Penerapan Klasifikasi *Tweets* Pada Berita Twitter Menggunakan metode *K-Nearest Neighbor* dan *Query Expansion* berbasis *Distributional Semantic*

SKRIPSI

Diajukan untuk memenuhi sebagian persyaratan memperoleh gelar Sarjana
Komputer

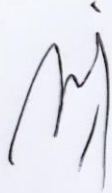
Disusun Oleh :

Galih Nuring Bagaskoro

NIM: 135150218113027

Skripsi ini telah diuji dan dinyatakan lulus pada
12 Januari 2018 Telah diperiksa dan disetujui
oleh:

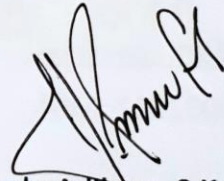
Dosen Pembimbing I



M. Ali Fauzi S.Kom, M.Kom

NIK. 201502 890101 1 001

Dosen Pembimbing II



Putra Pandu Adikara, S.Kom, M.Kom

NIP. 19850725 200812 1 002

Mengetahui

Ketua Jurusan Teknik Informatika



Tri Astoto Kurniawan, S.T, M.T, Ph.D

NIP. 19710518 200312 1 001

IDENTITAS TIM PENGUJI

PEMBIMBING 1

Mochammad Ali Fauzi, S.Kom, M.Kom

NIK. 201502 890101 1 001

Email: moch.ali.fauzi@ub.ac.id

PEMBIMBING 2

Putra Pandu Adikara, S.Kom, M.Kom

NIP. 19850725 200812 1 002

Email: adikara.putra@ub.ac.id

PENGUJI 1

Rizal Setya Perdana, S.Kom, M.Kom

NIP. 2016039101181001

Email: rizalespe@ub.ac.id

PENGUJI 2

Lailil Muflikhah, S.kom, M.Sc

NIP. 197411132005012001

Email: lailil@ub.ac.id

PERNYATAAN ORISINALITAS

Saya menyatakan dengan sebenar-benarnya bahwa sepanjang pengetahuan saya, di dalam naskah skripsi ini tidak terdapat karya ilmiah yang pernah diajukan oleh orang lain untuk memperoleh gelar akademik di suatu perguruan tinggi, dan tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain, kecuali yang secara tertulis disitasi dalam naskah ini dan disebutkan dalam daftar pustaka.

Apabila ternyata didalam naskah skripsi ini dapat dibuktikan terdapat unsurunsur plagiasi, saya bersedia skripsi ini digugurkan dan gelar akademik yang telah saya peroleh (sarjana) dibatalkan, serta diproses sesuai dengan peraturan perundang-undangan yang berlaku (UU No. 20 Tahun 2003, Pasal 25 ayat 2 dan Pasal 70).

Malang, 12 Januari 2018



Janji Nuring Bagaskoro

NIM: 135150218113027

RIWAYAT HIDUP



Penulis dilahirkan di Tulungagung pada tanggal 11 Juni 1995 dari Ibu yang bernama Sri Wulan Ros Indah dan Ayah bernama Mudjiono Emje. Penulis menyelesaikan pendidikan sekolah dasar di SDN KAMPUNG DALEM VI pada tahun 2007 dengan mengikuti organisasi pramuka dan berbagai ekstrakurikuler yang diikuti. Kemudian melanjutkan ke sekolah menengah pertama di SMP Islam Al-huda Kediri dengan tahun kelulusan 2010 dengan mengikuti organisasi OSIS. Dan melanjutkan studi ke jenjang sekolah menengah atas di SMAN 8 Kota Kediri dengan menyelesaikan masa studi pada tahun 2013

Pada tahun 2018 penulis telah berhasil menyelesaikan pendidikannya di Universitas Brawijaya pada Program Studi Teknik Informatika Jurusan Teknik Informatika Fakultas Ilmu Komputer. Selain itu penulis aktif menjadi Anggota Pengusaha Muda Kediri, serta menjadi aktivis sosial yang bergerak dengan Seni.

“Impossible Is Nothing”

Assalamualaikum wr.wb.

*Alhamdulillah, Terimakasih Ya Rabb atas Kehadirat Allah SWT,
serta junjungan besar Nabi Muhammad SAW.*

*Tugas Akhir ini saya persembahkan untuk Orang Tua dan
Keluarga Besar saya yang selalu mendukung saya serta kepada
seluruh teman-*

*temanku dan orang terdekat yang tak berhenti memberikan
semangat dan motivasi.*

ABSTRAK

Galih Nuring Bagaskoro, 2018. Penerapan Klasifikasi *Tweets* pada Berita Twitter Menggunakan Metode *K-Nearest Neighbor* dan *Query Expansion* Berbasis *Distributional Semantic*.

Fakultas Ilmu Komputer Universitas Brawijaya. Pembimbing : Ali Fauzi S.Kom, M.Sc dan Putra Pandu Adikara S.Kom, M.Kom

Penggunaan teks pendek berbasis digital sampai saat ini masih berkembang dan meluas hingga diberbagai media sosial. Media sosial Twitter memiliki fitur kategori jenis informasi melalui *tweets* yang di unggah. Setiap pengelompokan jenis informasi dilakukan agar mempermudah pengguna untuk memanfaatkannya. Tujuan dari penggunaan kategori dalam hal ini klasifikasi, untuk meningkatkan kualitas media sosial dalam pengelompokan kategori isi dari konten yang disediakan. Klasifikasi tradisional sampai saat ini masih digunakan, namun hasil yang diperoleh terkadang tidak maksimal, perlu dilakukan ekspansi kata untuk menambahkan kata kedalam teks agar dapat meningkatkan akurasi. Ekspansi kata digunakan dengan berbasis *distributional semantic* dengan teknik *Euclidean distance* untuk menemukan kata terdekat dari sumber eksternal agar menjadi kueri yang akan ditambahkan ke teks data uji. Dengan menggunakan data uji 105 dan data latih 400, klasifikasi yang menggunakan *K-Nearest Neighbor* dapat memperoleh hasil 90% dengan tetangga terdekat $K=5$. Hasil tersebut sama halnya dengan hasil pengujian yang dilakukan dengan tanpa menggunakan teknik ekspansi kata. Sedangkan pengujian yang dilakukan dengan menambahkan ekspansi kata dengan *threshold* 0,5 dan nilai tertangga terdekat *K-Nearest Neighbor* $K=5$ memperoleh hasil akurasi 92%.

Kata kunci: *twitter, tweet, ekspansi kata, distributional semantic, euclidean distance, klasifikasi, k-nearest neighbor*

ABSTRACT

Galih Nuring Bagaskoro, 2018. Application of Tweets Classification on Twitter News Using K-Nearest Neighbor Method and Query Expansion Based on Distributional Semantic.

Faculty of Computer Science University of Brawijaya. Advisor : Ali Fauzi, S.Kom. M.Sc and Putra Pandu Adikara, S.Kom. M.Kom

The use of short text based on digital to date is still growing and extending to various social media. Twitter has news features in tweets to represent information representing each type. Each categorization of this type is done to make it easier for users to use it. The purpose of the use of categories in this classification, to evaluate and improve the quality of social media in grouping categories of content of the content provided. Traditional classification is still used today, but the results are sometimes not maximal, it is necessary to expand the word to add words to the text in order to improve the accuracy. Word expansion is used with a semantic-based distributional euclidean distance technique to find the closest word from an external source to be a query to be added to the test data text. Using test data 105 and training data 400, the classification using K-Nearest Neighbor can obtain 90% results with nearest neighbor K=5. These results are similar to the results of tests conducted without using word expansion techniques. While the test is done by adding the expansion of words with threshold 0.5 and the nearest immediate value K-Nearest Neighbor K=5 obtained an accuracy of 92%.

Keywords: *twitter, tweet, word expansion, distributional semantic, euclidean distance, classification, k-nearest neighbor*

KATA PENGANTAR

Puji syukur penulis panjatkan kehadirat Allah SWT atas limpahan rahmat dan karunia-Nya sehingga penulis dapat menyelesaikan skripsi yang berjudul “Penerapan Klasifikasi Tweets pada Berita Twitter Menggunakan metode K-Nearest Neighbor dan Query Expansion Berbasis Distributional Semantic”. Skripsi ini diajukan untuk memenuhi prasyarat kelulusan dalam menyelesaikan pendidikan dan memperoleh gelar Sarjana Komputer konsentrasi Komputasi Cerdas dan Visualisasi di program studi Informatika Fakultas Ilmu Komputer Universitas Brawijaya.

Dalam penyelesaian skripsi ini, penulis telah mendapatkan banyak dukungan dan bantuan dari berbagai banyak kalangan pihak. Mungkin hanya ucapan terima kasih dari penulis yang bias membalas bantuan dan dukungan dari beberapa pihak. Atas bantuan yang telah diberikan, penulis ingin menyampaikan sebuah ucapan terima kasih yang sebesar-besarnya kepada :

1. Bapak Ali Fauzi, S.Kom., M.Kom., selaku dosen pembimbing utama yang telah meluangkan waktu untuk memberikan pengarahan dan masukan kepada penulis
2. Bapak Putra Pandu Adikara, S.Kom, M.Kom selaku dosen pendamping kedua yang juga telah meluangkan waktu untuk memeberikan pengarahan dan masukan pada penulis
3. Wayan Firdaus Mahmudy, S.Si., M.T., Ph.D., Ir. Heru Nurwasito, M.Kom., Drs. Mardji M.T, dan Edy Santoso, S.Si., M.Kom., selaku Dekan, wakil Dekan 1, Wakil Dekan 2, Wakil Dekan 3 Fakultas Ilmu Komputer Universitas Brawijaya
4. Tri Astoto Kurniawan, S.T, M.T, Ph.D dan M. Tanzil Furqon, S.Kom, M.CompSc., selaku Ketua dan Sekretaris Program Studi Teknik Informatika Universitas Brawijaya
5. Seluruh Dosen Fakultas Ilmu Komputer Universitas Brawijaya atas kesediaannya membagi ilmu kepada penulis
6. Bapak dan Ibu kandung serta Adik-Adik yang senantiasa selalu memeberikan dukungan dan doa kepada penulis demi kelancaran pengerjaan skripsi.
7. Seorang bernama Alfi Rizky Anita Fajarina yang tidak luput dari bantuan masukan dan inspirasi, serta dukungan terbaiknya yang membuat penulis terus bersemangat.
8. Teman-teman dan para sahabat penulis Wildan Afif, Nasrul Ashar, Cristian Herlando, Dhimas Tungga, Afif Ridwan, Erda Endika, dan Ilham

Zunaidy yang memberikan bantuan, semangat dan doa demi terselesaikannya skripsi ini.

9. Semua pihak yang tidak dapat disebutkan satu persatu oleh penulis yang memberikan dukungan dan semangat secara langsung maupun tidak langsung demi terselesainya skripsi ini.

Meskipun penulis berharap isi dari skripsi ini bebas dari kekurangan dan kesalahan, namun karena ini adalah buatan dari manusia yang dengan kata lain pasti tidaklah sempurna. Oleh karena itu, penulis mengharapkan sebuah kritik dan saran yang dapat membangun agar skripsi ini dapat lebih baik. Akhir kata penulis berharap skripsi ini dapat bermanfaat bagi segala kalangan pembaca.

Malang, 8 Januari 2018

penulis

galihnuring2@gmail.com

DAFTAR ISI

PENGESAHAN.....	ii
PERNYATAAN ORISINALITAS.....	iii
KATA PENGANTAR	iv
ABSTRAK	vi
ABSTRACT	vii
DAFTAR ISI.....	viii
DAFTAR TABEL	xi
DAFTAR GAMBAR.....	xiii
DAFTAR KODE PROGRAM	xiv
DAFTAR PERSAMAAN	xv
DAFTAR LAMPIRAN.....	xvi
BAB 1 PENDAHULUAN	1
1.1 Latar belakang.....	1
1.2 Rumusan masalah.....	3
1.3 Tujuan	3
1.4 Manfaat.....	3
1.5 Batasan masalah	4
1.6 Sistematika pembahasan	4
BAB 2 LANDASAN KEPUSTAKAAN	6
2.1 Kajian Pustaka	6
2.2 Twitter	7
2.2.1 Halaman Utama (<i>Home</i>)	7
2.2.2 Profil (<i>Profile</i>).....	7
2.2.3 Favorit (<i>Favorite</i>).....	7
2.2.4 <i>Mention</i>	7
2.2.6 Tagar (<i>Hashtag/#</i>).....	7
2.2.7 Topik populer(<i>Trending Topic</i>).....	7

2.2.8 <i>List</i>	7
2.2.9 Mengikuti (<i>Following</i>)	8
2.2.10 Pengikut (<i>Follower</i>)	8
2.3 Text Mining.....	8
2.3.1 <i>Text Preprocessing</i>	8
2.3.1.1 <i>Case Folding</i>	8
2.3.1.2 <i>Tokenizing</i>	9
2.3.1.3 <i>Filtering</i>	9
2.3.1.4 <i>Stemming</i>	9
2.4 <i>Query Expansion</i>	10
2.5 <i>Distributional Semantic</i>	10
2.5.1 <i>Euclidean Distance</i>	10
2.6 Klasifikasi Teks	11
2.6.1 <i>K-Nearest Neighbor</i>	11
2.6.1.1 <i>Cosine Similarity</i>	12
2.7 Evaluasi.....	13
BAB 3 METODOLOGI	14
3.1 Studi Kepustakaan	14
3.2 Pengumpulan data	15
3.3 Implementasi Perangkat Lunak.....	15
3.4 Pengujian Perangkat Lunak	15
3.5 Pengambilan kesimpulan dan Saran	16
BAB 4 ANALISIS DAN PERANCANGAN	17
4.1 Lingkungan Perancangan dan Implementasi	17
4.2 Perancangan Perangkat Lunak.....	17
4.2.1 Deskripsi Sistem Secara Umum.....	17
4.2.2 Perancangan Diagram alir sistem.....	19
4.3 Manualisasi Perhitungan Data	25
4.3.1 Manualisasi Klasifikasi Tanpa Ekspansi Kata	25
4.3.2 Manualisasi Menggunakan <i>Query Expansion</i>	32
4.3.3 Manualisasi Klasifikasi Menggunakan <i>Query Expansion</i>	36
4.4 Perancangan Antarmuka	43

4.4.1 Halaman Dokumen Kamus dan <i>Term Unik</i>	44
4.4.2 Halaman Klasifikasi Teks	44
4.4.3 Halaman Pengujian Klasifikasi	45
4.5 Perancangan Pengujian	45
4.5.1 Pengujian nilai <i>k</i> pada <i>KNN</i>	46
4.5.2 Perancangan Pengujian Variasi <i>Threshold</i> pada <i>Query Expansion</i>	46
4.6 Penarikan Kesimpulan	47
BAB 5 IMPLEMENTASI.....	48
5.1 Batasan Implementasi	48
5.2 Implementasi Algoritme	48
5.2.1 Implementasi Algoritme <i>Preprocessing Text</i>	48
5.2.2 <i>Query Expansion</i>	51
5.2.3 Klasifikasi Teks	53
5.3 Implementasi Antarmuka	57
5.3.1 Implementasi Antarmuka Dokumen Berita.....	58
5.3.2 Implementasi Antarmuka Klasifikasi Teks	58
5.3.3 Implementasi Antarmuka Hasil Klasifikasi.....	59
BAB 6 PENGUJIAN DAN ANALISIS	60
6.1 Data yang digunakan	60
6.2 Pengujian Nilai <i>k</i> Pada <i>K-Nearest Neighbor</i>	60
6.2.1. Skenario Pengujian Nilai <i>k</i> Pada <i>K-Nearest Neighbor</i>	60
6.2.2. Analisis Pengujian Nilai <i>k</i> Pada <i>K-Nearest Neighbor</i>	61
6.2 Pengujian Variasi <i>Threshold</i> nilai kedekatan pada <i>Query Expansion</i>	62
6.2.1 Skenario Pengujian Variasi <i>Threshold</i> pada <i>Query Expansion</i>	63
6.2.2 Analisis Pengujian Variasi <i>Threshold</i> pada <i>Query Expansion</i>	64
BAB 7 KESIMPULAN DAN SARAN	66
7.1 Kesimpulan.....	66
7.2 Saran	66
Daftar Pustaka.....	68
LAMPIRAN DATASET	70

DAFTAR TABEL

Tabel 4.1 Data Latih 1.....	24
Tabel 4.2 Data Uji Klasifikasi 1.....	24
Tabel 4.3 Perhitungan <i>Tf-Idf</i>	25
Tabel 4.4 Nilai Yang Digunakan Untuk menghitung Jarak Antar Dokumen.....	27
Tabel 4.5 Normalisasi akhir 1.....	28
Tabel 4.6 Contoh Dokumen Kamus Berita.....	31
Tabel 4.7 id Kata Ekspansi.....	31
Tabel 4.8 Jumlah Kata pada Setiap Dokumen	32
Tabel 4.9 Perhitungan Kemunculan Antar Kata.....	33
Tabel 4.10 Hasil Ekspansi Kata.....	34
Tabel 4.11 Data Latih 2.....	35
Tabel 4.12 Daftar Ekspansi Kata.....	35
Tabel 4.13 Data Uji Klasifikasi 2.....	35
Tabel 4.14 Pembobotan <i>TF-IDf</i> 2.....	36
Tabel 4.15 Tabel Untuk Perhitungan Normalisasi 2.....	38
Tabel 4.16 Hasil Perhitungan Normalisasi 2.....	39
Tabel 4.17 Perancangan Pengujian Nilai <i>k</i> Pada <i>KNN</i>	44
Tabel 4.17 Perancangan Pengujian variasi <i>Threshold</i> Pada <i>Query Expansion</i>	45
Tabel 6.1 Hasil Pengujian <i>K-Nearest Neighbor</i>	57
Tabel 6.2 Contoh Hasil Penentuan Dari Pemilihan Nilai <i>k</i> <i>tweet</i> ekonomi.....	58
Tabel 6.3 Contoh Hasil Penentuan Dari Pemilihan Nilai <i>k</i> Pada <i>Tweet</i> olahraga.....	58
Tabel 6.4 Hasil Pengujian <i>Query Expansion</i>	60
Tabel 6.5 Contoh Hasil <i>Tweets</i> Yang Diekspansi.....	61
Tabel 6.5 Contoh ekspansi kata menggunakan <i>Threshold</i>	61

DAFTAR GAMBAR

Gambar 2.1 Contoh Proses <i>Tokenizing</i>	8
Gambar 2.2 Contoh Proses <i>Filtering</i>	9
Gambar 3.1 Metodologi Penelitian.....	13
Gambar 4.1 Diagram Alur Kerja Sistem Secara Umum.....	17
Gambar 4.2 Diagram Alir Sistem.....	18
Gambar 4.3 Diagram Alir <i>Preprocessing</i>	19
Gambar 4.4 Diagram Alir <i>Case Folding</i>	19
Gambar 4.5 Diagram Alir <i>Tokenizing</i>	20
Gambar 4.6 Diagram Alir <i>Filtering</i>	20
Gambar 4.7 Diagram Alir <i>Stemming</i>	21
Gambar 4.8 Diagram alir Ekspansi Kata.....	22
Gambar 4.9 Diagram Alir Klasifikasi <i>K-Nearest Neighbor</i> dengan <i>Query Expansion</i>	23
Gambar 4.10 Halaman Perancangan Dokumen Kamus berita.....	42
Gambar 4.11 Perancangan Halaman Klasifikasi.....	43
Gambar 4.12 Perancangan Halaman Pengujian.....	44
Gambar 5.1 Tampilan Dokumen berita Ekspansi Kata.....	55
Gambar 5.2 Tampilan Halaman Klasifikasi Teks.....	56
Gambar 5.3 Tampilan Halaman Hasil Klasifikasi.....	56
Gambar 6.1. Grafik Perolehan Pengujian Nilai <i>k</i> Pada <i>KNN</i>	59
Gambar 6.2 Grafik Perolehan nilai Pengujian Variasi <i>Threshold</i> Pada <i>KNN</i>	60

DAFTAR KODE PROGRAM

Kode Program 5.1 Implementasi <i>Case Folding</i>	46
Kode Program 5.2 Implementasi <i>Tokenizing</i>	47
Kode Program 5.3 Implementasi <i>Filtering</i>	47
Kode Program 5.4 Implementasi <i>Stemming</i>	48
Kode Program 5.5 Implementasi <i>Distributional Semantic</i>	48
Kode Program 5.6 Implementasi Ekspansi kata.....	49
Kode Program 5.7 Implementasi Penentuan <i>Term</i> Setiap Dokumen.....	51
Kode Program 5.8 Implementasi Pembobotan <i>TF-Idf</i>	52
Kode Program 5.9 Implementasi Pemetaan Pembobotan Awal.....	52
Kode Program 5.10 Implementasi Normalisasi Ketetapan Bobot.....	53
Kode Program 5.11 Implementasi Perhitungan <i>Cosine Similarity</i>	54

DAFTAR PERSAMAAN

Persamaan 2.1 <i>Euclidean Distance</i>	10
Persamaan 2.2 <i>Cosine Similarity</i>	12
Persamaan 2.3 Perhitungan Akurasi.....	12
Persamaan 4.1 Perhitungan <i>Cosine Similarity</i>	41
Persamaan 4.1 Contoh Perhitungan <i>Euclidean Distance</i>	42

DAFTAR LAMPIRAN

LAMPIRAN A Data Uji.....	68
LAMPIRAN B Data Latih <i>Twitter</i> Kompas dan Detik	71
LAMPIRAN C Dokumen Berita.....	8

